

# Faster Algorithms for Convex and Combinatorial Optimization

by

Yin Tat Lee

Submitted to the Department of Mathematics  
in partial fulfillment of the requirements for the degree of

Doctor of Philosophy in Mathematics

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2016

© Massachusetts Institute of Technology 2016. All rights reserved.

Author .....  
Department of Mathematics  
May 13, 2016

Certified by .....  
Jonathan A. Kelner  
Associate Professor  
Thesis Supervisor

Accepted by .....  
Jonathan A. Kelner  
Chairman, Applied Mathematics Committee



# Faster Algorithms for Convex and Combinatorial Optimization

by  
Yin Tat Lee

Submitted to the Department of Mathematics  
on May 13, 2016, in partial fulfillment of the  
requirements for the degree of  
Doctor of Philosophy in Mathematics

## Abstract

In this thesis, we revisit three algorithmic techniques: sparsification, cutting and collapsing. We use them to obtain the following results on convex and combinatorial optimization:

- **Linear Programming:** We obtain the first improvement to the running time for linear programming in 25 years. The convergence rate of this randomized algorithm nearly matches the universal barrier for interior point methods. As a corollary, we obtain the first  $\tilde{O}(m\sqrt{n})$  time randomized algorithm for solving the maximum flow problem on directed graphs with  $m$  edges and  $n$  vertices. This improves upon the previous fastest running time of  $\tilde{O}(m \min(n^{2/3}, m^{1/2}))$ , achieved over 15 years ago by Goldberg and Rao.
- **Maximum Flow Problem:** We obtain one of the first almost-linear time randomized algorithms for approximating the maximum flow in undirected graphs. As a corollary, we improve the running time of a wide range of algorithms that use the computation of maximum flows as a subroutine.
- **Non-Smooth Convex Optimization:** We obtain the first nearly-cubic time randomized algorithm for convex problems under the black box model. As a corollary, this implies a polynomially faster algorithm for three fundamental problems in computer science: submodular function minimization, matroid intersection, and semidefinite programming.
- **Graph Sparsification:** We obtain the first almost-linear time randomized algorithm for spectrally approximating any graph by one with just a linear number of edges. This sparse graph approximately preserves all cut values of the original graph and is useful for solving a wide range of combinatorial problems. This algorithm improves all previous linear-sized constructions, which required at least quadratic time.
- **Numerical Linear Algebra:** Multigrid is an efficient method for solving large-scale linear systems which arise from graphs in low dimensions. It has been used extensively for 30 years in scientific computing. Unlike the previous approaches that make assumptions on the graphs, we give the first generalization of the multigrid that provably solves Laplacian systems of any graphs in nearly-linear expected time.

Thesis Supervisor: Jonathan A. Kelner  
Title: Associate Professor



---

---

## Acknowledgments

First, I would like to thank my advisor Jonathan Kelner. Jon is always full of ideas and provided me a lot of high-level inspiration. I am particularly thankful to him for encouraging me to learn interior point methods, which led to many of key results in this thesis. From the very beginning, Jon has given me lots of support, freedom and advice. Because of him, I have achieved what I never imagined. I will be forever grateful to Jon.

Next, I would like to thank my colleague, Aaron Sidford. Over the past 4 years, one of my most exciting activities has been discussing with Aaron in how to solve Maximum Flow faster. During thousands of hours of discussions, we obtained many of the results presented in this thesis. Although we still haven't gotten the ultimate result we want, I believe, someday, we will get it.

I would also like to thank Lap Chi Lau. Originally, I was planning to study scientific computing, but he made me change my mind. Since my undergraduate, I have been asking him for all kinds of career advice and, up to now, all of them are proven to be useful.

Next, I would like to thank Pravin Vaidya for his beautiful work on his interior point methods, cutting plane methods and Laplacian systems. Vaidya's results were definitely far ahead of his time and many results in this thesis are inspired by and built on his results.

Also, I would like to thank my academic friends and collaborators: Sébastien Bubeck, Michael Cohen, Laurent Demanet, Michel Goemans, Michael Kapralov, Tsz Chiu Kwok, Rasmus Kyng, Aleksander Mądry, Gary Miller, Cameron Musco, Christopher Musco, Lorenzo Orecchia, Shayan Oveis Gharan, Jakub Pachocki, Richard Peng, Satish Rao, Sushant Sachdeva, Ludwig Schmidt, Mohit Singh, Daniel Spielman, Nikhil Srivastava, Robert Strichartz, He Sun, Luca Trevisan, Santosh Vempala, Nisheeth Vishnoi, Matt Weinberg, Sam Chiu-wai Wong, and Zeyuan Allen Zhu.

Last but not the least, I would like to thank my parents Man Cheung and Sau Yuen, my siblings Ah Chi and Ah Tat, and my wife Yan Kit.

Yin Tat Lee  
Cambridge, Massachusetts  
May 2016

My work was partially supported by NSF awards 0843915 and 1111109. Some of my work was done while I was visiting the Simons Institute for the Theory of Computing, UC Berkeley.

# Contents

<b>1</b>	<b>Introduction</b>	<b>9</b>
1.1	Techniques . . . . .	9
1.2	Results . . . . .	11
<b>2</b>	<b>Preliminaries and Thesis Structure</b>	<b>15</b>
2.1	Problems . . . . .	15
2.2	Thesis Organization . . . . .	17
2.3	Notation . . . . .	24
<b>I</b>	<b>Sparsification</b>	<b>29</b>
<b>3</b>	<b><math>\ell_2</math> Regression In Input Sparsity Time</b>	<b>31</b>
3.1	Introduction . . . . .	31
3.2	Previous Work . . . . .	34
3.3	Generalized Leverage Score . . . . .	35
3.4	Expectation Bound on Leverage Score Estimation . . . . .	35
3.5	Coherence-Reducing Weighting . . . . .	36
3.6	High Probability Bound on Leverage Score Estimation . . . . .	38
3.7	Main Result . . . . .	39
3.8	Appendix . . . . .	42
<b>4</b>	<b>Linear-Sized Sparsifier In Almost-Linear Time</b>	<b>45</b>
4.1	Introduction . . . . .	45
4.2	Algorithm . . . . .	45
4.3	Analysis . . . . .	49
4.4	Appendix . . . . .	58
<b>5</b>	<b><math>\ell_p</math> Regression and John Ellipsoid</b>	<b>63</b>
5.1	Introduction . . . . .	63
5.2	Lewis Weights . . . . .	63
5.3	Approximate Lewis Weights . . . . .	64
5.4	Approximate John Ellipsoids . . . . .	67
<b>6</b>	<b>Faster Linear Programming</b>	<b>69</b>
6.1	Introduction . . . . .	69
6.2	Previous Work . . . . .	70
6.3	Simple Case . . . . .	72
6.4	Difficulties . . . . .	76

6.5	Weighted Path Finding . . . . .	77
6.6	Progressing Along Weighted Paths . . . . .	82
6.7	Weight Function . . . . .	91
6.8	The Algorithm . . . . .	100
6.9	Generalized Minimum Cost Flow . . . . .	104
6.10	Glossary . . . . .	106
6.11	Appendix: Technical Lemmas . . . . .	107
6.12	Appendix: Projection on Mixed Norm Ball . . . . .	109
6.13	Appendix: A $\tilde{O}(n)$ -Self-Concordant Barrier Function . . . . .	110
<b>7</b>	<b>Geometric Median In Nearly-Linear Time</b>	<b>123</b>
7.1	Introduction . . . . .	123
7.2	Preliminaries . . . . .	128
7.3	Properties of the Central Path . . . . .	129
7.4	Overview of the Algorithm . . . . .	130
7.5	Analysis of the Central Path . . . . .	133
7.6	Analysis of the Algorithm . . . . .	143
7.7	Pseudo Polynomial Time Algorithm . . . . .	148
7.8	Derivation of Penalty Function . . . . .	151
7.9	Technical Lemmas . . . . .	152
7.10	Weighted Geometric Median . . . . .	156
<b>II</b>	<b>Cutting</b>	<b>158</b>
<b>8</b>	<b>Convex Minimization In Nearly-Cubic Time</b>	<b>159</b>
8.1	Introduction . . . . .	159
8.2	Preliminaries . . . . .	163
8.3	Our Cutting Plane Method . . . . .	164
8.4	Technical Tools . . . . .	186
8.5	Glossary . . . . .	192
<b>9</b>	<b>Effective Use of Cutting Plane Method</b>	<b>195</b>
9.1	Introduction . . . . .	195
9.2	Preliminaries . . . . .	201
9.3	Convex Optimization . . . . .	202
9.4	Intersection of Convex Sets . . . . .	209
<b>10</b>	<b>Submodular Minimization In Nearly-Cubic Time</b>	<b>221</b>
10.1	Introduction . . . . .	221
10.2	Preliminaries . . . . .	224
10.3	Improved Weakly Polynomial Algorithms for SFM . . . . .	227
10.4	Improved Strongly Polynomial Algorithms for SFM . . . . .	229
10.5	Discussion and Comparison with Previous Algorithms . . . . .	252
<b>11</b>	<b>First Order Method vs Cutting Plane Method</b>	<b>255</b>
11.1	Introduction . . . . .	255
11.2	Intuition . . . . .	256
11.3	An Optimal Algorithm . . . . .	257

11.4	Politician . . . . .	259
11.5	Experiments . . . . .	264
11.6	Discussion . . . . .	268
11.7	Appendix: Convergence of SD+ . . . . .	270
<b>III Collapsing</b>		<b>276</b>
<b>12 Undirected MaxFlow In Almost-Linear Time</b>		<b>277</b>
12.1	Introduction . . . . .	277
12.2	Preliminaries . . . . .	282
12.3	Solving Max-Flow Using a Circulation Projection . . . . .	283
12.4	Oblivious Routing . . . . .	289
12.5	Flow Sparsifiers . . . . .	294
12.6	Removing Vertices in Oblivious Routing Construction . . . . .	304
12.7	Nonlinear Projection and Maximum Concurrent Flow . . . . .	309
12.8	Appendix . . . . .	315
<b>13 Connection Laplacian In Nearly-Linear Time</b>		<b>319</b>
13.1	Introduction . . . . .	319
13.2	Preliminaries . . . . .	321
13.3	Overview of the Algorithms . . . . .	322
13.4	Summary . . . . .	328
13.5	Background . . . . .	329
13.6	Block Diagonally Dominant Matrices . . . . .	330
13.7	Schur Complement Chains . . . . .	332
13.8	Finding $\alpha$ -bDD Subsets . . . . .	335
13.9	Jacobi Iteration on $\alpha$ -bDD Matrices . . . . .	337
13.10	Existence of Linear-Sized Approximate Inverses . . . . .	339
13.11	Spectral Vertex Sparsification Algorithm . . . . .	341
13.12	Estimating Leverage Scores by Undersampling . . . . .	356
13.13	Recursive Construction of Schur Complement Chains . . . . .	358
13.14	Weighted Expander Constructions . . . . .	363
<b>IV Geometry</b>		<b>374</b>
<b>14 Sampling In Sub-Quadratic Steps</b>		<b>375</b>
14.1	Introduction . . . . .	375
14.2	Preliminaries . . . . .	379
14.3	Intuition . . . . .	386
14.4	Convergence of the Geodesic Walk . . . . .	389
14.5	Logarithmic Barrier . . . . .	407
14.6	Collocation Method for ODE . . . . .	431
14.7	Derivative Estimations . . . . .	437



Convex optimization has been studied extensively and is a prominent tool used in various areas such as combinatorial optimization, data analysis, operations research, and scientific computing. Each field has developed specialized tools including data structures, sampling methods, and dimension reduction. In this thesis, we combine and improve the optimization techniques from different fields to design provably faster optimization algorithms for several important problems, including linear programming, general non-smooth convex minimization, and maximum flow problems.

### ■ 1.1 Techniques

The techniques used in this thesis can be categorized into three types: sparsification, cutting and collapsing. Here, we explain these three techniques and highlight their applications to convex optimization and combinatorial optimization.

#### ■ 1.1.1 Sparsification

Sparsification techniques originated from graph theory, and it refers to the task of approximating dense graphs by sparse graphs. There are different types of sparsifiers specifically designed for different graph problems, such as spanners for distance-related problems [49], cut sparsifiers for cut-related problems [33], spectral sparsifiers for spectral related problems [243, 242]. For instance, if we want to compute a minimum  $s$ - $t$  cut on a given dense graph, instead of running our favorite  $s$ - $t$  cut algorithm on the original graph, we can first compute a cut sparsifier of the original graph, then run the  $s$ - $t$  cut algorithm on the sparsifier. These sparsifiers allow us to speed up graph algorithms for many problems on dense graphs by paying some error introduced by the sparsifiers.

Beyond graph related problems, this notion of spectral sparsifier has been generalized and applied to solve different matrix-related problems [76, 170, 56]. In the linear regression problem, we are given a matrix  $\mathbf{A}$  and a vector  $\mathbf{b}$ , and we want to minimize  $\|\mathbf{Ax} - \mathbf{b}\|_2$  over all  $\mathbf{x} \in \mathbb{R}^n$ . The sparsification concept allows us to compress a tall matrix  $\mathbf{A}$  into an almost square matrix. This is useful in various contexts such as statistics problems with much more samples than variables. However, how can we find those sparsifiers in the first place without solving a dense problem? How sparse of a matrix can we effectively attain? How general is the idea of sparsification? In Part I, we answer these questions. In particular, we show how to solve overdetermined systems faster, how to find linear-sized sparsifiers of graphs in almost-linear time, and how to find geometric medians of point clouds in nearly-linear time. In one of the key results of this part, we extend the idea of sparsification to linear programming and give the first improvement to the running time for linear programming in 25 years.

### ■ 1.1.2 Cutting

Cutting refers to the task of iteratively refining a region that contains the solution. This technique is particularly useful for minimizing convex functions. In combinatorial optimization, Khachiyan’s seminal result in 1980 [142] proved that the ellipsoid method, a cutting plane method<sup>1</sup>, solves linear programs in polynomial time. Since then, cutting plane methods have been crucial to solving discrete problems in polynomial time [110]. In continuous optimization, cutting plane methods have long played a critical role in convex optimization, where they are fundamental to the theory of non-smooth optimization [102].

Despite the key role that cutting plane methods have played historically in both combinatorial and convex optimization, over the past two decades, the theoretical running time of cutting plane methods for convex optimization has not been improved since the breakthrough result by Vaidya in 1989 [253, 253]. Moreover, for many of the key combinatorial applications of cutting plane methods, such as submodular minimization, matroid intersection, and submodular flow, the running time improvements over the past two decades have been primarily combinatorial; that is they have been achieved by discrete algorithms that do not use numerical machinery such as cutting plane methods.

In Part II, we make progress towards cutting plane methods on multiple directions. Firstly, we show how to improve on the running time of cutting plane methods. Secondly, we provide several frameworks for applying the cutting plane methods and illustrate the efficacy of these frameworks by obtaining faster running times for semidefinite programming, matroid intersection, and submodular flow. Thirdly, we show how to couple our approach with the problem-specific structure and obtain faster weakly and strongly polynomial running times for submodular function minimization, a problem of tremendous importance in combinatorial optimization. In both cases, our algorithms are faster than the previous best running times by a factor of roughly  $\Omega(n^2)$ . Finally, we give some variants of cutting plane methods that is suitable for large-scale problems and demonstrate empirically that our new methods compare favorably with the state-of-the-art algorithms.

### ■ 1.1.3 Collapsing

Collapsing refers to the task of removing variables from problems while (approximately) preserving solutions. There are different collapsings for different situations, such as mathematical induction for proofs, Gaussian elimination and Schur complement for matrices, and Fourier–Motzkin elimination for linear programs. Since collapsing often increases the density of the problems, it often does not lead to efficient algorithms.

In a seminal work of Spielman and Teng, they showed that this density issue can be avoided for graph Laplacians via sparsification. In Part III, we generalize this idea and show how to combine collapsing and sparsification for the maximum flow problem and various linear systems. As a result, we obtain the first almost-linear time algorithms for approximating the maximum flow in undirected graphs<sup>2</sup> and the first nearly-linear-time algorithms for solving systems of equations in connection Laplacians, a generalization of Laplacian matrices that arises in many problems in image and signal processing.

In general, we assume  $T(m, n)$  to be the time needed to solve a certain graph problem on a graph with  $m$  edges and  $n$  vertices. Suppose that we know  $\mathcal{T}(m, n) = O(m) + O(\mathcal{T}(O(n), n))$  by sparsification and  $\mathcal{T}(m, n) = O(m) + O(\mathcal{T}(O(m), n/2))$  by collapsing, then it is easy to show that  $\mathcal{T}(m, n) = O(m)$

<sup>1</sup>Throughout this thesis, our focus is on algorithms for polynomial time solvable convex optimization problems given access to a linear separation oracle. Our usage of the term *cutting plane methods*, should not be confused with work on integer programming.

<sup>2</sup>We also note that independently, Jonah Sherman produced an almost-linear time algorithm for maximum flow.

by alternatively sparsifying and collapsing. Although we do not know how to do sparsifying and collapsing in general, we believe that the concept of combining sparsifying and collapsing is a powerful technique that it will provide a foundation of future works.

## ■ 1.2 Results

By applying and improving the three techniques mentioned above, we have substantially advanced the state-of-the-art algorithms for solving many fundamental problems in computer science. Here, we describe three of the key results in this thesis.

### ■ 1.2.1 Convex Programming

Many optimization problems in theoretical computer science are shown to be polynomial-time-solvable by reducing to one of the following two general problems:

- **The Feasibility Problem:** Given a set  $K \subset \mathbb{R}^n$  and a separation oracle that, given a point  $\mathbf{x}$ , either outputs that  $\mathbf{x}$  is in  $K$  or outputs a separating hyperplane. The problem is to find a point in  $K$  or prove that  $K$  is almost empty.
- **The Intersection Problem:** Given two sets  $K_1, K_2 \subset \mathbb{R}^n$ , any vector  $\mathbf{d}$  and an optimization oracle that, given a vector  $\mathbf{c}$ , outputs an  $\mathbf{x}_1$  that maximizes  $\mathbf{c}^\top \mathbf{x}_1$  over  $K_1$  and  $\mathbf{x}_2$  that maximizes  $\mathbf{c}^\top \mathbf{x}_2$  over  $K_2$ . The problem is to find a point  $\mathbf{x}$  nearly inside  $K_1 \cap K_2$  such that it approximately maximizes  $\mathbf{d}^\top \mathbf{x}$ .

For many problems, this reduction is either the best or the only known approach. Therefore, it is important to find an efficient algorithm for these two general problems. Since a breakthrough of Vaidya in 1989 [253], the feasibility problem can be solved in  $\tilde{O}(nT + n^{\omega+1})$  time (Table 1.1), where  $T$  is the cost of the oracle and  $\omega \sim 2.373$ . Furthermore, using the seminal works of Grötschel, Lovász, Schrijver and independently, Karp and Papadimitriou in 1981 and 1982 [110, 138], the intersection problem can be solved by solving  $n$  feasibility problems. In this thesis, we show that both problems can be solved in  $\tilde{O}(nT + n^3)$  expected time<sup>3</sup>, the first improvement for both problems over 25 years.

As a corollary of these two theorems, we obtained polynomial improvements for many important problems, such as submodular function minimization (Table 1.2), matroid intersection (Table 1.3, 1.4), submodular flow (Table 1.5), and semidefinite programming (Table 1.6).

Year	Algorithm	Complexity
1979	Ellipsoid Method [236, 274, 142]	$\tilde{O}(n^2T + n^4)$
1988	Inscribed Ellipsoid [144, 213]	$\tilde{O}(nT + n^{4.5})$
1989	Volumetric Center [253]	$\tilde{O}(nT + n^{\omega+1})$
1995	Analytic Center [20]	$\tilde{O}(nT + n^{\omega+1})$
2004	Random Walk [36]	$\rightarrow \tilde{O}(nT + n^7)$
-	Chapter 8	$\tilde{O}(nT + n^3)$

**Table 1.1:** Algorithms for the Feasibility Problem. The arrow,  $\rightarrow$ , indicates that it solves a more general problem where only a membership oracle is given.

<sup>3</sup>In this thesis, nearly all algorithms are randomized and nearly all time are expected running time.

Year	Author	Complexity
1981,1988	Grötschel, Lovász, Schrijver [110, 109]	$\tilde{O}(n^5 \cdot \text{EO} + n^7)$ [189]
1985	Cunningham [61]	$O(Mn^6 \log(nM) \cdot \text{EO})$
2000	Schrijver [231]	$O(n^8 \cdot \text{EO} + n^9)$
2000	Iwata, Fleischer, Fujishige [121]	$O(n^5 \cdot \text{EO} \log M)$ $O(n^7 \log n \cdot \text{EO})$
2000	Iwata, Fleischer [87]	$O(n^7 \cdot \text{EO} + n^8)$
2003	Iwata [118]	$O((n^4 \cdot \text{EO} + n^5) \log M)$ $O((n^6 \cdot \text{EO} + n^7) \log n)$
2003	Vygen [264]	$O(n^7 \cdot \text{EO} + n^8)$
2007	Orlin [216]	$O(n^5 \cdot \text{EO} + n^6)$
2009	Iwata, Orlin [125]	$O((n^4 \cdot \text{EO} + n^5) \log(nM))$ $O((n^5 \cdot \text{EO} + n^6) \log n)$
-	Chapter 9	$O(n^2 \log(nM) \cdot \text{EO} + n^3 \log^{O(1)}(nM))$ $O(n^3 \log^2 n \cdot \text{EO} + n^4 \log^{O(1)} n)$

**Table 1.2:** Algorithms for submodular function minimization. EO is the time for evaluation oracle of the submodular function and  $M$  is the maximum absolute function value.

Year	Author	Complexity
1968	Edmonds [81]	not stated
1971	Aigner, Dowling [6]	$O(n^3 \mathcal{T}_{\text{ind}})$
1975	Lawler [157]	$O(n^3 \mathcal{T}_{\text{ind}})$
1986	Cunningham [60]	$O(n^{2.5} \mathcal{T}_{\text{ind}})$
-	Chapter 9	$\tilde{O}(n^2 \mathcal{T}_{\text{ind}} + n^3)$

**Table 1.3:** Algorithms for (unweighted) matroid intersection.  $n$  is the size of the ground set and  $\mathcal{T}_{\text{ind}}$  is the time to check if a set is independent.

Year	Author	Complexity
1978	Fujishige [92]	not stated
1981	Grotschel, Lovasz, Schrijver[110]	weakly polynomial
1987	Frank, Tardos [91]	strongly polynomial
1993	McCormick, Ervolina [190]	$O(n^7 h^* \log(nCU))$
1994	Wallacher, Zimmermann [266]	$O(n^8 h \log(nCU))$
1997	Iwata [117]	$O(n^7 h \log U)$
1998	Iwata, McCormick, Shigeno [123]	$O(n^4 h \min \{\log(nC), n^2 \log n\})$
1999	Iwata, McCormick, Shigeno [124]	$O(n^6 h \min \{\log(nU), n^2 \log n\})$
1999	Fleischer, Iwata, McCormick[89]	$O(n^4 h \min \{\log U, n^2 \log n\})$
1999	Iwata, McCormick, Shigeno [122]	$O(n^4 h \min \{\log C, n^2 \log n\})$
2000	Fleischer, Iwata [88]	$O(mn^5 \log(nU) \cdot \text{EO})$
-	Chapter 9	$O(n^2 \log(nCU) \cdot \text{EO} + n^3 \log^{O(1)}(nCU))$

**Table 1.5:** Algorithms for minimum cost submodular flow with  $n$  vertices, maximum cost  $C$  and maximum capacity  $U$ . The factor  $h$  is the time for an exchange capacity oracle,  $h^*$  is the time for a more complicated exchange capacity oracle, and EO is the time for evaluation oracle.

Year	Author	Complexity
1968	Edmonds [81]	not stated
1975	Lawler [157]	$O(n^3\mathcal{T}_{\text{ind}} + n^4)$
1981	Frank [90]	$O(n^3(\mathcal{T}_{\text{circuit}} + n))$
1986	Brezovec, Cornuéjols, Glover [39]	$O(n^2(\mathcal{T}_{\text{circuit}} + n))$
1995	Fujishige, Zhang [96]	$O(n^{2.5} \log(nM) \cdot \mathcal{T}_{\text{ind}})$
1995	Shigeno, Iwata [234]	$O((n + \mathcal{T}_{\text{circuit}})n^{1.5} \log(nM))$
-	Chapter 9	$O(n^2 \log(nM)\mathcal{T}_{\text{ind}} + n^3 \log^{O(1)}(nM))$

**Table 1.4:** Algorithms for weighted matroid intersection. In addition to the notation in Table 1.3,  $\mathcal{T}_{\text{circuit}}$  is the time needed to find a fundamental circuit and  $M$  is the maximum weight.

Year	Author	Complexity
1992	Nesterov, Nemirovsky [210]	$\tilde{O}(\sqrt{m}(nm^\omega + n^{\omega-1}m^2))$
2000	Anstreicher [13]	$\tilde{O}((mn)^{1/4}(nm^\omega + n^{\omega-1}m^2))$
2003	Krishnan, Mitchell [152]	$\tilde{O}(n(n^\omega + m^\omega + S))$ (dual SDP)
-	Chapter 9	$\tilde{O}(n(n^2 + m^\omega + S))$

**Table 1.6:** Algorithms for solving a  $m \times m$  SDP with  $n$  constraints and  $S$  non-zero entries

## ■ 1.2.2 Linear Programming

For an arbitrary linear program  $\min_{x \text{ s.t. } Ax \geq b} c^\top x$  with  $n$  variables and  $m$  constraints, the fastest algorithm required either solving  $\tilde{O}(\sqrt{m})$  linear systems or inverting  $\tilde{O}((mn)^{1/4})$  dense matrices [227, 253]. However, since a breakthrough of Nesterov and Nemirovski in 1994 [211], it has been known that  $\tilde{O}(\sqrt{n})$  iterations in fact suffice. Unfortunately, each iteration of their algorithm requires computing the center of gravity of a polytope, which is a problem that is harder than linear programming. In this thesis, we resolved this 25-year-old computational question and developed an efficient algorithm that requires solving only  $\tilde{O}(\sqrt{n})$  linear systems (Table 1.7).

## ■ 1.2.3 Maximum Flow Problem

The maximum flow problem has been studied extensively for more than 60 years; it has a wide range of theoretical and practical applications and is widely used as a key subroutine in other algorithms. Until 5 years ago, work on the maximum flow and the more general problems of linear programming and first order convex optimization was typically performed independently. Nevertheless, all known

Year	Author	Number of Iterations	Nature of iterations
1984	Karmarkar [136]	$\tilde{O}(m)$	Linear system solve
1986	Renegar [227]	$\tilde{O}(\sqrt{m})$	Linear system solve
1989	Vaidya [253]	$\tilde{O}((mn)^{1/4})$	Expensive linear algebra
1994	Nesterov and Nemirovskii [211]	$\tilde{O}(\sqrt{n})$	Volume computation
-	Chapter 6	$\tilde{O}(\sqrt{n})$	$\tilde{O}(1)$ Linear system solves

**Table 1.7:** Algorithms for linear programs with  $n$  variables and  $m$  constraints

Year	Author	Complexity	Directed	Weighted
1998	Goldberg & Rao [104]	$\tilde{O}(m \min(m^{1/2}, n^{2/3}))$	Yes	Yes
1998	Karger [134]	$\tilde{O}(m\sqrt{n}\varepsilon^{-1})$	No	Yes
2002	Karger & Levine [135]	$O(nF)$	No	Yes
2011	Christiano et al. [51]	$\tilde{O}(mn^{1/3}\varepsilon^{-11/3})$	No	Yes
2012	Lee, Rao & Srivastava [161]	$\tilde{O}(mn^{1/3}\varepsilon^{-2/3})$	No	No
2013	Chapter 12	$\tilde{O}(m^{1+o(1)}\varepsilon^{-2})$	No	Yes
	Sherman [233]			
2013	Mądry [183]	$\tilde{O}(m^{10/7})$	Yes	No
-	Chapter 6	$\tilde{O}(m\sqrt{n})$	Yes	Yes

**Table 1.8:** Algorithms for maximum flow problem with  $n$  vertices,  $m$  edges and maximum flow value  $F$ . Results before 1998 are omitted for brevity.

techniques, whether combinatorial or numerical, achieved essentially the same (or worse) running time,  $\tilde{O}(m^{1.5})$  or  $\tilde{O}(mn^{2/3})$  for solving the maximum flow on a weighted graph with  $n$  vertices and  $m$  edges. Consequently, the maximum flow problem poses a problem for testing the state-of-the-art methods in both convex and combinatorial optimization.

In this thesis, we obtain the first  $\tilde{O}(m\sqrt{n})$  time algorithm for solving the maximum flow problem on directed graphs with  $m$  edges and  $n$  vertices. This is an improvement over the previous fastest running time of  $\tilde{O}(m \min\{\sqrt{m}, n^{2/3}\})$  achieved over 15 years ago by Goldberg and Rao [104] and the previous fastest running times for solving dense directed unit capacity graphs of  $O(m \min\{\sqrt{m}, n^{2/3}\})$  achieved by Even and Tarjan [85] over 35 years ago and of  $\tilde{O}(m^{10/7})$  achieved recently by Mądry [183]. For undirected graphs, we give one of the first algorithm with running time  $O(m^{1+o(1)}/\varepsilon^2)$  for approximating maximum flow.

## ■ 1.2.4 Philosophical Contributions

Besides the improvements on running time in many classical problems in computer science, the secondary goal of this thesis is to further promote the interaction between continuous optimization and discrete optimization.

For many decades, continuous optimization and discrete (combinatorial) optimization have been studied separately. In continuous optimization, researchers have developed optimal iterative methods under various black-box models, and in discrete optimization, researchers have developed various randomization techniques and data structures. Despite the different approaches, convex problems in both fields in fact share some common difficulties. In this thesis, we identify three important components from both fields and demonstrate the combination of both techniques is indeed useful for solving many fundamental problems in optimization.

Despite decades of research on convex optimization and combinatorial optimization, we still do not have nearly optimal time algorithm for the majority of non-trivial convex problems. We strongly believe that the combination of both fields will help advance future research in this area.

## Preliminaries and Thesis Structure

In this chapter, we first describe the problems we encounter throughout this thesis, then we discuss the structure of this thesis. At the end, we list some linear algebraic notations and some facts we commonly used throughout this thesis.

### ■ 2.1 Problems

#### ■ 2.1.1 Black-Box Problems

We say that a set  $K \subseteq \mathbb{R}^n$  is convex if for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  and all  $t \in [0, 1]$  it holds that  $t\mathbf{x} + (1-t)\mathbf{y} \in K$ . We say that a function  $f : K \rightarrow \mathbb{R}$  is convex if  $K$  is convex and for all  $\mathbf{x}, \mathbf{y} \in K$  and all  $t \in [0, 1]$  it holds that  $f(t\mathbf{x} + (1-t)\mathbf{y}) \leq t \cdot f(\mathbf{x}) + (1-t) \cdot f(\mathbf{y})$ . In this thesis, we are interested in the following minimization problem

$$\min_{\mathbf{x} \in K} f(\mathbf{x})$$

where both  $f$  and  $K$  are convex and we call this a convex problem. There are two types of convex problems, black-box and structural problems, depending on if we have full information about the problem.

In black-box problems, the function  $f$  and/or the set  $K$  are unknown and one can only access them through queries to an oracle. Here, we give a list of black-box problems we are interested in. All of these problems below are very general and cover a wide range of problems in computer science.

- **Convex Minimization:** Given a convex function  $f$  and a (sub)gradient oracle that computes  $\nabla f(\mathbf{x})$  at any point  $\mathbf{x}$ . The problem is to find a point that minimizes  $f$  over  $\mathbf{x} \in \mathbb{R}^n$ .
- **The Feasibility Problem:** Given a convex set  $K \subset \mathbb{R}^n$  and a separation oracle that, given a point  $\mathbf{x}$ , either outputs that  $\mathbf{x}$  is in  $K$  or outputs a hyperplane separates  $\mathbf{x}$  and  $K$ . The problem is to find a point in  $K$  or prove that  $K$  is almost empty.
- **The Intersection Problem:** Given two convex sets  $K_1, K_2 \subset \mathbb{R}^n$  and an optimization oracle that, given a vector  $\mathbf{c}$ , outputs a vector  $\mathbf{x}_1$  that maximizes  $\mathbf{c}^\top \mathbf{x}_1$  over  $K_1$  and  $\mathbf{x}_2$  that maximizes  $\mathbf{c}^\top \mathbf{x}_2$  over  $K_2$ . For any vector  $\mathbf{d}$ , the problem is to find a point  $\mathbf{x}$  nearly inside  $K_1 \cap K_2$  such that it approximately maximizes  $\mathbf{d}^\top \mathbf{x}$ .
- **Submodular Minimization:** Given an integer-valued function  $f$  defined on all subsets of  $n$  elements.  $f$  is submodular if it satisfies the diminishing returns, i.e.

$$f(X \cup \{\mathbf{x}\}) - f(X) \geq f(Y \cup \{\mathbf{x}\}) - f(Y)$$

for any  $X, Y$  such that  $X \subset Y$  and any  $\mathbf{x} \notin Y$ . We are given an evaluation oracle that computes  $f(S)$  for any subset  $S$ . The problem is to find a subset  $S$  that minimizes  $f$ .

### ■ 2.1.2 Structural Problems

In structural problems, the function  $f$  and the set  $K$  are given directly by explicit formulas. In this and the next section, we list some structure problems we are interested in this thesis.

- **$\ell_p$  regression:** Given a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and a vector  $\mathbf{b} \in \mathbb{R}^m$ , the  $\ell_p$  regression problem requires us to find the minimizer of

$$\min_x \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_p.$$

We are particularly interested in the case  $p = 2$  because many optimization algorithms need to solve a linear system per iteration. As we will see, techniques for solving  $\ell_2$  regression problems can be sometimes used or generalized to general convex problems.

- **Linear Programming:** Given a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and vectors  $\mathbf{b} \in \mathbb{R}^m, \mathbf{c} \in \mathbb{R}^n$ , linear programming requires us to find the minimizer of

$$\min_{\mathbf{A}\mathbf{x} \geq \mathbf{b}} \mathbf{c}^\top \mathbf{x}.$$

This is one of the key problems in computer science. In both theory and practice, many problems can be reformulated as linear programs to take advantage of fast algorithms.

- **Semidefinite Programming:** Given a symmetric matrices  $\mathbf{C} \in \mathbb{R}^{m \times m}, \mathbf{A}_i \in \mathbb{R}^{m \times m}$  for  $i = 1, \dots, n$ , the semidefinite program requires us to find the minimizer of

$$\max_{\mathbf{X} \succeq \mathbf{0}} \mathbf{C} \bullet \mathbf{X} \text{ s.t. } \mathbf{A}_i \bullet \mathbf{X} = b_i.$$

This is a generalization of linear programming and has a lot of applications in control theory, polynomial minimization, and approximation algorithms.

- **Geometric Median (Fermat-Weber problem):** Given  $n$  points  $\{\mathbf{a}_i\}_{i=1}^m$  in  $\mathbb{R}^n$ , the geometric median problem requires us to find the minimizer of

$$\min_x \sum_{i=1}^m \|\mathbf{a}_i - \mathbf{x}\|_2.$$

This is one of the oldest problems in computation geometry and has numerous applications such as facility location, clustering and statistics.

### ■ 2.1.3 Graph-Related Problems

Here, we define some graph-related problems that are used throughout in this thesis.

We let  $G = (V, E)$  denote a graph with  $n = |V|$  vertices,  $m = |E|$  edges. We let  $w_e \geq 0$  denote the weight of an edge. For all undirected graphs in this thesis, we assume an arbitrary orientation of the edges to clarify the meaning of vectors  $\mathbf{f} \in \mathbb{R}^E$ .

Let  $\mathbf{W} \in \mathbb{R}^{E \times E}$  denote the diagonal matrices associated with the weights and  $\mathbf{B} \in \mathbb{R}^{E \times V}$  denote the graphs incidence matrix where for all  $e = (a, b) \in E$  we have  $\mathbf{B}^\top \mathbf{1}_e = \mathbf{1}_a - \mathbf{1}_b$ . Let  $\mathcal{L} \in \mathbb{R}^{V \times V}$  denote the graph Laplacian, i.e.  $\mathcal{L} \stackrel{\text{def}}{=} \mathbf{B}^\top \mathbf{W} \mathbf{B}$ . It is easy to see that

$$\mathbf{x}^\top \mathcal{L} \mathbf{x} = \sum_{u \sim v} w_{u,v} (\mathbf{x}_u - \mathbf{x}_v)^2 \tag{2.1}$$

for any  $\mathbf{x} \in \mathbb{R}^n$ . Here are three important graph related problems:



**Shortest Path Problem:** The shortest path problem requires us to find a path  $\mathbf{f}$  between two vertices that minimizes the sum of the lengths  $l$  on its constituent edges. Let  $s$  be the starting point and  $t$  be the ending point. This problem can be written as

$$\min_{\mathbf{B}^\top \mathbf{f} = \boldsymbol{\chi}} \sum_e w_e^{-1} |\mathbf{f}_e|$$

where  $\boldsymbol{\chi} = \mathbf{1}_s - \mathbf{1}_t$  and  $w_e = 1/l_e$ . We say  $\mathbf{f} \in \mathbb{R}^E$  meets the demands  $\boldsymbol{\chi}$  if  $\mathbf{B}^\top \mathbf{f} = \boldsymbol{\chi}$ .

**Laplacian Systems:** Given a graph with weights  $\mathbf{w}$  and a vector  $\boldsymbol{\chi} \in \mathbb{R}^E$  with  $\sum_e \chi_e = 0$ , our goal is to find  $\mathbf{x}$  such that  $\mathcal{L}\mathbf{x} = \boldsymbol{\chi}$ . Setting  $\mathbf{f} = \mathbf{W}\mathbf{B}\mathbf{x}$ , this problem can be written as

$$\min_{\mathbf{B}^\top \mathbf{f} = \boldsymbol{\chi}} \sum_e w_e^{-1} \mathbf{f}_e^2.$$

**Undirected Maximum Flow Problem:** Given a graph with capacities  $\mathbf{w}$ , the maximum flow problem requires us to find a flow  $\mathbf{f}$  that routes as much flow as possible from a source vertex  $s$  to a sink vertex  $t$  while sending at most  $w_e$  units of flow over each edge  $e$ . Rescaling  $\mathbf{f}$ , this problem can be written as

$$\min_{\mathbf{B}^\top \mathbf{f} = \boldsymbol{\chi}} \max_e w_e^{-1} |\mathbf{f}_e|$$

where  $\boldsymbol{\chi} = \mathbf{1}_s - \mathbf{1}_t$ .

Interestingly, these three problems respectively require us to find a flow that satisfies a certain demand and minimizes  $l_1$ ,  $l_2$ , or  $l_\infty$  norm.

## ■ 2.2 Thesis Organization

This thesis presents a thorough study of convex optimization and their applications to combinatorial optimization. In figure 2-1, we show the relationship between the chapters. Although the ideas used and the questions addressed in each chapter are all directly or indirectly related, we present the results in each chapter in a modular way so that they may be read in any order. We now discuss the results of each chapter in further details.

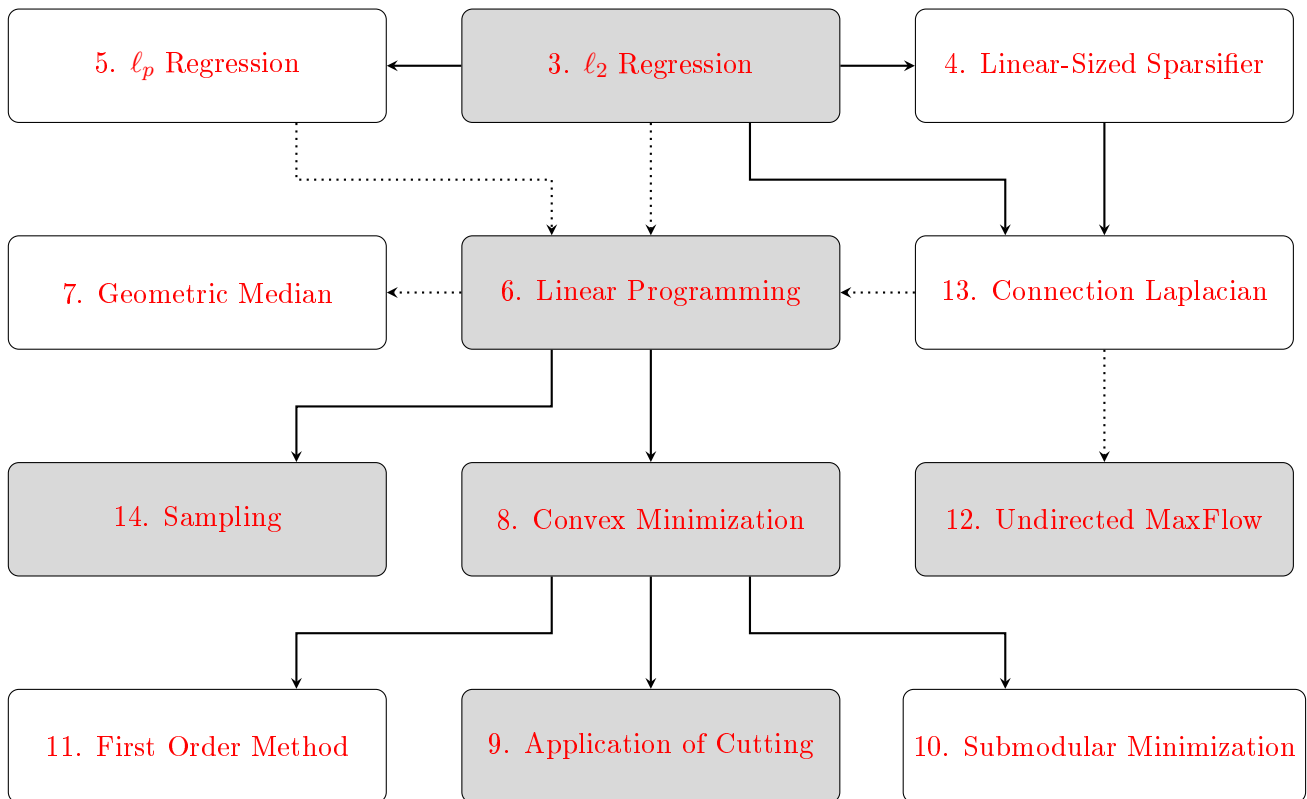
### ■ 2.2.1 Part I: Sparsification

In this part, we study several structural convex problems where the convex functions involve some tall matrices  $\mathbf{A}$ . Our goal is to find general techniques to reduce problems involve tall-and-skinny matrices  $\mathbf{A}$  to problems involve only square matrices.

#### 2.2.1.1 Chapter 3: $l_2$ Regression In Input Sparsity Time

In this chapter, we study the  $l_2$  regression problem  $\min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$  for tall matrices  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . In general, a small, manageable set of rows of  $\mathbf{A}$  can be randomly selected to approximate a tall-and-skinny data matrix, improving processing time significantly. In particular, if we sample each row with probability proportional to its importance (called *leverage score*),  $O(n \log n)$  rows suffices to approximate  $\mathbf{A}^\top \mathbf{A}$  (Lemma 2.3.3). Unfortunately, the previous best approach to compute leverage scores involves solving some linear systems of the form  $\mathbf{A}^\top \mathbf{A}$  (Lemma 2.3.4) and hence this is as difficult as the original problem. A simple alternative is to sample rows uniformly at random. While this often works, uniform sampling will eliminate critical row information for many natural instances.

We consider uniform sampling from a new perspective by examining what information it *does* preserve. Specifically, we show that uniform sampling yields a matrix that, in some sense, well ap-



**Figure 2-1:** This flow chart shows how the chapters are related between each other. A solid arrow indicates a direct usage of the result and a dotted arrow indicates a vague relation. The gray box indicates the key results in this thesis.

proximates a large fraction of the original matrix. While this weak form of approximation is not enough to solve linear regression directly, it *is* enough to compute a better approximation. This observation leads to simple iterative row sampling algorithms for matrix approximation that preserves row structure and sparsity at all intermediate steps. If it takes  $\mathcal{T}(m, n)$  time to solve a  $\ell_2$  regression problem with  $m$  rows and  $n$  columns, our result can be viewed as a reduction showing that

$$\mathcal{T}(m, n) = \tilde{O}(\text{nnz}(\mathbf{A})) + \tilde{O}(\mathcal{T}(O(n \log n), n)).$$

### 2.2.1.2 Chapter 4: Linear-Sized Sparsifier In Almost-Linear Time

Naturally, one may ask if  $\Omega(n \log n)$  rows is necessary for spectral approximations. In a seminal work of Batson, Spielman, Srivastava, they showed that

**Theorem 2.2.1** ([31]). *For any matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , there is a diagonal matrix  $\mathbf{D}$  with  $O(n/\varepsilon^2)$  non-zeros such that*

$$(1 - \varepsilon)\mathbf{A}^\top \mathbf{A} \preceq \mathbf{A}^\top \mathbf{D} \mathbf{A} \preceq (1 + \varepsilon)\mathbf{A}^\top \mathbf{A}.$$

However, all previous constructions of linear-sized spectral sparsification require solving  $\tilde{\Omega}(n)$  linear systems. In comparison, the sampling approaches use  $O(n \log n/\varepsilon^2)$  rows, but the construction only requires solving  $\tilde{O}(1)$  linear systems. In this chapter, we give a much better trade-off for getting linear-sized spectral sparsifiers. For any integer  $\rho$ , we present an algorithm that uses  $O(\rho n/\varepsilon^2)$  rows via solving  $\tilde{O}(n^{1/\rho})$  linear systems. Hence, this gives a smooth trade off between sparsity and running time. A key component in our algorithm is a novel combination of two techniques used in literature for constructing spectral sparsification: random sampling by leverage scores and adaptive constructions based on barrier functions. Similar to the previous chapter, one can view this result as a reduction showing that

$$\mathcal{T}(m, n) = O(\text{nnz}(\mathbf{A})n^{o(1)}) + \mathcal{T}(O(n), n).$$

We believe the term  $n^{o(1)}$  can be improved to  $\log^{O(1)} n$  and we left it as an open problem.

### 2.2.1.3 Chapter 4: $\ell_p$ Regression and John Ellipsoid

Similar to the  $\ell_2$  regression, it is known that if we sample each row with certain probability (called  $\ell_p$  Lewis weight), one can construct a matrix  $\tilde{\mathbf{A}}$  such that, for any  $\mathbf{x}$ ,  $\|\tilde{\mathbf{A}}\mathbf{x}\|_p \approx \|\mathbf{A}\mathbf{x}\|_p$ . This is a useful primitive for solving  $\ell_p$  regression problem. For  $1 \leq p < 4$ , Cohen and Peng showed how to compute  $\ell_p$  Lewis weight by solving  $\tilde{O}(1)$  linear systems [56] and they used it to construct an efficient approximation algorithm for  $\ell_p$  regression. However, for  $p \geq 4$ , the current best algorithm for computing  $\ell_p$  Lewis weight involves applying a general algorithm for convex problems.

In this chapter, we give a simple algorithm for  $\ell_p$  Lewis weight for  $p \geq 2$ . Similar to [56], the running time of our algorithm is dominated by solving  $\tilde{O}(1)$  linear systems. Then, we discuss the relation between  $\ell_\infty$  Lewis weight and John ellipsoid, the largest volume ellipsoid inside a given convex set. Finally, we show how to approximate John ellipsoid of symmetric polytopes by solving  $\tilde{O}(1)$  linear systems again.

### 2.2.1.4 Chapter 6: Faster Linear Programming

In this chapter, we show how to use leverage score and John ellipsoid to solve linear programs using only  $\tilde{O}(\sqrt{n}L)$  iterations where  $\mathbf{A}$  is the constraint matrix of a linear program with  $m$  constraints,  $n$  variables, and bit complexity  $L$ . Each iteration of our method consists of solving  $\tilde{O}(1)$  linear systems and additional nearly-linear time computation. Our method improves upon the previous best iteration

bounds by a factor of  $\tilde{\Omega}((m/n)^{1/4})$  for methods with polynomial time computable iterations [253] and by a factor of  $\tilde{\Omega}((m/n)^{1/2})$  for methods which solve at most  $\tilde{O}(1)$  linear systems in each iteration achieved over 20 years ago [227]. This result shows that one can “sparsify” the linear program via John ellipsoid and make sure that the complexity of solving a linear program mainly depends on the number of variables  $n$  instead of the number of constraints  $m$ .

Applying our techniques to the maximum flow yields an  $\tilde{O}(m\sqrt{n}\log^2 U)$  time algorithm for solving the maximum flow problem on directed graphs with  $m$  edges,  $n$  vertices, and capacity ratio  $U$ . This improves upon the previous fastest running time of  $O(m\min\{m^{1/2}, n^{2/3}\}\log(n^2/m)\log(U))$  achieved over 15 years ago by Goldberg and Rao [104] and improves upon the previous best running times for solving dense directed unit capacity graphs of  $O(m\min\{m^{1/2}, n^{2/3}\})$  achieved by Even and Tarjan [85] over 35 years ago and a running time of  $\tilde{O}(m^{10/7})$  achieved recently by Mądry [183].

### 2.2.1.5 Chapter 7: Geometric Median In Nearly-Linear Time

From the results above, we see that a wide range of problems can be sparsified and the running time of solving a convex optimization problem mainly depends on the number of variables. In this chapter, we would like to argue it sometimes depends on the number of true “dimensions”, which can be much smaller. For example, linear programs can be solved efficiently when the constraint matrix  $\mathbf{A}$  is an identity matrix. Then, it is natural to ask if we can solve a linear program faster if the constraint matrix of the linear program can be divided into different blocks and all blocks are identity matrices. Unfortunately, we do not have such result yet.

In this chapter, we illustrate this belief via the geometric median problem: given  $n$  points in  $\mathbb{R}^d$  compute a point that minimizes the sum of Euclidean distances to the points. This is one of the oldest non-trivial problems in computational geometry yet despite an abundance of research the previous fastest algorithms for computing a  $(1 + \varepsilon)$ -approximate geometric median were  $O(d \cdot n^{4/3}\varepsilon^{-8/3})$  by Chin et. al [50],  $\tilde{O}(d \exp(\varepsilon^{-4} \log \varepsilon^{-1}))$  by Badoiu et. al [25],  $O(nd + \text{poly}(d, \varepsilon^{-1}))$  by Feldman and Langberg [86], and  $O((nd)^{O(1)} \log \frac{1}{\varepsilon})$  by Parrilo and Sturmfels [222] and Xue and Ye [271].

Different from all previous approaches, we rely on the fact that the Hessian of the objective function can be approximated by an identity matrix minus a rank 1 matrix, namely, the number of true “dimensions” of this problem is 2. This allows us to develop a non-standard interior point method that computes a  $(1 + \varepsilon)$ -approximate geometric median in time  $O(nd \log^3 \frac{n}{\varepsilon})$ .<sup>1</sup> Our result is one of very few cases we are aware of outperforming traditional interior point theory and the only we are aware of using interior point methods to obtain a nearly linear time algorithm for a canonical optimization problem that traditionally requires superlinear time.

## ■ 2.2.2 Part II: Cutting

In this part, we give a faster cutting plane method and show how to apply this cutting plane method efficiently in various problems.

### 2.2.2.1 Chapter 8: Convex Minimization In Nearly-Cubic Time

The central problem we consider in this chapter is the feasibility problem. We are promised that a set  $K$  is contained a box of radius  $R$  and a separation oracle, that given a point  $\mathbf{x}$  in time  $\text{SO}$ , either outputs that  $\mathbf{x}$  is in  $K$  or outputs a separating hyperplane. We wish to either find a point in  $K$  or prove that  $K$  does not contain a ball of radius  $\varepsilon$ .

<sup>1</sup>We also show how to compute a  $(1 + \varepsilon)$ -approximate geometric median in sub-linear time  $O(d\varepsilon^{-2})$ .

We show how to solve this problem in  $O(nSO \log(nR/\varepsilon) + n^3 \log^{O(1)}(nR/\varepsilon))$  time. This is an improvement over the previous best running time of  $O(nSO \log(nR/\varepsilon) + n^{\omega+1} \log(nR/\varepsilon))$  for the current best known bound of  $\omega \sim 2.37$  [99]. Our key idea for achieving this running time improvement is a new straightforward technique for providing low variance unbiased estimates for changes in leverage scores that we hope will be of independent interest (See Section 8.4.1). We show how to use this technique along with ideas from [20, 257] and Chapter 6 to decrease the  $\tilde{O}(n^{\omega+1} \log(nR/\varepsilon))$  overhead in the previous fastest algorithm [253].

### 2.2.2.2 Chapter 9: Effective Use of Cutting Plane Method

In this chapter, we provide two techniques for applying our cutting plane method (and cutting plane methods in general) to optimization problems and provide several applications of these techniques.

The first technique concerns reducing the number of dimensions through duality. For many problems, their dual is significantly simpler than itself (primal). We use semidefinite programming as a concrete example to show how to improve upon the running time for finding both primal and dual solution by using the cutting planes maintained by our cutting plane method.

The second technique concerns how to minimize a linear function over the intersection of convex sets using optimization oracle. We analyze a simple potential function which allows us to bypass the typical reduction between separation and optimization to achieve faster running times. This reduction provides an improvement over the reductions used previously in [110]. Moreover, we show how this technique allows us to achieve improved running times for matroid intersection and minimum cost submodular flow.

### 2.2.2.3 Chapter 10: Submodular Minimization In Nearly-Cubic Time

In this chapter, we consider the problem of submodular minimization – a fundamental problem in combinatorial optimization with many diverse applications in theoretical computer science, operations research, machine learning and economics. We show that by considering the interplay between the guarantees of our cutting plane algorithm and the primal-dual structure of submodular minimization, we can achieve improved running times in various settings.

First, we show that a direct application of our method yields an improved weakly polynomial time algorithm for submodular minimization. Then, we present a simple geometric argument that submodular function can be solved with  $O(n^3 \log n)$  oracle calls but with exponential running time. Finally, we show that by further studying the combinatorial structure of submodular minimization and a modification to our cutting plane algorithm we can obtain a fully improved strongly polynomial time algorithm for submodular minimization.

### 2.2.2.4 Chapter 11: First Order Method vs Cutting Plane Method

Cutting plane methods are often regarded as inefficient in practice because each iteration of these methods is costly and their convergence rate are sometimes independent of how simple the problem is. In this chapter, we obtain a cutting plane method enjoys the following properties

1. Each iteration takes only linear time for simple problems.
2. It converges faster for simpler problems.
3. It is a polynomial time algorithm in the worst case.

We first propose a new method, a combination of gradient descent and the ellipsoid method, that satisfies both properties 1 and 2. This method attains the optimal rate of convergence of Nesterov’s accelerated gradient descent [203] and each iteration of this method takes exactly linear time. However, this method is not a polynomial time algorithm if the function is not smooth.

Then, we propose a new framework which leverages several concepts from convex optimization, from standard first-order methods to cutting plane methods. We show how to use our framework to derive a method that satisfies all three properties. We also demonstrate empirically that our new technique compares favorably with state of the art algorithms (such as accelerated gradient descent and BFGS). Therefore, this suggests that cutting plane methods can be as aggressive as first order methods if designed properly.

### ■ 2.2.3 Part III: Collapsing

In this part, we show how to collapse vertices in the maximum flow problem and the connection Laplacian problem.

#### 2.2.3.1 Chapter 12: Undirected MaxFlow In Almost-Linear Time

Flow problems on trees are simple because we often know how to collapse the leaves on trees while maintaining the solution. Using this fact and various routing and optimization techniques, we introduce a new framework for approximately solving flow problems in capacitated, undirected graphs and apply it to provide asymptotically faster algorithms for the maximum  $s$ - $t$  flow and maximum concurrent multicommodity flow problems. For graphs with  $n$  vertices and  $m$  edges, it allows us to find an  $\varepsilon$ -approximate maximum  $s$ - $t$  flow in time  $O(m^{1+o(1)}\varepsilon^{-2})^2$ , improving the previous best bound of  $\tilde{O}(mn^{1/3}\text{poly}(1/\varepsilon))$ . Applying the same framework in the multicommodity setting solves a maximum concurrent multicommodity flow problem with  $k$  commodities in  $O(m^{1+o(1)}\varepsilon^{-2}k^2)$  time, improving the existing bound of  $\tilde{O}(m^{4/3}\text{poly}(k, \varepsilon^{-1}))$ .

Our algorithms utilize several new technical tools that we believe may be of independent interest:

- We show how to reduce approximate maximum flow and maximum concurrent flow to oblivious routing.
- We define and provide an efficient construction of a new type of *flow sparsifier* that allowing one to transfer flows from the sparse graph back to the original graph.
- We give the first almost-linear-time construction of an  $O(m^{o(1)})$ -competitive oblivious routing scheme. No previous such algorithms ran in time better than  $\tilde{\Omega}(mn)$ . By reducing the running time to almost-linear, our work provides a powerful new primitive for constructing very fast graph algorithms.

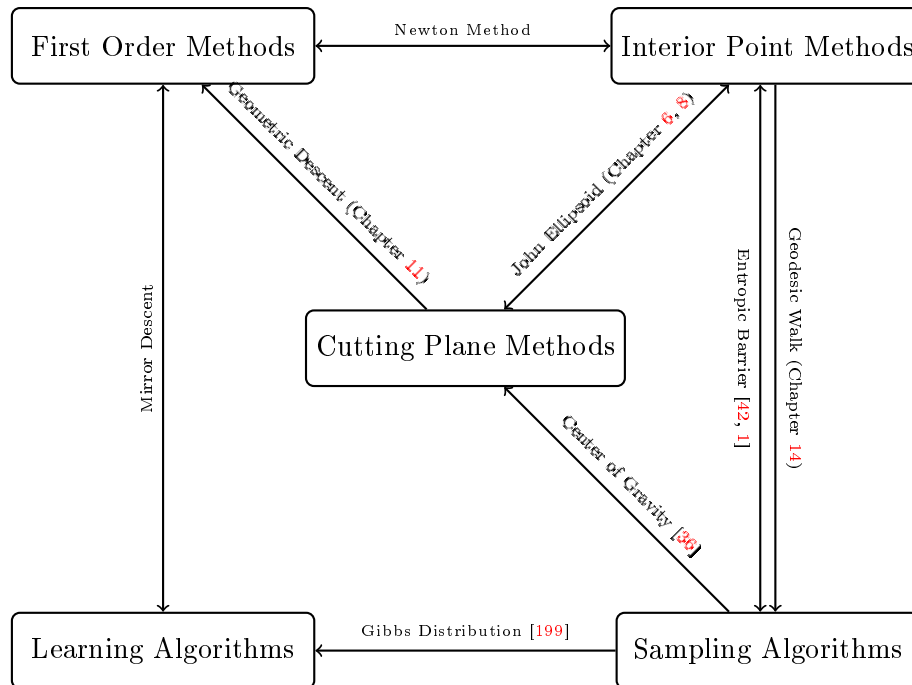
#### 2.2.3.2 Chapter 13: Connection Laplacian In Nearly-Linear Time

The idea of collapsing can be extended to general convex problems as follows: Given a convex function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . For any  $\mathbf{x} \in \mathbb{R}^n$ , we write  $\mathbf{x} = (\mathbf{y}, \mathbf{z})$  where  $\mathbf{y}$  is the first  $n/2$  variables and  $\mathbf{z}$  is the last  $n/2$  variables. We define the collapse of  $f$  by

$$g(\mathbf{y}) = \min_{\mathbf{z} \in \mathbb{R}^{n/2}} f(\mathbf{x}, \mathbf{z}).$$

---

<sup>2</sup>We also note that independently, Jonah Sherman produced an almost-linear time algorithm for maximum flow.



**Figure 2-2:** Connections between interior point methods, cutting plane methods, first-order methods, learning algorithms for bandit problems, and sampling algorithms for convex sets.

In general, a convex function is difficult to minimize if some variables are intimately related to each other. If not, we can simply solve the problem coordinate-wise. However, if some variables are related and if we can solve the half of these coordinates, the minimization problem often becomes trivial. Therefore, if we pick half of the variables in random as  $\mathbf{y}$ ,  $g$  should be easy to compute and this effectively reduce the dimension by half.

In this chapter, we show this is indeed true for connection Laplacian. As a result, we give the first nearly-linear time algorithms for solving systems of equations in connection Laplacians — a generalization of Laplacian matrices that arise in many problems in image and signal processing.

We also prove that every connection Laplacian has a linear sized approximate inverse. This is a LU factorization with a linear number of non-zero entries that is a spectral approximation of the original matrix. Using such a factorization, one can solve systems of equations in a connection Laplacian in linear time. Such a factorization was unknown even for ordinary graph Laplacians.

## ■ 2.2.4 Part IV: Geometry

In the end of this thesis, we would like to point out a large set of connections between interior point methods, cutting plane methods, first-order methods, learning algorithms for bandit problems, and sampling algorithms for convex sets as they relate to important problems across computer science, operation research, convex optimization and machine learning; see Figure 2-2 for a diagram of work that suggests this relation. One key technique spans all these algorithms is the geometry of convex set. Since these algorithms are widely used across various disciplines, we believe these connections can lead to breakthroughs in different areas. As an example, we only give an application to the sampling problem.

### 2.2.4.1 Chapter 13: Sampling In Sub-Quadratic Steps

The problems of sampling random points in convex set and computing volume of convex sets is a central problem in Markov chain Monte Carlo and leads to developments in geometric random walks. This problem has been studied for twenty years and the best algorithms take  $\tilde{O}(nm)$  iterations for explicit polytopes. We introduce the geodesic walk for sampling Riemannian manifolds and apply it to the problem of generating uniform random points from polytopes. The walk is a discrete-time simulation of a stochastic differential equation (SDE) on the Riemannian manifold equipped with the metric induced by the Hessian of a convex function used for solving linear programming; each step is the solution of an ordinary differential equation (ODE). The resulting sampling algorithm for polytopes mixes in  $\tilde{O}(mn^{\frac{3}{4}})$  steps. This is the first walk that breaks the quadratic barrier for mixing in high dimension. We also show that each step of the geodesic walk (solving an ODE) can be implemented efficiently, thus obtaining the current fastest complexity for sampling polytopes.

Advances in algorithms often come from insights from different areas; I believe these new connections will make multifaceted impacts on computation.

## ■ 2.3 Notation

Here, we define some notations that are used throughout in this thesis.

**Variables:** We use bold lowercase, e.g.  $\mathbf{x}$ , to denote a vector and bold uppercase, e.g.  $\mathbf{A}$ , to denote a matrix. For integers  $z \in \mathbb{Z}$  we use  $[z] \subseteq \mathbb{Z}$  to denote the set of integers from 1 to  $z$ . We let  $\mathbf{1}_i$  denote the vector that has value 1 in coordinate  $i$  and is 0 elsewhere. We let  $\mathbf{1}$  denote the vector that has value 1 in all coordinates. We use  $\text{nnz}(\mathbf{x})$  or  $\text{nnz}(\mathbf{A})$  to denote the number of nonzero entries in a vector or a matrix respectively. For a vector,  $\mathbf{x}$ , we let  $\|\mathbf{x}\|_M \stackrel{\text{def}}{=} \sqrt{\mathbf{x}^\top \mathbf{M} \mathbf{x}}$ . We use  $\mathbb{R}_{>0}^m$  to denote the vectors in  $\mathbb{R}^m$  where each coordinate is positive.

**Vector Operations:** We frequently apply scalar operations to vectors with the interpretation that these operations should be applied coordinate-wise. For example, for vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  we let  $\mathbf{x}/\mathbf{y} \in \mathbb{R}^n$  with  $[\mathbf{x}/\mathbf{y}]_i \stackrel{\text{def}}{=} (x_i/y_i)$  and  $\log(\mathbf{x}) \in \mathbb{R}^n$  with  $[\log(\mathbf{x})]_i = \log(x_i)$  for all  $i \in [n]$ .

**Spectral Approximations:** We call a symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  positive semidefinite (PSD) if  $\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$  for all  $\mathbf{x} \in \mathbb{R}^n$  and we call  $\mathbf{A}$  positive definite (PD) if  $\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0$  for all  $\mathbf{x} \in \mathbb{R}^n$ . For symmetric matrices  $\mathbf{N}, \mathbf{M} \in \mathbb{R}^{n \times n}$ , we write  $\mathbf{N} \preceq \mathbf{M}$  to denote that  $\mathbf{x}^\top \mathbf{N} \mathbf{x} \leq \mathbf{x}^\top \mathbf{M} \mathbf{x}$  for all  $\mathbf{x} \in \mathbb{R}^n$ , and we define  $\mathbf{N} \succeq \mathbf{M}$ ,  $\mathbf{N} \prec \mathbf{M}$  and  $\mathbf{N} \succ \mathbf{M}$  analogously. For symmetric matrices, we call  $\mathbf{B}$  is a  $\varepsilon$ -spectral approximation of  $\mathbf{A}$  if  $(1-\varepsilon)\mathbf{A} \preceq \mathbf{B} \preceq (1+\varepsilon)\mathbf{A}$ . We denote this by  $\mathbf{A} \approx_\varepsilon \mathbf{B}$ . In chapter 13, we instead define it as  $e^{-\varepsilon}\mathbf{A} \preceq \mathbf{B} \preceq e^\varepsilon\mathbf{A}$  for computational simplicity. For nonsymmetric matrices, we call  $\mathbf{B}$  is a  $\varepsilon$ -spectral approximation of  $\mathbf{A}$  if  $\mathbf{B}^\top \mathbf{B}$  is a  $\varepsilon$ -spectral approximation of  $\mathbf{A}^\top \mathbf{A}$ .

**Matrix Operations:** For  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times m}$ , we let  $\mathbf{A} \circ \mathbf{B}$  denote the Schur product, i.e.  $[\mathbf{A} \circ \mathbf{B}]_{ij} \stackrel{\text{def}}{=} A_{ij} \cdot B_{ij}$  for all  $i \in [n]$  and  $j \in [m]$ , and we let  $\mathbf{A}^{(2)} \stackrel{\text{def}}{=} \mathbf{A} \circ \mathbf{A}$ . For any norm  $\|\cdot\|$  and matrix  $\mathbf{M}$ , the operator norm of  $\mathbf{M}$  is defined by  $\|\mathbf{M}\| = \sup_{\|\mathbf{x}\|=1} \|\mathbf{M}\mathbf{x}\|$ . For any two matrices  $\mathbf{A}$  and  $\mathbf{B}$  of equal dimensions, let the dot product of  $\mathbf{A}$  and  $\mathbf{B}$  defined by  $\mathbf{A} \bullet \mathbf{B} = \text{Tr}(\mathbf{A}^\top \mathbf{B})$ . For any matrix  $\mathbf{M}$ , the Frobenius norm of  $\mathbf{M}$  is defined by  $\|\mathbf{M}\|_F = \sqrt{\sum_{i,j} M_{ij}^2} = \sqrt{\mathbf{M} \bullet \mathbf{M}}$ .

**Spectral Theory:** For  $\mathbf{A} \in \mathbb{R}^{m \times n}$  with rank  $r$ , the singular value decomposition  $\mathbf{A}$  is given by  $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$  where  $\mathbf{U} \in \mathbb{R}^{m \times r}$  and  $\mathbf{V} \in \mathbb{R}^{n \times r}$  have orthonormal columns and  $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$  is diagonal. We let  $(\mathbf{A}^\top \mathbf{A})^+$  denote the Moore-Penrose pseudoinverse of  $\mathbf{A}^\top \mathbf{A}$ , namely,  $(\mathbf{A}^\top \mathbf{A})^+ = \mathbf{V}(\mathbf{\Sigma}^{-1})^2 \mathbf{V}^\top$ . For symmetric  $\mathbf{A}$ , the singular value decomposition  $\mathbf{A}$  becomes  $\mathbf{U}\mathbf{\Sigma}\mathbf{U}^\top$  where the diagonal entries of  $\mathbf{\Sigma}$



are the eigenvalues of  $\mathbf{A}$ . We let  $\lambda_{\max}(\mathbf{A})$  and  $\lambda_{\min}(\mathbf{A})$  be the maximum and minimum eigenvalues of  $\mathbf{A}$ . The condition number of matrix  $\mathbf{A}$  is defined by  $\lambda_{\max}(\mathbf{A})/\lambda_{\min}(\mathbf{A})$ .

**Diagonal Matrices:** For  $\mathbf{A} \in \mathbb{R}^{n \times n}$  we let  $\mathbf{diag}(\mathbf{A}) \in \mathbb{R}^n$  denote the vector such that  $\mathbf{diag}(\mathbf{A})_i = \mathbf{A}_{ii}$  for all  $i \in [n]$ . For  $\mathbf{x} \in \mathbb{R}^n$  we let  $\mathbf{Diag}(\mathbf{x}) \in \mathbb{R}^{n \times n}$  be the diagonal matrix such that  $\mathbf{diag}(\mathbf{Diag}(\mathbf{x})) = \mathbf{x}$ . For  $\mathbf{A} \in \mathbb{R}^{n \times n}$  we let  $\mathbf{Diag}(\mathbf{A})$  be the diagonal matrix such that  $\mathbf{diag}(\mathbf{Diag}(\mathbf{A})) = \mathbf{diag}(\mathbf{A})$ . For  $\mathbf{x} \in \mathbb{R}^n$  when the meaning is clear from context we let  $\mathbf{X} \in \mathbb{R}^{n \times n}$  denote  $\mathbf{X} \stackrel{\text{def}}{=} \mathbf{Diag}(\mathbf{x})$ .

**Calculus:** For  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  differentiable at  $\mathbf{x} \in \mathbb{R}^n$ , we let  $\nabla f(\mathbf{x}) \in \mathbb{R}^n$  denote the gradient of  $f$  at  $\mathbf{x}$ , i.e.  $[\nabla f(\mathbf{x})]_i = \frac{\partial}{\partial x_i} f(\mathbf{x})$  for all  $i \in [n]$ . For  $f \in \mathbb{R}^n \rightarrow \mathbb{R}$  twice differentiable at  $\mathbf{x} \in \mathbb{R}^n$ , we let  $\nabla^2 f(\mathbf{x})$  denote the Hessian of  $f$  at  $\mathbf{x}$ , i.e.  $[\nabla^2 f(\mathbf{x})]_{ij} = \frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x})$  for all  $i, j \in [n]$ . Sometimes, we will consider functions of two vectors,  $g : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}$ , and wish to compute the gradient and Hessian of  $g$  restricted to one of the two vectors. For  $\mathbf{x} \in \mathbb{R}^{n_1}$  and  $\mathbf{y} \in \mathbb{R}^{n_2}$  we let  $\nabla_{\mathbf{x}} \mathbf{g}(\mathbf{a}, \mathbf{b}) \in \mathbb{R}^{n_1}$  denote the gradient of  $\mathbf{g}$  for fixed  $\mathbf{y}$  at point  $(\mathbf{a}, \mathbf{b}) \in \mathbb{R}^{n_1 \times n_2}$ . We define  $\nabla_{\mathbf{y}}$ ,  $\nabla_{\mathbf{xx}}^2$ , and  $\nabla_{\mathbf{yy}}^2$  similarly. For  $h : \mathbb{R}^n \rightarrow \mathbb{R}^m$  differentiable at  $\mathbf{x} \in \mathbb{R}^n$  we let  $\mathbf{J}(\mathbf{h}(\mathbf{x})) \in \mathbb{R}^{m \times n}$  denote the Jacobian of  $\mathbf{h}$  at  $\mathbf{x}$  where for all  $i \in [m]$  and  $j \in [n]$  we let  $[\mathbf{J}(\mathbf{h}(\mathbf{x}))]_{ij} \stackrel{\text{def}}{=} \frac{\partial}{\partial x_j} h(\mathbf{x})_i$ . For functions of multiple vectors we use subscripts, e.g.  $\mathbf{J}_{\mathbf{x}}$ , to denote the Jacobian of the function restricted to the  $\mathbf{x}$  variable.

**Sets:** We call a set  $U$  is symmetric if  $\mathbf{x} \in \mathbb{R}^n \Leftrightarrow -\mathbf{x} \in \mathbb{R}^n$ . For any  $\alpha > 0$  and set  $U \subseteq \mathbb{R}^n$  we let  $\alpha U \stackrel{\text{def}}{=} \{\mathbf{x} \in \mathbb{R}^n \mid \alpha^{-1} \mathbf{x} \in U\}$ . For any  $p \in [1, \infty]$  and  $r > 0$  we refer to the set  $B_p(r) \stackrel{\text{def}}{=} \{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\|_p \leq r\}$  as the  $\ell_p$  ball of radius  $r$ . For brevity, we refer to  $B_2(r)$  as a ball of radius  $r$  and  $B_\infty(r)$  as a box of radius  $r$ .

**Matrix Multiplication Constant:** We let  $\omega < 2.3728639$  [100] denote the matrix multiplication constant. It is known that the product a  $m \times n$  matrix with a  $n \times r$  matrix can be computed in  $mnr \min(m, n, r)^{\omega-3+o(1)}$  time and the inverse of a  $n \times n$  matrix can be computed in  $n^{\omega+o(1)}$  time. For simplicity, we always write  $\omega + o(1)$  as  $\omega$ .

**Running Time:** Unless mentioned specifically, we will disregard the issue of rounding error if the algorithm is obviously stable. Most of the algorithms we proposed are randomized algorithms that succeeds with high probability and the running time is often the expected running time. For any function  $f$ , we write  $\tilde{O}(f) \triangleq O(f \cdot \log^{O(1)} f)$ . Suppose  $m$  is the length of the input and  $n$  is the number of variables, we say that an algorithm runs in nearly-linear time if its running time is  $\tilde{O}(m)$ , runs in almost linear time if its running time is  $m^{1+o(1)}$ , and runs in input sparsity time if its running time is  $\tilde{O}(m + n^{O(1)})$ .

### ■ 2.3.1 Sparsification and Leverage Score

In this section, we define spectral approximation and leverage score. They are two key concepts in this thesis, especially in Part I. First, we define leverage score which is a measure of importance. Given a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . The leverage score of the  $i^{\text{th}}$  row  $\mathbf{a}_i$  of  $\mathbf{A}$  is:

$$\sigma_i(\mathbf{A}) \stackrel{\text{def}}{=} \mathbf{a}_i^\top (\mathbf{A}^\top \mathbf{A})^+ \mathbf{a}_i.$$

We also define the related *cross leverage score* as  $\sigma_{ij}(\mathbf{A}) \stackrel{\text{def}}{=} \mathbf{a}_i^\top (\mathbf{A}^\top \mathbf{A})^+ \mathbf{a}_j$ . Let  $\boldsymbol{\sigma}(\mathbf{A})$  be a vector containing  $\mathbf{A}$ 's  $m$  leverage scores.  $\boldsymbol{\sigma}(\mathbf{A})$  is the diagonal of  $\mathbf{A}(\mathbf{A}^\top \mathbf{A})^+ \mathbf{A}^\top$ , which is a projection matrix. Thus,  $\sigma_i(\mathbf{A}) = \mathbf{1}_i^\top \mathbf{A}(\mathbf{A}^\top \mathbf{A})^+ \mathbf{A}^\top \mathbf{1}_i \leq 1$ . Furthermore, since  $\mathbf{A}(\mathbf{A}^\top \mathbf{A})^+ \mathbf{A}^\top$  is a projection matrix, the sum of  $\mathbf{A}$ 's leverage scores is equal to the matrix's rank:

$$\sum_{i=1}^m \sigma_i(\mathbf{A}) = \text{Tr}(\mathbf{A}(\mathbf{A}^\top \mathbf{A})^+ \mathbf{A}^\top) = \text{rank}(\mathbf{A}(\mathbf{A}^\top \mathbf{A})^+ \mathbf{A}^\top) = \text{rank}(\mathbf{A}) \leq n. \quad (2.2)$$

A row's leverage score measures how important it is in composing the row space of  $\mathbf{A}$ . If a row has a component orthogonal to all other rows, its leverage score is 1. Removing it would decrease the rank of  $\mathbf{A}$ , completely changing its row space. The *coherence* of  $\mathbf{A}$  is  $\|\boldsymbol{\sigma}(\mathbf{A})\|_\infty$ . If  $\mathbf{A}$  has low coherence, no particular row is especially important. If  $\mathbf{A}$  has high coherence, it contains at least one row whose removal would significantly affect the composition of  $\mathbf{A}$ 's row space. Two characterizations that help with this intuition follow:

**Lemma 2.3.1.** *For all  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $i \in [m]$  we have that*

$$\sigma_i(\mathbf{A}) = \min_{\mathbf{A}^\top \mathbf{x} = \mathbf{a}_i} \|\mathbf{x}\|_2^2$$

where  $\mathbf{a}_i$  is the  $i^{\text{th}}$  row of  $\mathbf{A}$ . Let  $\mathbf{x}_i$  denote the optimal  $\mathbf{x}$  for  $\mathbf{a}_i$ . The  $j^{\text{th}}$  entry of  $\mathbf{x}_i$  is given by  $\sigma_{ij}(\mathbf{A})$ .

**Lemma 2.3.2.** *For all  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $i \in [m]$ , we have that  $\sigma_i(\mathbf{A})$  is the smallest  $t$  such that*

$$\mathbf{a}_i \mathbf{a}_i^\top \preceq t \cdot \mathbf{A}^\top \mathbf{A}. \quad (2.3)$$

Sampling rows from  $\mathbf{A}$  according to their exact leverage scores gives a spectral approximation for  $\mathbf{A}$  with high probability. Sampling by leverage score overestimates also suffices. The following lemma is a corollary of matrix concentration result [251, Cor 5.2] and the inequality (2.3).

**Lemma 2.3.3.** *Given an error parameter  $0 < \varepsilon < 1$  and a vector  $\mathbf{u}$  of leverage score overestimates, i.e.,  $\sigma_i(\mathbf{A}) \leq \mathbf{u}_i$  for all  $i$ . For any large enough universal constant  $c$ , we define a probability distribution  $\mathbf{p}_i = \min\{1, c\varepsilon^{-2}\mathbf{u}_i \log n\}$ . Let  $\mathbf{S}$  be the diagonal matrix with independently chosen entries.  $\mathbf{S}_{ii} = \frac{1}{\sqrt{\mathbf{p}_i}}$  with probability  $\mathbf{p}_i$  and 0 otherwise. Then, with probability at least  $1 - n^{-\Theta(c)}$ , we have*

$$(1 - \varepsilon)\mathbf{A}^\top \mathbf{A} \preceq \mathbf{A}^\top \mathbf{S}^2 \mathbf{A} \preceq (1 + \varepsilon)\mathbf{A}^\top \mathbf{A}.$$

Although (2.2) shows that  $\|\boldsymbol{\sigma}\|_1$  is small, computing  $\boldsymbol{\sigma}$  exactly is too expensive for many purposes. In [242], they showed that we can compute leverage scores,  $\boldsymbol{\sigma}$ , approximately by solving only polylogarithmically many regression problems. This result uses the fact that

$$\sigma_i(\mathbf{A}) = \|\mathbf{A}(\mathbf{A}^\top \mathbf{A})^+ \mathbf{A}^\top \mathbf{1}_i\|_2$$

and that by the Johnson-Lindenstrauss Lemma these lengths are preserved up to multiplicative error if we project these vectors onto certain random low dimensional subspaces.

**Lemma 2.3.4** ([242]). *Let  $\mathcal{T}$  be the time to apply  $(\mathbf{A}^\top \mathbf{A})^+$  to an arbitrary vector. Then, for any  $0 < \theta < 1$ , it is possible to compute an estimate of  $\sigma(\mathbf{A})$ ,  $\boldsymbol{\sigma}^{(app)}$ , in time  $O(\theta^{-1}(\mathcal{T} + \text{nnz}(\mathbf{A})))$  such that w.h.p. in  $n$ , for all  $i$ ,*

$$d^{-\max(\sqrt{\theta/\log n}, \theta)} \sigma_i(\mathbf{A}) \leq \sigma_i^{(app)} \leq d^{\max(\sqrt{\theta/\log n}, \theta)} \sigma_i(\mathbf{A}).$$

Setting  $\theta = \frac{\varepsilon^2}{\log n}$  gives  $(1 + \varepsilon)$  factor leverage score approximations in  $\tilde{O}((\mathcal{T} + \text{nnz}(\mathbf{A}))/\varepsilon^2)$  time. Therefore, approximating leverage scores is as simple as solving  $\tilde{O}(1)$  linear systems.

## ■ 2.3.2 Separation Oracle

In this section, we define separation oracles. They are the key concepts in Part II. Our definitions are possibly non-standard and chosen to handle the different settings that occur in this part.

**Definition 2.3.5** (Separation Oracle for a Set). Given a set  $K \subset \mathbb{R}^n$  and  $\delta \geq 0$ , a  $\delta$ -separation oracle for  $K$  is a function on  $\mathbb{R}^n$  such that for any input  $\mathbf{x} \in \mathbb{R}^n$ , it either outputs “successful” or a half space of the form  $H = \{\mathbf{z} : \mathbf{c}^T \mathbf{z} \leq \mathbf{c}^T \mathbf{x} + b\} \supseteq K$  with  $b \leq \delta \|\mathbf{c}\|_2$  and  $\mathbf{c} \neq \mathbf{0}$ . We let  $\text{SO}_\delta(K)$  be the time complexity of this oracle.

The parameter  $\delta$  indicates the accuracy of the oracle. For brevity we refer to a 0-separation oracle for a set as just a *separation oracle*.

Note that in Definition 2.3.5 we do not assume that  $K$  is convex. However, we remark that there is a separation oracle for a set if and only if it is convex and that there is a  $\delta$  separation oracle if and only if the set is close to convex in some sense.

**Definition 2.3.6** (Separation Oracle for a Function). For any function  $f$ ,  $\eta \geq 0$  and  $\delta \geq 0$ , a  $(\eta, \delta)$ -separation oracle on a set  $\Gamma$  for  $f$  is a function on  $\mathbb{R}^n$  such that for any input  $\mathbf{x} \in \Gamma$ , it either asserts  $f(\mathbf{x}) \leq \min_{\mathbf{y} \in \Gamma} f(\mathbf{y}) + \eta$  or outputs a half space  $H$  such that

$$\{\mathbf{z} \in \Gamma : f(\mathbf{z}) \leq f(\mathbf{x})\} \subset H \stackrel{\text{def}}{=} \{\mathbf{z} : \mathbf{c}^T \mathbf{z} \leq \mathbf{c}^T \mathbf{x} + b\} \quad (2.4)$$

with  $b \leq \delta \|\mathbf{c}\|$  and  $\mathbf{c} \neq \mathbf{0}$ . We let  $\text{SO}_{\eta, \delta}(f)$  be the time complexity of this oracle.

### ■ 2.3.3 Other Basic Facts

Here, we state some lemmas that are used throughout this thesis.

**Lemma 2.3.7** (Sherman-Morrison Formula, [193, Thm 3]). Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , and  $\mathbf{u}, \mathbf{v} \perp \ker(\mathbf{A})$ . Suppose that  $1 + \mathbf{v}^T \mathbf{A}^+ \mathbf{u} \neq 0$ . Then, it holds that

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^T)^+ = \mathbf{A}^+ - \frac{\mathbf{A}^+ \mathbf{u}\mathbf{v}^T \mathbf{A}^+}{1 + \mathbf{v}^T \mathbf{A}^+ \mathbf{u}}.$$

**Lemma 2.3.8.** Let  $\mathbf{B}_t$  be a differentiable family of invertible matrix. Then, we have

1.  $\frac{d}{dt} \ln \det(\mathbf{B}_t) = \text{Tr}(\mathbf{B}_t^{-1} \frac{d}{dt} \mathbf{B}_t)$ .
2.  $\frac{d}{dt} \mathbf{B}_t^{-1} = -\mathbf{B}_t^{-1} (\frac{d}{dt} \mathbf{B}_t) \mathbf{B}_t^{-1}$ .

**Lemma 2.3.9.** Let  $\mathbf{f} : \mathbb{R}^{n+m} \rightarrow \mathbb{R}^m$  be a continuously differentiable function. Suppose that  $\mathbf{f}(\mathbf{x}, \mathbf{y}) = \mathbf{0}$  for some  $\mathbf{x} \in \mathbb{R}^m$  and  $\mathbf{y} \in \mathbb{R}^n$  and that  $\frac{\partial \mathbf{f}}{\partial \mathbf{y}}(\mathbf{x}, \mathbf{y})$  is invertible. Then, there is an open set  $U$  containing  $\mathbf{x}$ , an open set  $V$  containing  $\mathbf{y}$ , and a unique continuously differentiable function  $\mathbf{g} : U \rightarrow V$  such that  $\mathbf{f}(\mathbf{x}, \mathbf{g}(\mathbf{x})) = \mathbf{0}$ . Also, we have that

$$\frac{\partial \mathbf{g}}{\partial \mathbf{x}}(\mathbf{x}) = - \left( \frac{\partial \mathbf{f}}{\partial \mathbf{y}}(\mathbf{x}, \mathbf{g}(\mathbf{x})) \right)^{-1} \frac{\partial \mathbf{f}}{\partial \mathbf{x}}(\mathbf{x}, \mathbf{g}(\mathbf{x})).$$

**Lemma 2.3.10** (Matrix Chernoff Bound, [251, Cor 5.2]). Let  $\{\mathbf{X}_k\}$  be a finite sequence of independent, random, and self-adjoint matrices with dimension  $n$ . Assume that each random matrix satisfies  $\mathbf{X}_k \succeq \mathbf{0}$ , and  $\lambda_{\max}(\mathbf{X}_k) \leq R$ . Let  $\mu_{\max} \geq \lambda_{\max}(\sum_k \mathbf{E}\mathbf{X}_k)$  and  $\mu_{\min} \leq \lambda_{\min}(\sum_k \mathbf{E}\mathbf{X}_k)$ . Then, it holds that

$$\begin{aligned} \Pr \left[ \lambda_{\max} \left( \sum_k \mathbf{X}_k \right) \geq (1 + \delta) \mu_{\max} \right] &\leq n \cdot \left( \frac{e^\delta}{(1 + \delta)^{1 + \delta}} \right)^{\mu_{\max}/R} \quad \text{for any } \delta \geq 0, \text{ and} \\ \Pr \left[ \lambda_{\min} \left( \sum_k \mathbf{X}_k \right) \leq (1 - \delta) \mu_{\min} \right] &\leq n \cdot \left( \frac{e^{-\delta}}{(1 - \delta)^{1 - \delta}} \right)^{\mu_{\min}/R} \quad \text{for any } \delta \in [0, 1] \end{aligned}$$

**Lemma 2.3.11** (Lieb Thirring Inequality, [171]). *Let  $\mathbf{A}$  and  $\mathbf{B}$  be positive definite matrices and  $q \geq 1$ . Then it holds that*

$$\mathrm{Tr}(\mathbf{B}\mathbf{A}\mathbf{B})^q \leq \mathrm{Tr}(\mathbf{B}^q\mathbf{A}^q\mathbf{B}^q).$$

**Theorem 2.3.12** (Simple Constrained Minimization for Twice Differentiable Function [206]). *Let  $\mathbf{H}$  be a positive definite matrix,  $K$  be a convex set and  $L \geq \mu \geq 0$ . Let  $f(\mathbf{x}) : K \rightarrow \mathbb{R}$  be a twice differentiable function such that  $\mu\mathbf{H} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbf{H}$  for all  $\mathbf{x} \in K$ . If for some  $\mathbf{x}^{(0)} \in K$  and all  $k \geq 0$ , we apply the update rule*

$$\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x} \in K} \left\langle \nabla f(\mathbf{x}^{(k)}), \mathbf{x} - \mathbf{x}^{(k)} \right\rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{x}^{(k)}\|_{\mathbf{H}}^2$$

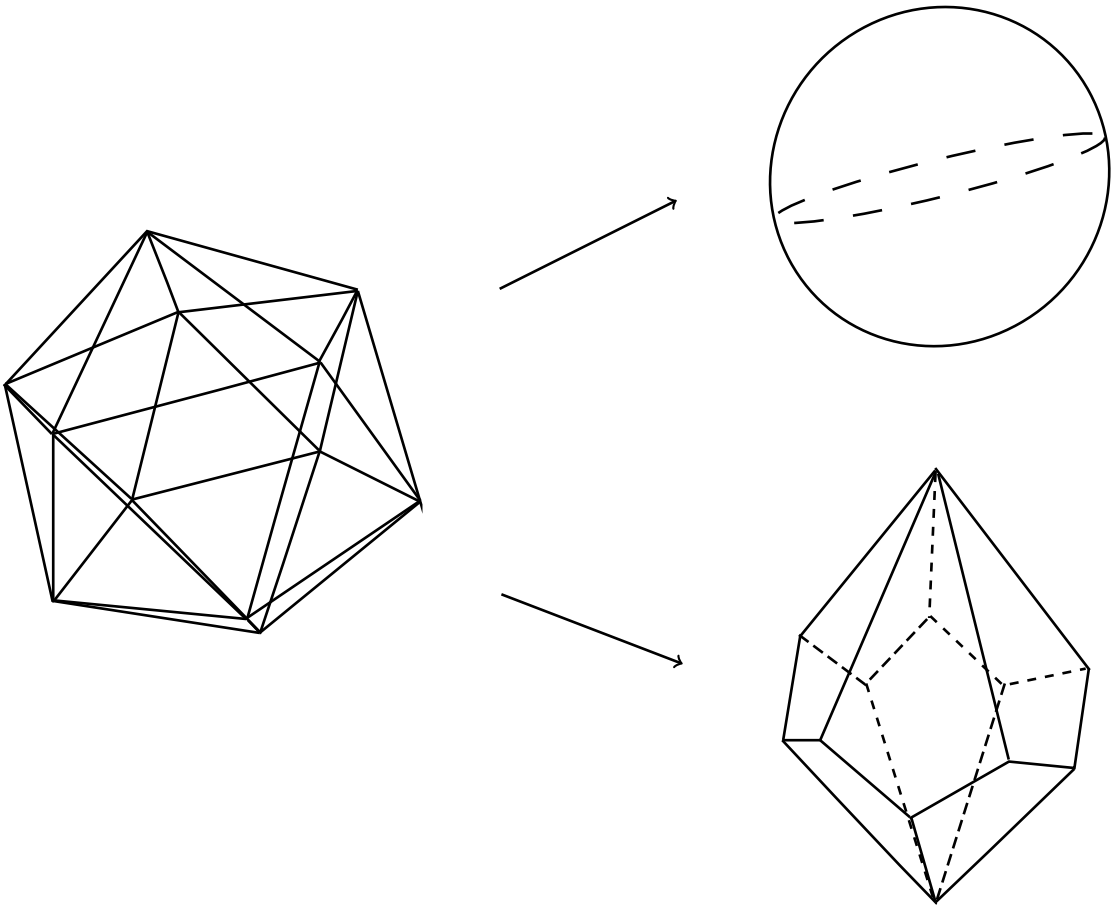
*then for all  $k \geq 0$  we have*

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\mathbf{H}}^2 \leq \left(1 - \frac{\mu}{L}\right)^k \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_{\mathbf{H}}^2.$$

# Part I

---

# Sparsification





## $\ell_2$ Regression In Input Sparsity Time

### ■ 3.1 Introduction

Many convex optimization problems can be written in the form of

$$\min_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^m f_i(\mathbf{a}_i^\top \mathbf{x})$$

where  $f_i$  are 1 dimension convex functions and  $\mathbf{a}_i \in \mathbb{R}^n$ . The goal of this part is to show that the running time for solving this type of problems, in general, should depends on  $m$  linearly.

In this chapter, we first study the simplest case,  $f_i$  being quadratic functions. In this case, the problem becomes the  $\ell_2$  regression problem  $\min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|$ . Currently, there are two types of methods on this problem for the case  $m \gg n$ .

The first type of methods rely on *spectral approximation*. For a tall, narrow data matrix  $\mathbf{A}$ , these methods find a shorter approximate data matrix,  $\tilde{\mathbf{A}}$ , such that, for all vectors  $\mathbf{x}$ ,  $\|\tilde{\mathbf{A}}\mathbf{x}\|_2 \approx \|\mathbf{A}\mathbf{x}\|_2$ . All of which rely on variations of Johnson-Lindenstrauss random projections for constructing  $\tilde{\mathbf{A}}$  [53, 185, 200, 170]. Such random projections do not maintain sparsity and row structure of  $\mathbf{A}$ . The second type of methods rely on uniform (or non-adaptive) sampling, such as stochastic descent. These methods works only for *low coherence* data and hence are less useful for theory development.

By re-examining uniform sampling, we give spectral approximation algorithms that avoid random projection entirely. Our methods are the first to match state-of-the-art runtimes while preserving row structure and sparsity in all matrix operations. Since it preserves row structure, this allows us to take advantage of the problem structure. Let  $\mathcal{T}(m, n)$  be the time needed to solve a  $\ell_2$  regression problem  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|$  where  $\mathbf{A}$  is a  $m \times n$  matrix and each rows of  $\mathbf{A}$  comes from a certain row dictionary. Roughly speaking, our method shows that

$$\mathcal{T}(m, n) = \tilde{O}(\text{nnz}(\mathbf{A})) + \tilde{O}(\mathcal{T}(\tilde{O}(n), n)).$$

This powerful reduction allows us to develop the first nearly linear time algorithm for connection Laplacian (Chapter 13). In comparison, the random projection methods only shows that  $\mathcal{T}(m, n) = \tilde{O}(\text{nnz}(\mathbf{A}) + n^\omega)$ .

As we explained in Section 2.3.1, it is known that for a data matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , a spectral approximation can be obtained by sampling  $O(n \log n)$  rows, each with probability proportional to its *statistical leverage score*. Recall that the leverage score of  $\mathbf{A}$ 's  $i^{\text{th}}$  row,  $\mathbf{a}_i$ , is  $\sigma_i = \mathbf{a}_i^\top (\mathbf{A}^\top \mathbf{A})^+ \mathbf{a}_i$ , where  $\mathbf{M}^+$  denotes the Moore-Penrose pseudoinverse of  $\mathbf{M}$ . A higher leverage score indicates that  $\mathbf{a}_i$  is more important in composing the spectrum of  $\mathbf{A}$ .

Unfortunately, leverage scores are difficult to calculate – finding them involves computing the pseudoinverse  $(\mathbf{A}^\top \mathbf{A})^+$ , which is as slow as solving our regression problem in the first place. In practice, data is often assumed to have *low coherence* [195], in which case simply selecting rows uniformly at random works [22, 156]. However, uniform sampling could be disastrous — if  $\mathbf{A}$  contains

a row with some component orthogonal to all other rows, removing it will reduce the rank of  $\mathbf{A}$  and thus we cannot possibly preserve all vector products ( $\|\tilde{\mathbf{A}}\mathbf{x}\|_2$  will start sending some vectors to 0). Any uniform sampling scheme is likely to drop any such single row. On the other hand, when leverage score sampling, such a row would have the highest possible leverage score of 1.

Possible fixes include randomly “mixing” data points to avoid degeneracies [22]. However, this approach sacrifices sparsity and structure in our data matrix, increasing storage and runtime costs. Is there a more elegant fix? First note that sampling  $\mathbf{A}$  by *approximate* leverage scores is fine, but we may need to select more than the optimal  $O(n \log n)$  rows. With that in mind, consider the following straightforward algorithm for iterative sampling, inspired by [170]:

1. **Reduce  $\mathbf{A}$  significantly by sampling rows uniformly.**
2. **Use the smaller matrix to approximate  $(\mathbf{A}^\top \mathbf{A})^+$  and leverage scores for  $\mathbf{A}$ .**
3. **Resample rows of  $\mathbf{A}$  using these estimates to obtain a spectral approximation  $\tilde{\mathbf{A}}$ .**

While intuitive, this scheme was not previously known to work. Our main result is proving that it does. This process will quickly converge on a small spectral approximation to  $\mathbf{A}$ , i.e. with  $O(n \log n)$  rows.

A few results come close to an analysis of such a routine – in particular, two iterative sampling schemes are analyzed in [170]. However, the first ultimately requires Johnson-Lindenstrauss projections that mix rows, something we were hoping to avoid. The second almost maintains sparsity and row structure (except for possibly including rows of the identity in  $\tilde{\mathbf{A}}$ ), but its convergence rate depends on the condition number of  $\mathbf{A}$ .

More importantly, both of these results are similar in that they rely on the primitive that a (possibly poor) spectral approximation to  $\mathbf{A}$  is sufficient for approximately computing leverage scores, which are in turn good enough for obtaining an even better spectral approximation. As mentioned, uniform sampling will not in general give a spectral approximation – it does not preserve information about all singular values. Our key contribution is a better understanding of what information uniform sampling *does* preserve. It turns out that, although weaker than a spectral approximation, the matrix obtained from uniform sampling can nonetheless give leverage score estimates that are good enough to obtain increasingly better approximations to  $\mathbf{A}$ .

### ■ 3.1.1 Our Approach

Suppose we compute a set of leverage score estimates,  $\{\tilde{\sigma}_i\}$ , using  $(\tilde{\mathbf{A}}^\top \tilde{\mathbf{A}})^+$  in place of  $(\mathbf{A}^\top \mathbf{A})^+$  for some already obtained matrix approximation  $\tilde{\mathbf{A}}$ . As long as our leverage score approximations are *upper bounds* on the true scores ( $\tilde{\sigma}_i \geq \sigma_i$ ) we can use them for sampling and still obtain a spectral approximation to  $\mathbf{A}$ . The number of samples we take will be

$$c \cdot \log n \cdot \sum_{i=1}^m \tilde{\sigma}_i$$

where  $c$  is some fixed constant. When sampling by exact leverage scores, it can be shown that  $\sum_i \sigma_i \leq n$  so we take  $c \cdot n \log n$  rows.

Thus, to prove that our proposed iterative algorithm works, we need to show that, if we uniformly sample a relatively small number of rows from  $\mathbf{A}$  (Step 1) and estimate leverage scores using these rows (Step 2), then the sum of our estimates will be small. Then, when we sample by these estimated leverage scores in Step 3, we can sufficiently reduce the size of  $\mathbf{A}$ .



In prior work, the sum of overestimates was bounded by estimating *each* leverage score to within a multiplicative factor. This requires a spectral approximation, which is why previous iterative sampling schemes could only boost poor spectral approximations to better spectral approximations. Of course, a “for each” statement is not required, and we will not get one through uniform sampling. Thus, our core result avoids this technique. Specifically, we show,

**Theorem 3.1.1.** *For any  $k$ , we can select  $k$  rows uniformly at random from  $\mathbf{A}$  to obtain  $\tilde{\mathbf{A}}$ . Let  $\{\tilde{\sigma}_i\}$  be a set of leverage score estimates for  $\mathbf{A}$  computed using  $\tilde{\mathbf{A}}$ <sup>1</sup>. Then, we have that  $\tilde{\sigma}_i \geq \sigma_i$  for all  $i$  and*

$$\mathbf{E} \left[ \sum_{i=1}^m \tilde{\sigma}_i \right] \leq \frac{mn}{k}.$$

The validity of our proposed iterative sampling scheme immediately follows from Theorem 3.1.1. For example, if we uniformly sample  $k = m/2$  rows then  $c \log n \sum \tilde{\sigma}_i \leq O(n \log n)$ , so we can cut our matrix down to  $O(n \log n)$  rows. There is a convenient tradeoff – the more rows uniformly sampled in Step 1, the more we can cut  $\mathbf{A}$  down by in Step 3. This tradeoff leads to natural recursive and iterative algorithms for row sampling.

We give a proof of Theorem 3.1.1 using a simple expectation argument that bounds  $\mathbf{E} \sum \tilde{\sigma}_i$ . We also prove alternative versions of Theorem 3.1.1 with slightly different guarantees (Theorems 3.6.1) using a technique that we believe is of independent interest. It is well known that, if  $\mathbf{A}$  has low coherence – that is, has a low maximum leverage score – then uniform sampling from the matrix is actually sufficient for obtaining a full spectral approximation. The uniform rate will upper bound the leverage score rate for every row. With this in mind, we show a powerful fact: while many matrices do not have low coherence, for any  $\mathbf{A}$ , we can decrease the weight on a small subset of rows to make the matrix have low coherence. Specifically,

**Lemma 3.1.2** (Coherence Reducing Reweighting). *For any  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and any coherence upper bound  $\alpha > 0$ , there exists a diagonal reweighting matrix  $\mathbf{W} \in \mathbb{R}^{m \times m}$  with all entries in  $[0, 1]$  and just  $(n/\alpha)$  entries not equal to 1, such that:*

$$\forall i, \sigma_i(\mathbf{W}\mathbf{A}) \leq \alpha.$$

Intuitively, this lemma shows that uniform sampling gives a matrix that spectrally approximates a large sub-matrix of the original data. It follows from our more general Theorem 3.5.1, which describes exactly how leverage scores of  $\mathbf{A}$  can be manipulated through row reweighting.

We never actually find  $\mathbf{W}$  explicitly – simply its existence implies our uniform sampling theorems. As explained, since  $\mathbf{W}\mathbf{A}$  has low coherence, uniform sampling would give a spectral approximation to the reweighted matrix and thus a multiplicatively good approximation to *each* leverage score. Thus, the sum of estimated leverage scores for  $\mathbf{W}\mathbf{A}$  will be low, i.e.  $< O(n)$ . It can be shown that, for any row that is not reweighted, the leverage score in  $\mathbf{A}$  computed using a uniformly sampled  $\tilde{\mathbf{A}}$ , is never greater than the corresponding leverage score in  $\mathbf{W}\mathbf{A}$  computed using a uniformly sampled  $\tilde{\mathbf{W}\mathbf{A}}$ . Thus, the sum of approximate leverage scores for rows in  $\mathbf{A}$  that are not reweighted is small by comparison to their corresponding leverage scores in  $\mathbf{W}\mathbf{A}$ . How about the rows that *are* reweighted in  $\mathbf{W}\mathbf{A}$ ? Lemma 3.1.2 claims there are not too many of these – we can trivially bound their leverage score estimates by 1 and even then the total sum of estimated leverage scores will be small.

This argument gives the result we need: even if a uniformly sampled  $\tilde{\mathbf{A}}$  cannot be used to obtain good per row leverage score upper bounds, it is sufficient for ensuring that the sum of all leverage score estimates is not too high.

<sup>1</sup>We describe exactly how each  $\tilde{\sigma}_i$  is computed when we prove Theorem 3.1.1 in Section 3.4.

## ■ 3.2 Previous Work

### ■ 3.2.1 Randomized Numerical Linear Algebra

In the past decade, fast randomized algorithms for matrix problems have risen to prominence. Numerous results give improved running times for matrix multiplication, linear regression, and low rank approximation – helpful surveys of this work include [184] and [112]. In addition to asymptotic runtime gains, randomized alternatives to standard linear algebra tools tend to offer significant gains in terms of data access patterns and required working memory.

Algorithms for randomized linear algebra often work by generically reducing problem size – large matrices are compressed (using randomness) to smaller approximations which are processed deterministically via standard linear algebraic methods. Methods for matrix reduction divide roughly into two categories – random projection methods [54, 53, 185, 200, 229] and sampling methods [70, 71, 72, 74, 73, 170].

Random projection methods recombine rows or columns from a large matrix to form a much smaller problem that approximates the original. Descending from the Johnson-Lindenstrauss Lemma [130] and related results, these algorithms are impressive for their simplicity and speed – reducing a large matrix simply requires multiplication by an appropriately chosen random matrix.

Sampling methods, on the other hand, seek to approximate large matrices by judiciously selecting (and reweighting) few rows or columns. Sampling itself is even simpler and faster than random projection – the challenge becomes efficiently computing the correct measure of “importance” for rows or columns. More important rows or columns are selected with higher probability.

### ■ 3.2.2 Approximate Linear Regression

We focus on linear regression, i.e. solving overdetermined systems, which requires our matrix reduction step to produce a spectral approximation  $\tilde{\mathbf{A}}$  to the data matrix  $\mathbf{A}$ . One possibility is to obtain a  $\varepsilon$ -approximation with  $O(n \log n / \varepsilon^2)$  rows and to solve regression on the smaller problem to give an approximate solution. To improve stability and achieve  $\log(1/\varepsilon)$  dependence, randomized schemes can be combined with known iterative regression algorithms. These methods only require a constant factor spectral approximation with  $O(n \log n)$  rows [22, 55, 53, 192, 228].

When random projections are used,  $\tilde{\mathbf{A}} = \mathbf{\Pi A}$  for some randomly generated matrix  $\mathbf{\Pi}$  which is known as a *subspace embedding*. Recent progress has significantly sped up the process of computing  $\mathbf{\Pi A}$ , leading to the first *input-sparsity time* algorithms for linear regression (or nearly input-sparsity time if iterative methods are employed) [53, 185, 200].

### ■ 3.2.3 Row Sampling

An alternative route to spectral approximation is importance sampling. Specifically,  $O(n \log n / \varepsilon^2)$  rows can be sampled with probability proportional to their leverage scores [75, 76, 242]. In [242], Spielman and Srivastava specifically focus on spectral approximations for the edge-vertex incidence matrix of a graph. This is more commonly referred to as *spectral graph sparsification*, a primitive that has become important in research on graph algorithms. Each row in a graph’s (potentially tall) edge-vertex incident matrix corresponds to an edge and the row’s leverage score is exactly the edge’s *weighted effective resistance*, which is used as the sampling probability in [242].

This application illustrates an important point: for spectral sparsification, it is critical that  $\mathbf{A}$  is compressed via sampling instead of random projection. Sampling ensures that  $\tilde{\mathbf{A}}$  contains only reweighted rows from  $\mathbf{A}$ , so it remains an edge-vertex incidence matrix. In general, sampling is

interesting because it preserves row structure. Even if that structure is just a certain level of sparsity, it can reduce memory requirements and accelerate matrix operations.

While leverage scores for the edge-vertex incidence matrix of a graph can be computed quickly [147, 243], in general, computing leverage scores requires evaluating  $(\mathbf{A}^\top \mathbf{A})^+$ , which is as difficult as solving regression in the first place. Li, Miller, and Peng address this issue with methods for iteratively computing good row samples [170]. Their algorithms achieve input-sparsity time regression, but are fairly involved and rely on intermediate operations that ultimately require Johnson-Lindenstrauss projections, mixing rows and necessitating dense matrix operations. An alternative approach from [170] does preserve row structure (except for possible additions of rows from the identity to intermediate matrices) but converges in a number of steps that depends on  $\mathbf{A}$ 's condition number.

### ■ 3.3 Generalized Leverage Score

We often approximate the leverage scores of  $\mathbf{A}$  by computing them with respect to some other matrix  $\mathbf{B} \in \mathbb{R}^{m' \times n}$ . We define the *generalized leverage score*:

$$\sigma_i^{\mathbf{B}}(\mathbf{A}) \stackrel{\text{def}}{=} \begin{cases} \mathbf{A}_i^\top (\mathbf{B}^\top \mathbf{B})^+ \mathbf{A}_i & \text{if } \mathbf{A}_i \perp \ker(\mathbf{B}), \\ \infty & \text{otherwise.} \end{cases} \quad (3.1)$$

If  $\mathbf{A}_i$  has an component in  $\ker(\mathbf{B})$ , we set its generalized leverage score to  $\infty$ , since it might be the only row in  $\mathbf{A}$  pointing in this direction. Thus, when sampling rows, we cannot remove it. We could set the generalized leverage score to 1, but using  $\infty$  simplifies notation in some of our proofs. If  $\mathbf{B}$  is a spectral approximation for  $\mathbf{A}$ , then every generalized leverage score is a good multiplicative approximation to its corresponding true leverage score:

**Lemma 3.3.1** ([170, Lem 4.3]). *If  $\frac{1}{\lambda} \mathbf{A}^\top \mathbf{A} \preceq \mathbf{B}^\top \mathbf{B} \preceq \mathbf{A}^\top \mathbf{A}$ , then  $\sigma_i(\mathbf{A}) \leq \sigma_i^{\mathbf{B}}(\mathbf{A}) \leq \lambda \cdot \sigma_i(\mathbf{A})$ .*

### ■ 3.4 Expectation Bound on Leverage Score Estimation

In this section, we use a simple expectation argument to prove Theorem 3.1.1, which is restated in full below:

*Theorem 3.1.1* (Full Statement). Given any  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Let  $S$  denote a uniformly random sample of  $k$  rows from  $\mathbf{A}$  and let  $\mathbf{S} \in \mathbb{R}^{m \times m}$  be its diagonal indicator matrix (i.e.  $\mathbf{S}_{ii} = 1$  for  $i \in S$ ,  $\mathbf{S}_{ii} = 0$  otherwise). Define

$$\tilde{\sigma}_i \stackrel{\text{def}}{=} \begin{cases} \sigma_i^{\mathbf{S}\mathbf{A}}(\mathbf{A}) & \text{if } i \in S, \\ \frac{1}{1 + \frac{1}{\sigma_i^{\mathbf{S}\mathbf{A}}(\mathbf{A})}} & \text{otherwise.} \end{cases}$$

Then,  $\tilde{\sigma}_i \geq \sigma_i(\mathbf{A})$  for all  $i$  and

$$\mathbf{E} \left[ \sum_{i=1}^m \tilde{\sigma}_i \right] \leq \frac{mn}{k}.$$

*Proof.* First we show that our estimates are valid leverage score upper bounds, i.e.  $\tilde{\sigma}_i \geq \sigma_i(\mathbf{A})$ . Let  $\mathbf{S}^{(i)}$  be the diagonal indicator matrix for  $S \cup \{i\}$ . We claim that, for all  $i$ ,

$$\tilde{\sigma}_i = \sigma_i^{\mathbf{S}^{(i)}\mathbf{A}}(\mathbf{A}). \quad (3.2)$$

This is proved case-by-case:

When  $i \in S$ ,  $\mathbf{S} = \mathbf{S}^{(i)}$  so it holds trivially.

When  $i \notin S$  and  $\mathbf{a}_i \notin \ker(\mathbf{S}\mathbf{A})$ , then by definition,  $\sigma_i^{\mathbf{S}\mathbf{A}}(\mathbf{A}) = \infty$  and  $\tilde{\sigma}_i = \frac{1}{1+\frac{1}{\infty}} = 1 = \sigma_i^{\mathbf{S}^{(i)}\mathbf{A}}(\mathbf{A})$ .

When  $i \notin S$  and  $\mathbf{a}_i \perp \ker(\mathbf{S}\mathbf{A})$  then by the Sherman-Morrison formula (Lemma 2.3.7),

$$\begin{aligned} \sigma_i^{\mathbf{S}^{(i)}\mathbf{A}}(\mathbf{A}) &= \mathbf{a}_i^\top \left( \mathbf{A}^\top \mathbf{S}^2 \mathbf{A} + \mathbf{a}_i \mathbf{a}_i^\top \right)^+ \mathbf{a}_i \\ &= \mathbf{a}_i^\top \left( \left( \mathbf{A}^\top \mathbf{S}^2 \mathbf{A} \right)^+ - \frac{\left( \mathbf{A}^\top \mathbf{S}^2 \mathbf{A} \right)^+ \mathbf{a}_i \mathbf{a}_i^\top \left( \mathbf{A}^\top \mathbf{S}^2 \mathbf{A} \right)^+}{1 + \mathbf{a}_i^\top \left( \mathbf{A}^\top \mathbf{S}^2 \mathbf{A} \right)^+ \mathbf{a}_i} \right) \mathbf{a}_i \\ &= \sigma_i^{\mathbf{S}\mathbf{A}}(\mathbf{A}) - \frac{\sigma_i^{\mathbf{S}\mathbf{A}}(\mathbf{A})^2}{1 + \sigma_i^{\mathbf{S}\mathbf{A}}(\mathbf{A})} = \frac{1}{1 + \frac{1}{\sigma_i^{\mathbf{S}\mathbf{A}}(\mathbf{A})}} = \tilde{\sigma}_i. \end{aligned}$$

By (3.2) and the fact that  $\mathbf{A}^\top \mathbf{S}^{(i)2} \mathbf{A} \preceq \mathbf{A}^\top \mathbf{A}$  (see Lemma 3.3.1), we have  $\tilde{\sigma}_i = \sigma_i^{\mathbf{S}^{(i)}\mathbf{A}}(\mathbf{A}) \geq \sigma_i(\mathbf{A})$ , so our estimates are upper bounds as desired. It remains to upper bound the expected sum of  $\tilde{\sigma}_i$ . We can break down the sum as:

$$\sum_{i=1}^m \tilde{\sigma}_i = \sum_{i \in S} \tilde{\sigma}_i + \sum_{i \notin S} \tilde{\sigma}_i.$$

The first term is simply the sum of  $\mathbf{S}\mathbf{A}$ 's leverage scores, so it is equal to  $\text{rank}(\mathbf{S}\mathbf{A}) \leq n$  by (2.2). To bound the second term, consider a random process that first selects  $\mathbf{S}$ , then selects a random row  $i \notin S$  and returns  $\tilde{\sigma}_i$ . There are always exactly  $m - k$  rows  $i \notin S$ , so the value returned by this random process is, in expectation, exactly equal to  $\frac{1}{m-k} \cdot \mathbf{E} \sum_{i \notin S} \tilde{\sigma}_i$ .

This random process is *also* equivalent to randomly selecting a set  $S'$  of  $k+1$  rows, then randomly choosing a row  $i \in S'$  and returning its leverage score! In expectation it is therefore equal to the average leverage score in  $\mathbf{S}'\mathbf{A}$ .  $\mathbf{S}'\mathbf{A}$  has  $k+1$  rows and its leverage scores sum to its rank, so we can bound its average leverage score by  $\frac{n}{k+1}$ . Overall:

$$\mathbf{E} \left[ \sum_{i=1}^m \tilde{\sigma}_i \right] \leq n + (m - k) \cdot \frac{n}{k+1} \leq \frac{n(m+1)}{k+1} \leq \frac{mn}{k}.$$

□

## ■ 3.5 Coherence-Reducing Weighting

In this section, we prove Theorem 3.5.1, which shows that we can reweight a small number of rows in any matrix  $\mathbf{A}$  to make it have low coherence. This structural result may be of independent interest. It is also fundamental in proving Theorem 3.6.1, a slightly stronger version of Theorem 3.1.1 that we will prove in Section 3.6.

Actually, for Theorem 3.5.1 we prove a more general statement, studying how to select a diagonal row reweighting matrix  $\mathbf{W}$  to arbitrarily control the leverage scores of  $\mathbf{W}\mathbf{A}$ . One simple conjecture would be that, given a vector  $\mathbf{u}$ , there always exists a  $\mathbf{W}$  such that  $\sigma_i(\mathbf{W}\mathbf{A}) = \mathbf{u}_i$ . This conjecture is unfortunately not true - if  $\mathbf{A}$  is the identity matrix, then  $\sigma_i(\mathbf{W}\mathbf{A}) = 0$  if  $\mathbf{W}_{ii} = 0$  and  $\sigma_i(\mathbf{W}\mathbf{A}) = 1$  otherwise. Instead, we show the following:

**Theorem 3.5.1.** *For any  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and any vector  $\mathbf{u} \in \mathbb{R}^m$  with  $\mathbf{u}_i > 0$  for all  $i$ , there exists a diagonal matrix  $\mathbf{W} \in \mathbb{R}^{m \times m}$  with  $\mathbf{0} \preceq \mathbf{W} \preceq \mathbf{I}$  such that:*

$$\forall i, \sigma_i(\mathbf{W}\mathbf{A}) \leq \mathbf{u}_i, \tag{3.3}$$

and

$$\sum_{i:W_{ii} \neq 1} \mathbf{u}_i \leq n. \quad (3.4)$$

Note that (3.3) is easy to satisfy – it holds if we set  $\mathbf{W} = \mathbf{0}$ . Hence, the main result is the second claim. Not only does a  $\mathbf{W}$  exist that gives the desired leverage score bounds, but it is only necessary to reweight rows in  $\mathbf{A}$  with a low total weight in terms of  $\mathbf{u}$ .

For any incoherence parameter  $\alpha$ , if we set  $\mathbf{u}_i = \alpha$  for all  $i$ , then this theorem shows the existence of a reweighting that reduces coherence to  $\alpha$ . Such a reweighting has  $\sum_{i:W_{ii} \neq 1} \alpha \leq n$  and therefore  $|\{i : W_{ii} \neq 1\}| \leq \frac{n}{\alpha}$ . So, we see that Lemma 3.1.2 follows as a special case of Theorem 3.5.1.

In order to prove Theorem 3.5.1, we first give two technical lemmas which are proved Section 3.8. Lemma 3.5.2 describes how the leverage scores of  $\mathbf{A}$  evolve when a single row of  $\mathbf{A}$  is reweighted. We show that, when we decrease the weight of a row, that row's leverage score decreases and the leverage score of all other rows increases.

**Lemma 3.5.2** (Leverage Score Changes Under Rank 1 Updates). *Given any  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\gamma \in (0, 1)$ , and  $i \in [m]$ , let  $\mathbf{W}$  be a diagonal matrix such that  $W_{ii} = \sqrt{1 - \gamma}$  and  $W_{jj} = 1$  for all  $j \neq i$ . Then,*

$$\sigma_i(\mathbf{W}\mathbf{A}) = \frac{(1 - \gamma) \sigma_i(\mathbf{A})}{1 - \gamma \sigma_i(\mathbf{A})} \leq \sigma_i(\mathbf{A}),$$

and for all  $j \neq i$ ,

$$\sigma_j(\mathbf{W}\mathbf{A}) = \sigma_j(\mathbf{A}) + \frac{\gamma \sigma_{ij}(\mathbf{A})^2}{1 - \gamma \sigma_i(\mathbf{A})} \geq \sigma_j(\mathbf{A}).$$

Next we claim that, with respect to weightings of  $\mathbf{A}$ 's rows, leverage scores are lower semi-continuous.

**Lemma 3.5.3** (Leverage Scores are Lower Semi-continuous).  *$\sigma(\mathbf{W}\mathbf{A})$  is lower semi-continuous in the diagonal matrix  $\mathbf{W}$ , i.e. for any sequence  $\mathbf{W}^{(k)} \rightarrow \overline{\mathbf{W}}$  with  $W_{ii}^{(k)} \geq 0$  for all  $k$  and  $i$ , we have*

$$\sigma_i(\overline{\mathbf{W}}\mathbf{A}) \leq \liminf_{k \rightarrow \infty} \sigma_i(\mathbf{W}^{(k)}\mathbf{A}).$$

With Lemmas 3.5.2 and 3.5.3 in place, we are ready to prove the main reweighting theorem.

*Proof of Theorem 3.5.1.* We prove the existence of the required  $\mathbf{W}$  by considering the limit of the Algorithm 1 for computing a reweighting matrix.

For all  $k \geq 0$ , let  $\mathbf{W}^{(k)}$  be the value of  $\mathbf{W}$  after the  $k^{\text{th}}$  update to the weight. We show that  $\overline{\mathbf{W}} = \lim_{k \rightarrow \infty} \mathbf{W}^{(k)}$  meets the conditions of Theorem 3.5.1. First note that Algorithm 1 is well defined and that all entries of  $\mathbf{W}^{(k)}$  are non-negative for all  $k \geq 0$ . To see this, suppose we need to decrease  $W_{ii}^{(k)}$  so that  $\sigma_i(\mathbf{W}^{(k+1)}\mathbf{A}) = \mathbf{u}_i$ . Note that the condition  $\sigma_i(\mathbf{W}^{(k)}\mathbf{A}) < 1$  gives

$$\lim_{\gamma \rightarrow 1} \frac{(1 - \gamma) \sigma_i(\mathbf{W}^{(k)}\mathbf{A})}{1 - \gamma \sigma_i(\mathbf{W}^{(k)}\mathbf{A})} = 0.$$

Therefore, Lemma 3.5.2 shows that we can make  $\sigma_i(\mathbf{W}^{(k+1)}\mathbf{A})$  arbitrary small by setting  $\gamma$  close enough to 1. Since the leverage score for row  $i$  is continuous, this implies that  $\mathbf{W}^{(k+1)}$  exists as desired.

---

**Algorithm 1:** ComputeReweighting (a.k.a the whack-a-mole algorithm)
 

---

```

Initialize  $\mathbf{W} = \mathbf{I}$ .
while true do
  for  $i = 1$  to  $n$  do
    if  $\sigma_i(\mathbf{W}\mathbf{A}) \geq \mathbf{u}_i$  then
      if  $\sigma_i(\mathbf{W}\mathbf{A}) < 1$  then
        | Decrease  $\mathbf{W}_{ii}$  so that  $\sigma_i(\mathbf{W}\mathbf{A}) = \mathbf{u}_i$ .
      else
        | Set  $\mathbf{W}_{ii} = 0$ .
      end
    end
  end
end
Output:  $\mathbf{W}$ .

```

---

Since, the entries of  $\mathbf{W}^{(k)}$  are non-negative and decrease monotonically by construction, clearly  $\overline{\mathbf{W}}$  exists. Furthermore, since setting  $\mathbf{W}_{ii} = 0$  makes  $\sigma_i(\mathbf{W}\mathbf{A}) = 0$ , we see that, by construction,

$$\liminf_{k \rightarrow \infty} \sigma_i(\mathbf{W}^{(k)}\mathbf{A}) \leq \mathbf{u}_i \text{ for all } i \in [n].$$

Therefore, by Lemma 3.5.3 we have that  $\sigma_i(\overline{\mathbf{W}}\mathbf{A}) \leq \mathbf{u}_i$ .

It only remains to show that  $\sum_{i: \overline{\mathbf{W}}_{ii} \neq 1} \mathbf{u}_i \leq n$ . Let  $k$  be the first iteration such that  $\mathbf{W}_{ii}^{(k)} \neq 1$  for any  $i$  such that  $\overline{\mathbf{W}}_{ii} \neq 1$ . Let  $S \subseteq [m]$  be the set of rows such that  $\mathbf{W}_{ii}^{(k)} = 0$  and let  $T = \{i : \overline{\mathbf{W}}_{ii} \neq 1\} - S$ . Since decreasing the weight of one row increases the leverage scores of all other rows, we have

$$\begin{aligned} \sum_{i \in T \cup S} \mathbf{u}_i &\leq \sum_{i \in T} \sigma_i(\mathbf{W}^{(k)}\mathbf{A}) + \sum_{i \in S} 1 \\ &\leq \text{rank}(\mathbf{W}^{(k)}\mathbf{A}) + |S|. \end{aligned}$$

When we set  $\mathbf{W}_{ii} = 0$ , it must be the case that  $\sigma_i(\mathbf{W}\mathbf{A}) = 1$ . In this case, removing the  $i^{\text{th}}$  row decreases the rank of  $\mathbf{W}\mathbf{A}$  by 1 and hence  $\text{rank}(\mathbf{W}^{(k)}\mathbf{A}) \leq n - |S|$ . Therefore,

$$\sum_{i: \overline{\mathbf{W}}_{ii} \neq 1} \mathbf{u}_i = \sum_{i \in T \cup S} \mathbf{u}_i \leq n.$$

□

## ■ 3.6 High Probability Bound on Leverage Score Estimation

Theorem 3.1.1 alone is enough to prove that a variety of iterative methods for spectral approximation work. In particular, we can use Lemma 2.3.3 to show that sampling each row with probability  $\Theta(\min\{1, \tilde{\sigma}_i \log n\})$  independently with suitable reweighting gives a spectral approximation of the original matrix. Since  $\tilde{\sigma}_i$  is a constant approximation of  $\min\{1, \sigma_i^{SA}(\mathbf{A})\}$ , the total number of rows it takes is  $\Theta(\sum_i \min\{1, \sigma_i^{SA}(\mathbf{A}) \log n\})$  which can be bounded by the  $\sum_i \min\{1, \sigma_i^{SA}(\mathbf{A})\}$  bound proved in Theorem 3.1.1.

In this section, we bound  $\sum_i \min\{1, \sigma_i^{SA}(\mathbf{A}) \log n\}$  directly and it allows us to bound the quantity with a high probability instead of merely in expectation. The proof relies on Theorem 3.5.1, which

intuitively shows that a large fraction of our matrix  $\mathbf{A}$  has low coherence. Sampling rows uniformly will give a spectral approximation for this portion of our matrix. Then, since few rows are reweighted in  $\mathbf{W}\mathbf{A}$ , even loose upper bounds on the leverage scores for those rows will allow us to bound the total sum of estimated leverage scores when we sample uniformly.

**Theorem 3.6.1.** *Given any  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Let  $S$  denote a uniformly random sample of  $k$  rows from  $\mathbf{A}$  and let  $\mathbf{S} \in \mathbb{R}^{m \times m}$  be its diagonal indicator matrix (i.e.  $\mathbf{S}_{ii} = 1$  for  $i \in S$ ,  $\mathbf{S}_{ii} = 0$  otherwise). With high probability in  $n$ , we have that*

$$\sum_{i=1}^n \min \{1, \sigma_i^{SA}(\mathbf{A}) \log n\} \leq O\left(\frac{mn \log n}{k}\right).$$

Resampling from  $\mathbf{A}$  by these estimates will, with high probability, return a  $\varepsilon$ -spectral approximation to  $\mathbf{A}$  with at most  $O\left(\frac{mn \log n}{k\varepsilon^2}\right)$  rows.

*Proof.* Let  $c$  be the constant  $c$  in Lemma 2.3.3. By Theorem 3.5.1, there is a diagonal matrix  $\mathbf{W} \preceq \mathbf{I}$  such that  $\sigma_i(\mathbf{W}\mathbf{A}) \leq \frac{k}{cm \log n}$  for all  $i$  and  $\sum_{i:\mathbf{W}_{ii} \neq 1} \frac{k}{cm \log n} \leq n$ . For this  $\mathbf{W}$ , we have

$$\begin{aligned} \sum_{i=1}^n \min \{1, \sigma_i^{SA}(\mathbf{A}) \log n\} &\leq \sum_{i:\mathbf{W}_{ii} \neq 1} 1 + \log n \sum_{i:\mathbf{W}_{ii}=1} \sigma_i^{SA}(\mathbf{A}) \\ &\leq \frac{cmn \log n}{k} + \log n \sum_{i:\mathbf{W}_{ii}=1} \sigma_i^{SA}(\mathbf{A}) \\ &= \frac{cmn \log n}{k} + \log n \sum_{i:\mathbf{W}_{ii}=1} \sigma_i^{SA}(\mathbf{W}\mathbf{A}). \end{aligned} \quad (3.5)$$

Since  $\sigma_i(\mathbf{W}\mathbf{A}) \leq \frac{k}{cm \log n}$  and  $S$  is a uniformly random sample of  $k$  rows, Lemma 2.3.3 shows that

$$\frac{1}{2} \mathbf{A}^\top \mathbf{W}^2 \mathbf{A} \preceq \frac{m}{k} \mathbf{A}^\top \mathbf{W} \mathbf{S}^2 \mathbf{W} \mathbf{A}. \quad (3.6)$$

Hence (3.6) along with Lemma 3.3.1 shows that

$$\sigma_i^{SA}(\mathbf{W}\mathbf{A}) \leq \frac{2m}{k} \sigma_i(\mathbf{W}\mathbf{A}).$$

Combining with (3.5), we have

$$\begin{aligned} \sum_{i=1}^n \min \{1, \sigma_i^{SA}(\mathbf{A}) \log n\} &\leq \frac{cmn \log n}{k} + \frac{2m \log n}{k} \sum_{i:\mathbf{W}_{ii}=1} \sigma_i(\mathbf{W}\mathbf{A}) \\ &\leq \frac{cmn \log n}{k} + \frac{2mn \log n}{k} = O\left(\frac{mn \log n}{k}\right). \end{aligned}$$

The last statement follows from Lemma 3.3.1 with sampling probabilities  $\min \left\{1, \frac{c\sigma_i^{SA}(\mathbf{A}) \log n}{\varepsilon^2}\right\}$ .  $\square$

## ■ 3.7 Main Result

As discussed in the introduction, Theorems 3.1.1 and 3.6.1 immediately yield new, extremely simple algorithms for spectral matrix approximation. For clarity, we initially present versions running in *nearly* input-sparsity time. However, we later explain how our first algorithm can be modified with

standard techniques to remove log factors, achieving input-sparsity time and thus matching state-of-the-art results [53, 185, 200]. Our algorithms rely solely on row sampling, which preserves matrix sparsity and structure, possibly improving space usage and runtime for intermediate system solves.

### ■ 3.7.1 Repeated Halving

The algorithm we present, `RepeatedHalving`, is a simple recursive procedure. We uniformly sample  $\frac{m}{2}$  rows from  $\mathbf{A}$  to obtain  $\mathbf{A}'$ . By Theorems 3.1.1 and 3.6.1, estimating leverage scores of  $\mathbf{A}$  with respect to this sample allows us to immediately find a spectral approximation to  $\mathbf{A}$  with  $O(n \log n)$  rows. Of course,  $\mathbf{A}'$  is still large, so computing these estimates would be slow. Thus, we *recursively* find a spectral approximation of  $\mathbf{A}'$  and use this to compute the estimated leverage scores.

---

#### Algorithm 2: RepeatedHalving

---

**Input:**  $m \times n$  matrix  $\mathbf{A}$

**Output:** spectral approximation  $\tilde{\mathbf{A}}$  consisting of  $O(n \log n)$  rescaled rows of  $\mathbf{A}$

Uniformly sample  $\frac{m}{2}$  rows of  $\mathbf{A}$  to form  $\mathbf{A}'$ .

If  $\mathbf{A}'$  has  $> O(n \log n)$  rows, **recursively** compute a spectral approximation  $\tilde{\mathbf{A}}'$  of  $\mathbf{A}'$ .

Compute approximate generalized leverage scores of  $\mathbf{A}$  w.r.t.  $\tilde{\mathbf{A}}'$

Use these estimates to sample rows of  $\mathbf{A}$  to form  $\tilde{\mathbf{A}}$ .

**Output**  $\tilde{\mathbf{A}}$ .

---

In analyzing the runtimes of these algorithms, we assume  $m = O(\text{poly}(n))$ , which is a reasonable assumption for any practical regression problem.<sup>2</sup> Furthermore, we use the fact that a  $n \times n$  system can be solved in time  $n^\omega$ , where  $\omega$  is the matrix multiplication exponent. However, we emphasize that, depending on the structure and sparsity of  $\mathbf{A}$ , alternative system solving methods may yield faster results or runtimes with different trade offs. For example, if the rows of  $\mathbf{A}$  are sparse, solving a system in  $\tilde{\mathbf{A}}$ , where  $\tilde{\mathbf{A}}$  consists of  $O(n \log n)$  rescaled rows from  $\mathbf{A}$  may be accelerated by using iterative conjugate gradient, or other Krylov subspace methods (which can also avoid explicitly computing  $\tilde{\mathbf{A}}^\top \tilde{\mathbf{A}}$ ). It is best to think of  $n^\omega$  as the runtime of the fastest available system solver in your domain, and the quoted runtimes as general guidelines that will change somewhat depending on exactly how the above algorithms are implemented.

First, we give an important primitive showing that estimates of generalized leverage scores can be computed efficiently. Computing exact generalized leverage scores is slow and we only need constant factor approximations, which will only increase our sampling rates and hence number of rows sampled by a constant factor.

**Lemma 3.7.1.** *Given  $\mathbf{B}$  containing  $O(n \log n)$  rescaled rows of  $\mathbf{A}$ , for any  $1 > \theta > 0$ , it is possible to compute an estimate of  $\boldsymbol{\sigma}^{\mathbf{B}}(\mathbf{A})$ ,  $\tilde{\boldsymbol{\sigma}}$ , in  $O(n^\omega \theta^{-1} + \text{nnz}(\mathbf{A})\theta^{-1})$  time such that, w.h.p. in  $n$ , for all  $i$ ,  $\tilde{\sigma}_i \geq \sigma_i^{\mathbf{B}}(\mathbf{A})$  and  $\tilde{\sigma}_i \leq O(n^\theta)\sigma_i^{\mathbf{B}}(\mathbf{A})$ .*

Setting  $\theta = O(\frac{1}{\log n})$  gives a constant factor approximation to generalized leverage score in  $O(n^\omega \log n + \text{nnz}(\mathbf{A}) \log n)$  time.

*Proof.* The proof is basically same as Lemma 2.3.4 except that we use fast matrix multiplication to solve the corresponding linear systems and that we need to handle the case  $\mathbf{a}_i \notin \ker(\mathbf{B})$ . When  $\mathbf{a}_i \notin \ker(\mathbf{B})$ , its generalized leverage score should be  $\infty$  – see (3.1). So, we need to check whether each  $\mathbf{a}_i$  has a component in the null-space of  $\mathbf{B}$ . This can be done in a variety of ways. For example, we can choose a random gaussian vector  $\mathbf{g}$  and compute  $\mathbf{g} - \mathbf{B}(\mathbf{B}^\top \mathbf{B})^+ \mathbf{B}^\top \mathbf{g}$ , which is the same as

---

<sup>2</sup>A simple method for handling even larger values of  $n$  is outlined in [170].



$\mathbf{g} - (\mathbf{B}^\top \mathbf{B})^+ \mathbf{B}^\top \mathbf{B} \mathbf{g}$ . This gives a random vector in the null space of  $\mathbf{B}$ , so computing its dot product with any row  $\mathbf{a}_i$  will tell us (with probability 1) whether  $\mathbf{a}_i$  is orthogonal to the null space or not.  $\square$

With Lemma 3.7.1 in place, we can analyze the runtime of our algorithms. For simplicity, we give runtimes for computing a constant factor spectral approximation to  $\mathbf{A}$ , which can be used as a preconditioner in iterative regression algorithms [22, 53, 228] or used to compute leverage scores of  $\mathbf{A}$  up to a constant factor. We could sample  $O(n \log n \varepsilon^{-2})$  rows to directly obtain a  $\varepsilon$ -approximation. By Lemma 3.7.1 the runtime of this final refinement is just  $O(\text{nnz}(\mathbf{A}) \log n + n^\omega \log n)$ .

**Lemma 3.7.2.** *RepeatedHalving (Algorithm 2) runs in time  $O(\text{nnz}(\mathbf{A}) \log n + n^\omega \log(m/n) \log n)$ , outputting a matrix with  $\tilde{\mathbf{A}}$  with  $O(n \log n)$  rows such that  $\tilde{\mathbf{A}}^\top \tilde{\mathbf{A}} \approx_2 \mathbf{A}^\top \mathbf{A}$ .*

*Proof.* The proof is by induction – it suffices to show that the work at the top level is  $O(\text{nnz}(\mathbf{A}) \log n + n^\omega \log n)$ . At each of the  $O(\log(m/n))$  levels of recursion, we cut our matrix in half uniformly so  $\text{nnz}(\mathbf{A})$  will also be cut approximately in half with high probability.

By Theorem 3.6.1, sampling by  $\sigma_i^{A'}(\mathbf{A})$  allows us to obtain  $\tilde{\mathbf{A}}$  with  $O(n \log n)$  rows. If we instead use  $\tilde{\mathbf{A}}'$ , our estimated leverage scores increase by at most a constant factor (since  $\tilde{\mathbf{A}}'$  is a constant factor spectral approximation to  $\mathbf{A}'$ ). Furthermore, using Lemma 3.7.1 to approximate generalized leverage scores increases our estimates by another constant factor at most. Overall,  $\tilde{\mathbf{A}}$  will have  $O(n \log n)$  rows as desired and our runtime at the top level is just the runtime of estimating leverage scores from Lemma 3.7.1 –  $O(\text{nnz}(\mathbf{A}) \log n + n^\omega \log n)$ .  $\square$

### ■ 3.7.2 Achieving Input Sparsity Time

We briefly note that, using techniques from [170], it is possible to remove the  $\log n$  factor on the  $\text{nnz}(\mathbf{A})$  term to achieve true input-sparsity time with **RepeatedHalving**. Instead of using Lemma 3.7.1 to estimate generalized leverage scores from up to a constant factor using  $\mathbf{A}'$ , we only estimate them up to a  $O(n^\theta)$  factor for some constant  $0 < \theta < \omega - 2$ . Using these rough estimates, we obtain  $\tilde{\mathbf{A}}$  with  $O(n^{1+\theta} \log n)$  rows. Then, for the rows in  $\tilde{\mathbf{A}}$ , we can again compute generalized leverage scores with respect to  $\mathbf{A}'$ , now up to constant factors, and reduce down to just  $O(n \log n)$  rows. In total, each iteration will take time  $O(\theta^{-1}(\text{nnz}(\mathbf{A}) + n^\omega \log n))$ , so obtaining a constant factor approximation to  $\mathbf{A}$  takes time  $O(\theta^{-1} \text{nnz}(\mathbf{A}) + \theta^{-1} n^\omega \log^2 n)$ . Recall that we assume  $m = \text{poly}(n)$ , so we have  $\log(m/n) = O(\log n)$  iterations.

In order to obtain a  $\varepsilon$ -spectral approximation with only  $O(n \log n \varepsilon^{-2})$  rows, we first obtain a constant factor approximation,  $\tilde{\mathbf{A}}$ , with  $O(n \log n)$  rows. We then use leverage scores estimated with  $\tilde{\mathbf{A}}$  to compute a  $\varepsilon/2$ -approximation to  $\mathbf{A}$  with  $O(n^{1+\theta} \log n \varepsilon^{-2})$  rows. Finally, we again use leverage scores estimated with  $\tilde{\mathbf{A}}$  and Lemma 3.7.1 with  $\theta = O(\varepsilon^2 / \log n)$  to compute a  $\varepsilon/2$ -approximation to this smaller matrix with only  $O(n \log n \varepsilon^{-2})$  rows. This takes total time  $O(\theta^{-1} \text{nnz}(\mathbf{A}) + \theta^{-1} n^\omega \log^2 n + n^{2+\theta} \log^2 n \varepsilon^{-2})$ . The  $n^{2+\theta} \log^3 n \varepsilon^{-2}$  comes from applying Lemma 3.7.1 to refine our second approximation, which has  $O(n^{1+\theta} \log n \varepsilon^{-2})$  rows and thus at most  $O(n^{2+\theta} \log n \varepsilon^{-2})$  nonzero entries. Overall, the technique yields:

**Lemma 3.7.3.** *Given any constant  $0 < \theta \leq \omega - 2$ , and any error  $0 \leq \varepsilon < 1$ , w.h.p. in  $n$  we can compute a matrix  $\tilde{\mathbf{A}}$  with  $O(n \log n \varepsilon^{-2})$  rows such that  $\tilde{\mathbf{A}}^\top \tilde{\mathbf{A}} \approx_{1+\varepsilon} \mathbf{A}^\top \mathbf{A}$  in  $O(\theta^{-1} \text{nnz}(\mathbf{A}) + \theta^{-1} n^\omega \log^2 n + n^{2+\theta} \varepsilon^{-2})$  time.*

*Proof.* As is standard,  $\log n$  factors on the  $n^{2+\theta}$  term are ‘hidden’ as we can just slightly increase the value of  $\theta$  to subsume them.  $\square$

## ■ 3.8 Appendix

### ■ 3.8.1 Rank 1 Updates

Here we prove Lemma 3.5.2, making critical use of the Sherman-Morrison formula for the Moore-Penrose pseudoinverse (Lemma 2.3.7).

*Lemma 3.5.2* (Leverage Score Changes Under Rank 1 Updates). Given any  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\gamma \in (0, 1)$ , and  $i \in [m]$ , let  $\mathbf{W}$  be a diagonal matrix such that  $\mathbf{W}_{ii} = \sqrt{1-\gamma}$  and  $\mathbf{W}_{jj} = 1$  for all  $j \neq i$ . Then, we have

$$\sigma_i(\mathbf{W}\mathbf{A}) = \frac{(1-\gamma)\sigma_i(\mathbf{A})}{1-\gamma\sigma_i(\mathbf{A})} \leq \sigma_i(\mathbf{A}),$$

and for all  $j \neq i$ ,

$$\sigma_j(\mathbf{W}\mathbf{A}) = \sigma_j(\mathbf{A}) + \frac{\gamma\sigma_{ij}(\mathbf{A})^2}{1-\gamma\sigma_i(\mathbf{A})} \geq \sigma_j(\mathbf{A}).$$

*Proof.*

$$\begin{aligned} \sigma_i(\mathbf{W}\mathbf{A}) &= \mathbf{1}_i \mathbf{W} \mathbf{A} \left( \mathbf{A}^\top \mathbf{W}^2 \mathbf{A} \right)^+ \mathbf{A}^\top \mathbf{W}^\top \mathbf{1}_i^\top && \text{(definition of leverage scores)} \\ &= (1-\gamma) \mathbf{a}_i^\top \left( \mathbf{A}^\top \mathbf{A} - \gamma \mathbf{a}_i \mathbf{a}_i^\top \right)^+ \mathbf{a}_i && \text{(definition of } \mathbf{W} \text{)} \\ &= (1-\gamma) \mathbf{a}_i^\top \left( (\mathbf{A}^\top \mathbf{A})^+ + \gamma \frac{(\mathbf{A}^\top \mathbf{A})^+ \mathbf{a}_i \mathbf{a}_i^\top (\mathbf{A}^\top \mathbf{A})^+}{1-\gamma \mathbf{a}_i^\top (\mathbf{A}^\top \mathbf{A})^+ \mathbf{a}_i} \right) \mathbf{a}_i && \text{(Sherman-Morrison formula)} \\ &= (1-\gamma) \left( \sigma_i(\mathbf{A}) + \frac{\gamma\sigma_i(\mathbf{A})^2}{1-\gamma\sigma_i(\mathbf{A})} \right) \\ &= \frac{(1-\gamma)\sigma_i(\mathbf{A})}{1-\gamma\sigma_i(\mathbf{A})} \leq \sigma_i(\mathbf{A}). \end{aligned}$$

Similarly,

$$\begin{aligned} \sigma_j(\mathbf{W}\mathbf{A}) &= \mathbf{a}_j^\top \left( \mathbf{A}^\top \mathbf{A} - \gamma \mathbf{a}_i \mathbf{a}_i^\top \right)^+ \mathbf{a}_j \\ &= \mathbf{a}_j^\top \left( (\mathbf{A}^\top \mathbf{A})^+ + \gamma \frac{(\mathbf{A}^\top \mathbf{A})^+ \mathbf{a}_i \mathbf{a}_i^\top (\mathbf{A}^\top \mathbf{A})^+}{1-\gamma \mathbf{a}_i^\top (\mathbf{A}^\top \mathbf{A})^+ \mathbf{a}_i} \right) \mathbf{a}_j \\ &= \sigma_j(\mathbf{A}) + \frac{\gamma\sigma_{ij}(\mathbf{A})^2}{1-\gamma\sigma_i(\mathbf{A})} \geq \sigma_j(\mathbf{A}). \end{aligned}$$

□

### ■ 3.8.2 Lower Semi-continuity of Leverage Scores

Here we prove Lemma 3.5.3 by providing a fairly general inequality, Lemma 3.8.1, for relating leverage scores under one set of weights to leverage scores under another.

*Lemma 3.5.3* (Leverage Scores are Lower Semi-continuous).  $\sigma(\mathbf{W}\mathbf{A})$  is lower semi-continuous in the diagonal matrix  $\mathbf{W}$ , i.e. for any sequence  $\mathbf{W}^{(k)} \rightarrow \overline{\mathbf{W}}$  with  $\mathbf{W}_{ii}^{(k)} \geq 0$  for all  $k$  and  $i$ , we have

$$\sigma_i(\overline{\mathbf{W}}\mathbf{A}) \leq \liminf_{k \rightarrow \infty} \sigma_i(\mathbf{W}^{(k)}\mathbf{A}). \quad (3.7)$$

**Lemma 3.8.1** (Comparing Leverage Scores). *Let  $\mathbf{W}, \overline{\mathbf{W}} \in \mathbb{R}^{m \times m}$  be non-negative diagonal matrices and suppose that  $\mathbf{W}_{ii} > 0$  and  $\overline{\mathbf{W}}_{ii} > 0$  for some  $i \in [m]$ . Then*

$$\sigma_i(\overline{\mathbf{W}}\mathbf{A}) \leq \frac{\overline{\mathbf{W}}_{ii}^2}{\mathbf{W}_{ii}^2} \left( 1 + \sqrt{\lambda_{\max} \left( \mathbf{A} \left( \mathbf{A}^\top \overline{\mathbf{W}}^2 \mathbf{A} \right)^+ \mathbf{A}^\top \right) \|\mathbf{W} - \overline{\mathbf{W}}\|_\infty} \right)^2 \sigma_i(\mathbf{W}\mathbf{A}). \quad (3.8)$$

*Proof.* Scaling the variables in Lemma 2.3.1 we have that there exists  $\mathbf{x} \in \mathbb{R}^m$  such that

$$\mathbf{A}^\top \mathbf{W}\mathbf{x} = \mathbf{a}_i \quad \text{and} \quad \|\mathbf{x}\|_2^2 = \frac{\sigma_i(\mathbf{W}\mathbf{A})}{\mathbf{W}_{ii}^2}. \quad (3.9)$$

Note that  $\mathbf{A}^\top (\mathbf{W} - \overline{\mathbf{W}})\mathbf{x}$  is in the image of  $\mathbf{A}^\top \overline{\mathbf{W}}$  as  $\mathbf{A}^\top \mathbf{W}\mathbf{x} = \mathbf{a}_i$  and  $\overline{\mathbf{W}}_{ii} \neq 0$ . Consequently,

$$\mathbf{A}^\top \overline{\mathbf{W}}\mathbf{y} = \mathbf{A}^\top (\mathbf{W} - \overline{\mathbf{W}})\mathbf{x}^{(k)} \quad \text{for} \quad \mathbf{y} \stackrel{\text{def}}{=} \overline{\mathbf{W}}\mathbf{A} \left( \mathbf{A}^\top \overline{\mathbf{W}}^2 \mathbf{A} \right)^+ \mathbf{A}^\top (\mathbf{W} - \overline{\mathbf{W}})\mathbf{x}.$$

Since  $\mathbf{A}^\top \overline{\mathbf{W}}(\mathbf{x} + \mathbf{y}) = \mathbf{a}_i$ , Lemma 2.3.1 implies

$$\sigma_i(\overline{\mathbf{W}}\mathbf{A}) \leq \overline{\mathbf{W}}_{ii}^2 \|\mathbf{x} + \mathbf{y}\|_2^2. \quad (3.10)$$

We can bound the contribution of  $\mathbf{y}$  by

$$\begin{aligned} \|\mathbf{y}\|_2^2 &\leq \left\| \overline{\mathbf{W}}\mathbf{A} \left( \mathbf{A}^\top \overline{\mathbf{W}}^2 \mathbf{A} \right)^+ \mathbf{A}^\top (\mathbf{W} - \overline{\mathbf{W}})\mathbf{x} \right\|_2^2 \\ &\leq \lambda_{\max} \left( \mathbf{A} \left( \mathbf{A}^\top \overline{\mathbf{W}}^2 \mathbf{A} \right)^+ \mathbf{A}^\top \right) \|\mathbf{W} - \overline{\mathbf{W}}\|_\infty^2 \|\mathbf{x}\|_2^2. \end{aligned} \quad (3.11)$$

Applying triangle inequality to (3.9), (3.10), and (3.11) yields the result.  $\square$

*Proof of Lemma 3.5.3.* For any  $i \in [n]$  such that  $\overline{\mathbf{W}}_{ii} = 0$  (3.7) follows trivially from the fact that leverage scores are non-negative. For any  $i \in [m]$  such that  $\overline{\mathbf{W}}_{ii} > 0$ , since  $\mathbf{W}^{(k)} \rightarrow \overline{\mathbf{W}}$  we know that, for all sufficiently large  $k \geq N$  for some fixed value  $N$ , it is the case that  $\mathbf{W}_{ii}^{(k)} > 0$ . Furthermore, this implies that as  $k \rightarrow \infty$  we have  $\overline{\mathbf{W}}_{ii}^2 / (\mathbf{W}_{ii}^{(k)})^2 \rightarrow 1$  and  $\|\mathbf{W}^{(k)} - \overline{\mathbf{W}}\|_\infty \rightarrow 0$ . Applying Lemma 3.8.1 with  $\mathbf{W} = \mathbf{W}^{(k)}$  and taking  $\liminf_{k \rightarrow \infty}$  on both sides of (3.8) gives the result.  $\square$



## Linear-Sized Sparsifier In Almost-Linear Time

### ■ 4.1 Introduction

In the last chapter, we showed how to compute  $O(n \log n)$  sized spectral sparsifier in input sparsity time. In this chapter, we show how to construct a spectral sparsifier with even smaller size. In particular, we present an almost-linear time algorithm for constructing linear-sized spectral sparsifiers for graphs, which improves all previous constructions that either require  $\Omega(n^{2+\varepsilon})$  time in order to produce linear-sized sparsifiers [9, 31, 278], or  $O(m \log^{O(1)} n/\varepsilon^2)$  time but the number of edges in the sparsifiers is sub-optimal [242]. Our algorithm is conceptually simple, and is based on a novel combination of two techniques used in literature for constructing spectral sparsifiers: random sampling by leverage scores we used in the last chapter, and adaptive construction based on barrier functions [9, 31]. Our result is summarized as follows:

**Theorem 4.1.1.** *Let  $q \geq 1$  and  $0 < \varepsilon \leq 1$  be constants, and  $\mathbf{I} = \sum_{i=1}^m \mathbf{v}_i \mathbf{v}_i^T$  be the sum of  $m$  rank-1 PSD matrices. Then, there is an algorithm that outputs scalars  $\{s_i\}_{i=1}^m$  with  $|\{s_i : s_i \neq 0\}| = O(qn/\varepsilon^2)$  such that*

$$(1 - \varepsilon) \cdot \mathbf{I} \preceq \sum_{i=1}^m s_i \mathbf{v}_i \mathbf{v}_i^T \preceq (1 + \varepsilon) \cdot \mathbf{I}.$$

*The algorithm runs in  $\tilde{O}(Z + n^{\omega+1/q} \varepsilon^{-4})$  time, where  $Z$  is the total number of non-zeros in  $\mathbf{v}_i$ .*

**Theorem 4.1.2.** *Let  $q \geq 1$  and  $0 < \varepsilon \leq 1$  be constants, and  $G = (V, E, \mathbf{w})$  be an undirected and weighted graph with  $n$  vertices and  $m$  edges. Then, there is an algorithm that outputs a  $\varepsilon$ -spectral sparsifier of  $G$  with  $O(qn/\varepsilon^2)$  edges. The algorithm runs in  $\tilde{O}(m + n^{1+1/q} \varepsilon^{-(6+1/q)})$  time.*

### ■ 4.2 Algorithm

In this section we study the algorithm of sparsifying the sum of rank-1 PSD matrices. Remember that our goal is to, for any vectors  $\mathbf{v}_1, \dots, \mathbf{v}_m$  with  $\sum_{i=1}^m \mathbf{v}_i \mathbf{v}_i^T = \mathbf{I}$ , find scalars  $\{s_i\}_{i=1}^m$  satisfying

$$|\{s_i : s_i \neq 0\}| = O\left(\frac{qn}{\varepsilon^2}\right),$$

such that

$$(1 - \varepsilon) \cdot \mathbf{I} \preceq \sum_{i=1}^m s_i \mathbf{v}_i \mathbf{v}_i^T \preceq (1 + \varepsilon) \cdot \mathbf{I}.$$

#### ■ 4.2.1 Overview of our approach

At a high level, our algorithm can be viewed as an improved and randomized version of the algorithm presented in Batson et al. [31]. We refer their algorithm BSS for short, and first give a brief overview of the BSS algorithm.

The BSS algorithm proceeds by iterations, and the matrix  $\mathbf{A}_j$  in iteration  $j \geq 1$  is defined as the sum of  $\mathbf{A}_{j-1}$  and a proper chosen rank-1 matrix. To control the spectral properties of matrix  $\mathbf{A}_j$ , the algorithm introduces two barrier values  $u_j$  and  $\ell_j$ , where  $u_0 > 0, \ell_0 < 0$  initially. It was proven that one can always find a vector in  $\{\mathbf{v}_i\}_{i=1}^m$  and update  $u_j, \ell_j$  in a proper manner in each iteration, such that the invariant

$$\ell_j \mathbf{I} \prec \mathbf{A}_j \prec u_j \mathbf{I} \quad (4.1)$$

always holds [31]. To quantitatively measure “how close the eigenvalues of  $\mathbf{A}$  are to the barriers  $u$  and  $\ell$ ”, Batson et al. [31] analyzed the potential function defined by

$$\Phi_{u,\ell}(\mathbf{A}) \triangleq \text{Tr}(u\mathbf{I} - \mathbf{A})^{-1} + \text{Tr}(\mathbf{A} - \ell\mathbf{I})^{-1}. \quad (4.2)$$

Notice that the value of  $\Phi_{u,\ell}(\mathbf{A})$  is large if and only if some eigenvalue of  $\mathbf{A}$  is close to  $u$  or  $\ell$ . With the help of this potential function, it was shown that, when updating  $\mathbf{A}_j$  and barrier values  $u_j, \ell_j$  properly, it holds after  $k = \Theta(n/\varepsilon^2)$  iterations that  $\ell_k \geq cu_k$  for some constant  $c$ . This implies that the resulting matrix  $\mathbf{A}_k$  is a linear-sized and  $\mathbf{A}_k \approx_{O(\varepsilon)} \mathbf{I}$ .

However, the BSS algorithm is deterministic, and its time complexity for finding a desired rank-1 matrix in each iteration is high. To improve the runtime and further discuss our approach, let us study the following randomized variant of the BSS algorithm: In each iteration, we choose a vector  $\mathbf{v}_i$  with probability  $p_i$ , and add a rank-1 matrix

$$\Delta_{\mathbf{A}} \triangleq \frac{\varepsilon}{\text{Tr}(u\mathbf{I} - \mathbf{A})^{-1} + \text{Tr}(\mathbf{A} - \ell\mathbf{I})^{-1}} \cdot \frac{1}{p_i} \cdot \mathbf{v}_i \mathbf{v}_i^T$$

to the current matrix  $\mathbf{A}$ . See Algorithm 3 for formal description.

---

**Algorithm 3:** Randomized BSS algorithm

---

Set  $j = 0, \ell_0 = -8n/\varepsilon, u_0 = 8n/\varepsilon, \mathbf{A}_0 = \mathbf{0}$ .

**while**  $u_j - \ell_j < 8n/\varepsilon$  **do**

Let  $t = \text{Tr}(u_j \mathbf{I} - \mathbf{A}_j)^{-1} + \text{Tr}(\mathbf{A}_j - \ell_j \mathbf{I})^{-1}$ .

Sample a vector  $\mathbf{v}_i$  with probability

$$p_i \triangleq \left( \mathbf{v}_i^T (u_j \mathbf{I} - \mathbf{A}_j)^{-1} \mathbf{v}_i + \mathbf{v}_i^T (\mathbf{A}_j - \ell_j \mathbf{I})^{-1} \mathbf{v}_i \right) / t.$$

$$\mathbf{A}_{j+1} = \mathbf{A}_j + \frac{\varepsilon}{t} \cdot \frac{1}{p_i} \cdot \mathbf{v}_i \mathbf{v}_i^T;$$

$$u_{j+1} = u_j + \frac{\varepsilon}{t \cdot (1-\varepsilon)} \text{ and } \ell_{j+1} = \ell_j + \frac{\varepsilon}{t \cdot (1+\varepsilon)};$$

$j \leftarrow j + 1$ ;

**end**

**Output:**  $\mathbf{A}_j$ .

---

Let us look at a fixed iteration  $j$ , and analyze how the added  $\Delta_{\mathbf{A}}$  impacts the potential function. We drop the subscript representing the iteration  $j$  for simplicity. After adding  $\Delta_{\mathbf{A}}$ , the first-order approximation of  $\Phi_{u,\ell}(\mathbf{A})$  gives that

$$\Phi_{u,\ell}(\mathbf{A} + \Delta_{\mathbf{A}}) \sim \Phi_{u,\ell}(\mathbf{A}) + (u\mathbf{I} - \mathbf{A})^{-2} \bullet \Delta_{\mathbf{A}} - (\mathbf{A} - \ell\mathbf{I})^{-2} \bullet \Delta_{\mathbf{A}}. \quad (4.3)$$

Since

$$\mathbf{E}[\Delta_{\mathbf{A}}] = \sum_{i=1}^m p_i \cdot \left( \frac{\varepsilon}{t} \cdot \frac{1}{p_i} \cdot \mathbf{v}_i \mathbf{v}_i^T \right) = \frac{\varepsilon}{t} \cdot \sum_{i=1}^m \mathbf{v}_i \mathbf{v}_i^T = \frac{\varepsilon}{t} \cdot \mathbf{I},$$

we have that

$$\begin{aligned} \mathbf{E} [\Phi_{u,\ell}(\mathbf{A} + \Delta_{\mathbf{A}})] &\sim \Phi_{u,\ell}(\mathbf{A}) + \frac{\varepsilon}{t} \cdot (u\mathbf{I} - \mathbf{A})^{-2} \bullet I - \frac{\varepsilon}{t} \cdot (\mathbf{A} - \ell\mathbf{I})^{-2} \bullet I \\ &= \Phi_{u,\ell}(\mathbf{A}) + \frac{\varepsilon}{t} \cdot \text{Tr}(u\mathbf{I} - \mathbf{A})^{-2} - \frac{\varepsilon}{t} \cdot \text{Tr}(\mathbf{A} - \ell\mathbf{I})^{-2} \\ &= \Phi_{u,\ell}(\mathbf{A}) - \frac{\varepsilon}{t} \cdot \frac{d}{du} \Phi_{u,\ell}(\mathbf{A}) - \frac{\varepsilon}{t} \cdot \frac{d}{d\ell} \Phi_{u,\ell}(\mathbf{A}). \end{aligned}$$

Notice that if we increase  $u$  by  $\varepsilon/t$  and  $\ell$  by  $\varepsilon/t$ ,  $\Phi_{u,\ell}$  approximately increases by

$$\frac{\varepsilon}{t} \cdot \frac{d}{du} \Phi_{u,\ell}(\mathbf{A}) + \frac{\varepsilon}{t} \cdot \frac{d}{d\ell} \Phi_{u,\ell}(\mathbf{A}).$$

Hence, comparing  $\Phi_{u+\varepsilon/t, \ell+\varepsilon/t}(\mathbf{A} + \Delta_{\mathbf{A}})$  with  $\Phi_{u,\ell}(\mathbf{A})$ , the increase of the potential function due to the change of barrier values is approximately compensated by the decrease of the potential function by the effect of  $\Delta_{\mathbf{A}}$ . For a more rigorous analysis, we need to analyze the higher-order terms and increase  $u$  slightly more than  $\ell$ . Batson et al. [31] gives the following estimate:

**Lemma 4.2.1** ([31], proof of Lemma 3.3 and 3.4). *Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , and  $u, \ell$  be parameters satisfying  $\ell\mathbf{I} \prec \mathbf{A} \prec u\mathbf{I}$ . Suppose that  $\mathbf{w} \in \mathbb{R}^n$  satisfies  $\mathbf{w}\mathbf{w}^T \preceq \delta(u\mathbf{I} - \mathbf{A})$  and  $\mathbf{w}\mathbf{w}^T \preceq \delta(\mathbf{A} - \ell\mathbf{I})$  for some  $0 < \delta < 1$ . Then, it holds that*

$$\Phi_{u,\ell}(\mathbf{A} + \mathbf{w}\mathbf{w}^T) \leq \Phi_{u,\ell}(\mathbf{A}) + \frac{\mathbf{w}^T(u\mathbf{I} - \mathbf{A})^{-2}\mathbf{w}}{1 - \delta} - \frac{\mathbf{w}^T(\mathbf{A} - \ell\mathbf{I})^{-2}\mathbf{w}}{1 + \delta}.$$

The estimate above shows that the first-order approximation (4.3) is good as long as  $\mathbf{w}\mathbf{w}^T \preceq \delta(u\mathbf{I} - \mathbf{A})$  and  $\mathbf{w}\mathbf{w}^T \preceq \delta(\mathbf{A} - \ell\mathbf{I})$  for small  $\delta$ . By setting  $\delta = \varepsilon$ , it is easy to see that the added matrix  $\Delta_{\mathbf{A}}$  satisfies the preconditions in Lemma 4.2.1, since

$$\frac{\varepsilon}{t} \cdot \frac{1}{p_i} \cdot \mathbf{v}_i \mathbf{v}_i^T = \frac{\varepsilon \cdot \mathbf{v}_i \mathbf{v}_i^T}{\mathbf{v}_i^T (u\mathbf{I} - \mathbf{A})^{-1} \mathbf{v}_i + \mathbf{v}_i^T (\mathbf{A} - \ell\mathbf{I})^{-1} \mathbf{v}_i} \preceq \frac{\varepsilon \cdot \mathbf{v}_i \mathbf{v}_i^T}{\mathbf{v}_i^T (u\mathbf{I} - \mathbf{A})^{-1} \mathbf{v}_i} \preceq \varepsilon (u\mathbf{I} - \mathbf{A}).$$

Here we used the fact that  $\mathbf{v}\mathbf{v}^T \preceq (\mathbf{v}^T B^{-1} \mathbf{v}) B$  for any vector  $\mathbf{v}$  and PSD matrix  $B$ . Similarly, we have that

$$\frac{\varepsilon}{t} \cdot \frac{1}{p_i} \cdot \mathbf{v}_i \mathbf{v}_i^T \preceq \varepsilon (\mathbf{A} - \ell\mathbf{I}).$$

Hence, if the initial value of the potential function is small,  $\Phi_{u,\ell}(\mathbf{A})$  is small throughout the executions of the whole algorithm. Up to a constant factor, this gives the same result as [31], and Algorithm 3 constructs an  $\Theta(n/\varepsilon^2)$ -sized  $O(\varepsilon)$ -spectral sparsifier.

However, Algorithm 3 runs for  $\Theta(n/\varepsilon^2)$  iterations, and in each iteration the algorithm re-computes the probability distribution for sampling vectors. To improve the runtime of Algorithm 3, we need to overcome two bottlenecks: (1) we need a fast algorithm to approximate the sampling probability  $\{p_i\}_{i=1}^m$  of vectors; (2) we need to reduce the number of iterations required by the algorithm, i.e. re-using the sampling probability  $\{p_i\}_{i=1}^m$  for few iterations.

For fast approximation of the sampling probabilities, we adopt the idea proposed in [9]: instead of defining the potential function by (4.2), we define the potential function by

$$\Phi_{u,\ell}(\mathbf{A}) \triangleq \text{Tr}(u\mathbf{I} - \mathbf{A})^{-q} + \text{Tr}(\mathbf{A} - \ell\mathbf{I})^{-q}. \quad (4.4)$$

Since  $q$  is a large constant, the value of the potential function becomes larger when some eigenvalue of  $\mathbf{A}$  is closer to  $u$  or  $\ell$ . Hence, a bounded value of  $\Phi_{u,\ell}(\mathbf{A})$  insures that the eigenvalues of  $\mathbf{A}$  never get too close to  $u$  or  $\ell$ , which allows us to compute the sampling probabilities  $\{p_i\}_{i=1}^m$  efficiently simply by Taylor expansion. Moreover, with our new potential function (4.4) one can prove a similar result

as Lemma 4.2.1. This gives an alternative analysis of the algorithm presented in [9], which is the first almost-quadratic time algorithm for constructing linear-sized spectral sparsifiers.

To overcome the second bottleneck, we re-compute  $\{p_i\}_{i=1}^m$  after every  $\Theta(n^{1-1/q})$  iterations: we show that as long as the sampling probability satisfies

$$p_i \geq C \cdot \frac{\mathbf{v}_i^T (\mathbf{u}\mathbf{I} - \mathbf{A})^{-1} \mathbf{v}_i + \mathbf{v}_i^T (\mathbf{A} - \ell\mathbf{I})^{-1} \mathbf{v}_i}{\sum_{i=1}^m \left( \mathbf{v}_i^T (\mathbf{u}\mathbf{I} - \mathbf{A})^{-1} \mathbf{v}_i + \mathbf{v}_i^T (\mathbf{A} - \ell\mathbf{I})^{-1} \mathbf{v}_i \right)}$$

for some constant  $C > 0$ , we can still sample  $v_i$  with probability  $p_i$  and get the same guarantee on the potential function. The reason is as follows: assume that  $\Delta_{\mathbf{A}} = \sum_{i=1}^T \Delta_{\mathbf{A},i}$  is the sum of the sampled matrices within  $T = O(n^{1-1/q})$  iterations. If a randomly chosen matrix  $\Delta_{\mathbf{A},i}$  satisfies  $\Delta_{\mathbf{A},i} \preceq \frac{1}{Cq} (\mathbf{u}\mathbf{I} - \mathbf{A})$ , then by the matrix Chernoff bound  $\Delta_{\mathbf{A}} \preceq \frac{1}{2} (\mathbf{u}\mathbf{I} - \mathbf{A})$  holds with high probability. By scaling every sampled rank-1 matrix  $q$  times smaller, the sampling probability only changes by a constant factor within  $T$  iterations. Since we choose  $\Theta(n/\varepsilon^2)$  vectors in total, our algorithm only recomputes the sampling probabilities  $\Theta(n^{1/q}/\varepsilon^2)$  times.

## ■ 4.2.2 Algorithm description

Our actual algorithm follows the same framework as Algorithm 3, and proceeds by iterations. Initially, the algorithm sets

$$u_0 \triangleq (2n)^{1/q}, \quad \ell_0 \triangleq -(2n)^{1/q}, \quad A_0 \triangleq \mathbf{0}.$$

After iteration  $j$  the algorithm updates  $u_j, \ell_j$  by  $\Delta_{u,j}, \Delta_{\ell,j}$  respectively, i.e.,

$$u_{j+1} \triangleq u_j + \Delta_{u,j}, \quad \ell_{j+1} \triangleq \ell_j + \Delta_{\ell,j},$$

and updates  $A_j$  with respect to the chosen matrix in iteration  $j$ . The choice of  $\Delta_{u,j}$  and  $\Delta_{\ell,j}$  insures that the invariant

$$\ell_j \mathbf{I} \prec \mathbf{A}_j \prec u_j \mathbf{I}$$

always holds for any iteration  $j$ . In iteration  $j$ , the algorithm computes the *relative leverage score* of vectors  $\{v_i\}_{i=1}^m$  defined by

$$R_i(\mathbf{A}_j, u_j, \ell_j) \triangleq \mathbf{v}_i^T (\mathbf{u}_j \mathbf{I} - \mathbf{A}_j)^{-1} \mathbf{v}_i + \mathbf{v}_i^T (\mathbf{A}_j - \ell_j \mathbf{I})^{-1} \mathbf{v}_i,$$

and samples  $N_j$  vectors independently with replacement, where vector  $v_i$  is chosen with probability proportional to  $R_i(\mathbf{A}_j, u_j, \ell_j)$ , and

$$N_j \triangleq \frac{1}{n^{2/q}} \left( \sum_{i=1}^m R_i(\mathbf{A}_j, u_j, \ell_j) \right) \min \{ \lambda_{\min}(\mathbf{u}_j \mathbf{I} - \mathbf{A}_j), \lambda_{\min}(\mathbf{A}_j - \ell_j \mathbf{I}) \}.$$

The algorithm sets  $\mathbf{A}_{j+1}$  to be the sum of  $\mathbf{A}_j$  and sampled  $\mathbf{v}_i \mathbf{v}_i^T$  with proper reweighting. For technical reasons, we define  $\Delta_{u,j}$  and  $\Delta_{\ell,j}$  by

$$\Delta_{u,j} \triangleq (1 + 2\varepsilon) \cdot \frac{\varepsilon \cdot N_j}{q \cdot \sum_{i=1}^m R_i(\mathbf{A}_j, u_j, \ell_j)}, \quad \Delta_{\ell,j} \triangleq (1 - 2\varepsilon) \cdot \frac{\varepsilon \cdot N_j}{q \cdot \sum_{i=1}^m R_i(\mathbf{A}_j, u_j, \ell_j)}.$$

See Algorithm 4 for formal description.

We remark that, although exact values of  $N_j$  and relative leverage scores are difficult to compute in almost-linear time, we can use approximated values of  $\{R_i\}_{i=1}^m$  and  $N_j$  instead. It is easy to see that in each iteration an over estimate of the relative leverage score for every vector  $v_i$ , and an under estimate of  $N_j$  with constant-factor approximation suffice for our purpose.



**Algorithm 4:** Algorithm for constructing spectral sparsifiers**Assume:**  $0 < \varepsilon \leq 1/120, q \geq 10$ .Set  $j = 0, \ell_0 = -(2n)^{1/q}, u_0 = (2n)^{1/q}, \mathbf{A}_0 = \mathbf{0}$ .**while**  $u_j - \ell_j < 4 \cdot (2n)^{1/q}$  **do**    Let  $\mathbf{W}_j = \mathbf{0}$ .    Compute  $R_i(\mathbf{A}_j, u_j, \ell_j)$  for all vectors  $v_i$ .    Sample  $N_j$  vectors independently with replacement, where every  $v_i$  is chosen with probability proportional to  $R_i(\mathbf{A}_j, u_j, \ell_j)$ .    For every sampled  $v$ , add  $\varepsilon/q \cdot (R_i(\mathbf{A}_j, u_j, \ell_j))^{-1} \cdot vv^T$  to  $\mathbf{W}_j$ .     $\mathbf{A}_{j+1} = \mathbf{A}_j + \mathbf{W}_j$ .     $u_{j+1} = u_j + \Delta_{u,j}, \ell_{j+1} = \ell_j + \Delta_{\ell,j}$ .     $j = j + 1$ .**end****Output:**  $\mathbf{A}_j$ .**■ 4.3 Analysis**

In this section we prove that the output matrix returned by Algorithm 4 is a  $\varepsilon$ -spectral sparsifier, and analyze the algorithm's runtime. To simplify our analysis, we always assume in the rest of the chapter that  $\varepsilon$  and  $q$  are constants such that  $0 < \varepsilon \leq 1/120$ , and  $q \geq 10$ .

This section is organized as follows: we will first show how the potential function (4.4) evolves after each iteration in Section 4.3.1. Combing this with the ending condition of the algorithm, we will prove in Section 4.3.2 that the algorithm outputs a linear-sized spectral sparsifier. Finally, we will prove Theorem 4.1.2 and Theorem 4.1.1 in Section 4.3.3. In appendix, we show almost-linear time algorithms for approximately computing all required quantities of Algorithm 4 within each iteration when constructing graph sparsifiers.

**■ 4.3.1 Analysis of a single iteration**

We analyze the sampling scheme within a single iteration, and drop the subscript representing the iteration  $j$  for simplicity. Recall that in each iteration the algorithm samples  $N$  vectors independently from  $\{\mathbf{v}_i\}_{i=1}^m$  satisfying  $\sum_{i=1}^m \mathbf{v}_i \mathbf{v}_i^T = \mathbf{I}$ , where every vector  $\mathbf{v}_i$  is sampled with probability

$$\frac{R_i(\mathbf{A}, u, \ell)}{\sum_{j=1}^m R_j(\mathbf{A}, u, \ell)}.$$

We use  $\mathbf{v}_1, \dots, \mathbf{v}_N$  to denote these  $N$  sampled vectors, and define the reweighted vectors by

$$\mathbf{w}_i \triangleq \sqrt{\frac{\varepsilon}{q \cdot R_i(\mathbf{A}, u, \ell)}} \cdot \mathbf{v}_i,$$

for any  $1 \leq i \leq N$ . Let

$$\mathbf{W} \triangleq \sum_{i=1}^N \mathbf{w}_i \mathbf{w}_i^T,$$

and we use  $\mathbf{W} \sim \mathcal{D}(\mathbf{A}, u, \ell)$  to represent that  $\mathbf{W}$  is sampled in this way with parameters  $\mathbf{A}, u$  and  $\ell$ . We will show that with high probability matrix  $\mathbf{W}$  satisfies  $\mathbf{0} \preceq \mathbf{W} \preceq \frac{1}{2}(u\mathbf{I} - \mathbf{A})$ . This is basically follows from Matrix Chernoff Bound (Lemma 2.3.10).

**Lemma 4.3.1.** *Assume that the number of sampled vectors satisfies*

$$N < \frac{2}{n^{2/q}} \left( \sum_{i=1}^m R_i(\mathbf{A}, u, \ell) \right) \cdot \lambda_{\min}(u\mathbf{I} - \mathbf{A}).$$

Then, it holds that  $\mathbf{E}[\mathbf{W}] = \frac{\varepsilon}{q} \cdot \frac{N}{\sum_{i=1}^m R_i(\mathbf{A}, u, \ell)} \cdot \mathbf{I}$  and

$$\Pr \left[ \mathbf{0} \preceq \mathbf{W} \preceq \frac{1}{2} \cdot (u\mathbf{I} - \mathbf{A}) \right] \geq 1 - \frac{\varepsilon^2}{100qn}.$$

*Proof.* By the description of the sampling procedure, the first statement follows from the calculation

$$\mathbf{E}[\mathbf{w}_i \mathbf{w}_i^T] = \sum_{j=1}^m \frac{R_j(\mathbf{A}, u, \ell)}{\sum_{t=1}^m R_t(\mathbf{A}, u, \ell)} \cdot \frac{\varepsilon}{q} \cdot \frac{\mathbf{v}_j \mathbf{v}_j^T}{R_j(\mathbf{A}, u, \ell)} = \frac{\varepsilon}{q} \cdot \frac{1}{\sum_{t=1}^m R_t(\mathbf{A}, u, \ell)} \cdot \mathbf{I},$$

and

$$\mathbf{E}[\mathbf{W}] = \mathbf{E} \left[ \sum_{i=1}^N \mathbf{w}_i \mathbf{w}_i^T \right] = \frac{\varepsilon}{q} \cdot \frac{N}{\sum_{i=1}^m R_i(\mathbf{A}, u, \ell)} \cdot \mathbf{I}.$$

Now for the second statement. Let  $\mathbf{z}_i = (u\mathbf{I} - \mathbf{A})^{-1/2} \mathbf{w}_i$ . It holds that

$$\begin{aligned} \text{Tr}(\mathbf{z}_i \mathbf{z}_i^T) &= \frac{\varepsilon}{q} \cdot \frac{\text{Tr}((u\mathbf{I} - \mathbf{A})^{-1/2} \mathbf{v}_i \mathbf{v}_i^T (u\mathbf{I} - \mathbf{A})^{-1/2})}{R_i(\mathbf{A}, u, \ell)} \\ &\leq \frac{\varepsilon}{q} \cdot \frac{\mathbf{v}_i^T (u\mathbf{I} - \mathbf{A})^{-1} \mathbf{v}_i}{\mathbf{v}_i^T (u\mathbf{I} - \mathbf{A})^{-1} \mathbf{v}_i + \mathbf{v}_i^T (\mathbf{A} - \ell \mathbf{I})^{-1} \mathbf{v}_i} \leq \frac{\varepsilon}{q}, \end{aligned}$$

and  $\lambda_{\max}(\mathbf{z}_i \mathbf{z}_i^T) \leq \frac{\varepsilon}{q}$ . Moreover, it holds that

$$\mathbf{E} \left[ \sum_{i=1}^N \mathbf{z}_i \mathbf{z}_i^T \right] \preceq \frac{\varepsilon}{q} \cdot \frac{N}{\sum_{t=1}^m R_t(\mathbf{A}, u, \ell)} \cdot \lambda_{\max} \left( \frac{1}{u\mathbf{I} - \mathbf{A}} \right) \cdot \mathbf{I}. \quad (4.5)$$

This implies that

$$\lambda_{\max} \left( \mathbf{E} \left[ \sum_{i=1}^N \mathbf{z}_i \mathbf{z}_i^T \right] \right) \leq \frac{\varepsilon}{q} \cdot \frac{N}{\sum_{t=1}^m R_t(\mathbf{A}, u, \ell)} \cdot \lambda_{\max} \left( \frac{1}{u\mathbf{I} - \mathbf{A}} \right).$$

Setting  $\mu = \frac{\varepsilon}{q} \cdot \frac{N}{\sum_{i=1}^m R_i(\mathbf{A}, u, \ell)} \cdot \lambda_{\max} \left( \frac{1}{u\mathbf{I} - \mathbf{A}} \right)$ , the Matrix Chernoff Bound (Lemma 2.3.10) shows that

$$\Pr \left[ \lambda_{\max} \left( \sum_{i=1}^N \mathbf{z}_i \mathbf{z}_i^T \right) \geq (1 + \delta) \mu \right] \leq n \cdot \left( \frac{e^\delta}{(1 + \delta)^{1 + \delta}} \right)^{\mu \cdot q / \varepsilon}.$$

Set the value of  $1 + \delta$  to be

$$\begin{aligned} 1 + \delta &= \frac{1}{2\mu} = \frac{q}{2\varepsilon N} \cdot \left( \sum_{j=1}^m R_j(\mathbf{A}, u, \ell) \right) \cdot \frac{1}{\lambda_{\max} \left( \frac{1}{u\mathbf{I} - \mathbf{A}} \right)} \\ &= \frac{q}{2\varepsilon N} \cdot \left( \sum_{j=1}^m R_j(\mathbf{A}, u, \ell) \right) \cdot \lambda_{\min}(u\mathbf{I} - \mathbf{A}) \geq \frac{q}{4\varepsilon} \cdot n^{2/q}, \end{aligned}$$

where the last inequality follows from the condition on  $N$ . Hence, with probability at least

$$1 - n \cdot \left( \frac{e^\delta}{(1+\delta)^{1+\delta}} \right)^{\mu \cdot q/\varepsilon} \geq 1 - n \cdot \left( \frac{e}{1+\delta} \right)^{(1+\delta) \cdot \mu \cdot q/\varepsilon} \geq 1 - n \left( \frac{e}{1+\delta} \right)^{\frac{q}{2\varepsilon}} \geq 1 - \frac{\varepsilon^2}{100qn},$$

we have that

$$\lambda_{\max} \left( \sum_{i=1}^N z_i z_i^T \right) \leq (1+\delta) \cdot \mu = \frac{1}{2},$$

which implies that  $\mathbf{0} \preceq \sum_{i=1}^N z_i z_i^T \preceq \frac{1}{2} \cdot I$  and  $\mathbf{0} \preceq W \preceq \frac{1}{2} \cdot (uI - A)$ .  $\square$

Now we analyze the change of the potential function after each iteration, and show that the expected value of the potential function decreases over time. By Lemma 4.3.1, with probability at least  $1 - \frac{\varepsilon^2}{100qn}$ , it holds that

$$\mathbf{0} \preceq \mathbf{W} \preceq \frac{1}{2}(u\mathbf{I} - \mathbf{A}).$$

We define

$$\tilde{\mathbf{E}}[f(\mathbf{W})] \triangleq \sum_{\mathbf{W} \sim \mathcal{D}(\mathbf{A}, u, \ell)} \Pr \left[ \mathbf{W} \text{ is chosen and } \mathbf{W} \preceq \frac{1}{2}(u\mathbf{I} - \mathbf{A}) \right] \cdot f(\mathbf{W}).$$

Lemma 4.3.2 below shows how the potential function changes after each iteration, and plays a key role in our analysis. This lemma was first proved in [31] for the case of  $q = 1$ , and similar lemma was proved in [9].

**Lemma 4.3.2** ([9]). *Let  $q \geq 10$  and  $\varepsilon \leq 1/10$ . Suppose that  $\mathbf{w}^T(u\mathbf{I} - \mathbf{A})^{-1}\mathbf{w} \leq \varepsilon/q$  and  $\mathbf{w}^T(\mathbf{A} - \ell\mathbf{I})^{-1}\mathbf{w} \leq \varepsilon/q$ . Then, it holds that*

$$\text{Tr}(\mathbf{A} + \mathbf{w}\mathbf{w}^T - \ell\mathbf{I})^{-q} \leq \text{Tr}(\mathbf{A} - \ell\mathbf{I})^{-q} - q(1 - \varepsilon) \mathbf{w}^T(\mathbf{A} - \ell\mathbf{I})^{-(q+1)}\mathbf{w},$$

and

$$\text{Tr}(u\mathbf{I} - \mathbf{A} - \mathbf{w}\mathbf{w}^T)^{-q} \leq \text{Tr}(u\mathbf{I} - \mathbf{A})^{-q} + q(1 + \varepsilon) \mathbf{w}^T(u\mathbf{I} - \mathbf{A})^{-(q+1)}\mathbf{w}.$$

*Proof.* Let  $\mathbf{Y} = \mathbf{A} - \ell\mathbf{I}$ . By the Sherman-Morrison Formula (Lemma 2.3.7), it holds that

$$\text{Tr}(\mathbf{Y} + \mathbf{w}\mathbf{w}^T)^{-q} = \text{Tr} \left( \mathbf{Y}^{-1} - \frac{\mathbf{Y}^{-1}\mathbf{w}\mathbf{w}^T\mathbf{Y}^{-1}}{1 + \mathbf{w}^T\mathbf{Y}^{-1}\mathbf{w}} \right)^q. \quad (4.6)$$

By the assumption of  $\mathbf{w}^T\mathbf{Y}^{-1}\mathbf{w} \leq \varepsilon/q$ , we have that

$$\text{Tr}(\mathbf{Y} + \mathbf{w}\mathbf{w}^T)^{-q} \leq \text{Tr} \left( \mathbf{Y}^{-1} - \frac{\mathbf{Y}^{-1}\mathbf{w}\mathbf{w}^T\mathbf{Y}^{-1}}{1 + \varepsilon/q} \right)^q \quad (4.7)$$

$$= \text{Tr} \left( \mathbf{Y}^{-1/2} \left( \mathbf{I} - \frac{\mathbf{Y}^{-1/2}\mathbf{w}\mathbf{w}^T\mathbf{Y}^{-1/2}}{1 + \varepsilon/q} \right) \mathbf{Y}^{-1/2} \right)^q$$

$$\leq \text{Tr} \left( \mathbf{Y}^{-q/2} \left( \mathbf{I} - \frac{\mathbf{Y}^{-1/2}\mathbf{w}\mathbf{w}^T\mathbf{Y}^{-1/2}}{1 + \varepsilon/q} \right)^q \mathbf{Y}^{-q/2} \right) \quad (4.8)$$

$$= \text{Tr} \left( \mathbf{Y}^{-q} \left( \mathbf{I} - \frac{\mathbf{Y}^{-1/2}\mathbf{w}\mathbf{w}^T\mathbf{Y}^{-1/2}}{1 + \varepsilon/q} \right)^q \right), \quad (4.9)$$

where (4.7) uses the fact that  $\mathbf{A} \preceq \mathbf{B}$  implies that  $\text{Tr}(\mathbf{A}^q) \leq \text{Tr}(\mathbf{B}^q)$ , (4.8) follows from the Lieb-Thirring inequality (Lemma 2.3.11), and (4.9) uses the fact that the trace is invariant under cyclic

permutations.

Let  $\mathbf{D} = \frac{\mathbf{Y}^{-1/2}ww^T\mathbf{Y}^{-1/2}}{1+\varepsilon/q}$ . Note that  $0 \preceq \mathbf{D} \preceq \frac{\varepsilon}{q} \cdot \mathbf{I}$ , and

$$\begin{aligned} (\mathbf{I} - \mathbf{D})^q &\preceq \mathbf{I} - q\mathbf{D} + \frac{q(q-1)}{2}\mathbf{D}^2 \\ &\preceq \mathbf{I} - \left(q - \frac{\varepsilon(q-1)}{2}\right)\mathbf{D}. \end{aligned}$$

Therefore, we have that

$$\begin{aligned} \left(\mathbf{I} - \frac{\mathbf{Y}^{-1/2}ww^T\mathbf{Y}^{-1/2}}{1+\varepsilon/q}\right)^q &\preceq \mathbf{I} - \left(q - \frac{\varepsilon(q-1)}{2}\right) \frac{\mathbf{Y}^{-1/2}ww^T\mathbf{Y}^{-1/2}}{1+\varepsilon/q} \\ &\preceq \mathbf{I} - \left(q - \frac{\varepsilon(q-1)}{2}\right) \left(1 - \frac{\varepsilon}{q}\right) \mathbf{Y}^{-1/2}ww^T\mathbf{Y}^{-1/2} \\ &\preceq \mathbf{I} - q \left(1 - \frac{\varepsilon(q+1)}{2q}\right) \mathbf{Y}^{-1/2}ww^T\mathbf{Y}^{-1/2} \\ &\preceq \mathbf{I} - q(1-\varepsilon) \mathbf{Y}^{-1/2}ww^T\mathbf{Y}^{-1/2}. \end{aligned}$$

This implies that

$$\begin{aligned} \text{Tr}(\mathbf{Y} + ww^T)^{-q} &\leq \text{Tr}\left(\mathbf{Y}^{-q} \left(\mathbf{I} - q(1-\varepsilon)\mathbf{Y}^{-1/2}ww^T\mathbf{Y}^{-1/2}\right)\right) \\ &\leq \text{Tr}(\mathbf{Y}^{-q}) - q(1-\varepsilon) w^T\mathbf{Y}^{-(q+1)}w, \end{aligned}$$

which proves the first statement.

Now for the second inequality. Let  $\mathbf{Z} = u\mathbf{I} - \mathbf{A}$ . By the Sherman-Morrison Formula (Lemma 2.3.7), it holds that

$$\text{Tr}(\mathbf{Z} - ww^T)^{-q} = \text{Tr}\left(\mathbf{Z}^{-1} + \frac{\mathbf{Z}^{-1}ww^T\mathbf{Z}^{-1}}{1 - w^T\mathbf{Z}^{-1}w}\right)^q.$$

By the assumption of  $w^T\mathbf{Z}^{-1}w \leq \varepsilon/q$ , it holds that

$$\text{Tr}(\mathbf{Z} - ww^T)^{-q} \leq \text{Tr}\left(\mathbf{Z}^{-1} + \frac{\mathbf{Z}^{-1}ww^T\mathbf{Z}^{-1}}{1 - \varepsilon/q}\right)^q \quad (4.10)$$

$$\begin{aligned} &= \text{Tr}\left(\mathbf{Z}^{-1/2} \left(\mathbf{I} + \frac{\mathbf{Z}^{-1/2}ww^T\mathbf{Z}^{-1/2}}{1 - \varepsilon/q}\right) \mathbf{Z}^{-1/2}\right)^q \\ &\leq \text{Tr}\left(\mathbf{Z}^{-q/2} \left(\mathbf{I} + \frac{\mathbf{Z}^{-1/2}ww^T\mathbf{Z}^{-1/2}}{1 - \varepsilon/q}\right)^q \mathbf{Z}^{-q/2}\right) \end{aligned} \quad (4.11)$$

$$= \text{Tr}\left(\mathbf{Z}^{-q} \left(\mathbf{I} + \frac{\mathbf{Z}^{-1/2}ww^T\mathbf{Z}^{-1/2}}{1 - \varepsilon/q}\right)^q\right), \quad (4.12)$$

where (4.10) uses the fact that  $\mathbf{A} \preceq \mathbf{B}$  implies  $\text{Tr}(\mathbf{A}^q) \leq \text{Tr}(\mathbf{B}^q)$ , (4.11) follows from the Lieb-Thirring inequality (Lemma 2.3.11), and (4.12) uses the fact that the trace is invariant under cyclic permutations.

Let  $\mathbf{E} = \mathbf{Z}^{-1/2}ww^T\mathbf{Z}^{-1/2}$ . Combing  $\mathbf{E} \preceq \frac{\varepsilon}{q} \cdot \mathbf{I}$  with the assumption that  $q \geq 10$  and  $\varepsilon \leq 1/10$ , we

have that

$$\begin{aligned}
\left(I + \frac{E}{1 - \varepsilon/q}\right)^q &\preceq I + \frac{qE}{1 - \varepsilon/q} + \frac{q(q-1)}{2} \left(1 + \frac{\varepsilon/q}{1 - \varepsilon/q}\right)^{q-2} \left(\frac{E}{1 - \varepsilon/q}\right)^2 \\
&\preceq I + q \left(1 + 1.1 \frac{\varepsilon}{q}\right) E + 1.4 \frac{q(q-1)}{2} E^2 \\
&\preceq I + q(1 + 0.3\varepsilon) E + 0.7\varepsilon q E \\
&\preceq I + q(1 + \varepsilon) E.
\end{aligned}$$

Therefore, we have that

$$\mathrm{Tr}(Z - ww^T)^{-q} \leq \mathrm{Tr}(Z^{-q}) + q(1 + \varepsilon) w^T Z^{-(q+1)} w,$$

which proves the second statement.  $\square$

The following lemma shows that, with our choice of updated matrices and barrier values, the value of the potential function will never increase.

**Lemma 4.3.3.** *It holds for any iteration  $j$  that  $\Phi_{u_{j+1}, \ell_{j+1}}(A_{j+1}) \leq \Phi_{u_j, \ell_j}(A_j)$ .*

*Proof.* Let  $w_1 w_1^T, \dots, w_{N_j} w_{N_j}^T$  be the matrices picked in iteration  $j$ , and define for any  $0 \leq i \leq N_j$  that  $B_i = A_j + \sum_{t=1}^i w_t w_t^T$ . We study the change of the potential function after adding a rank-1 matrix within each iteration. For this reason, we use

$$\bar{\Delta}_u = \frac{\Delta_{u,j}}{N_j} = (1 + 2\varepsilon) \cdot \frac{\varepsilon}{q \cdot \sum_{t=1}^m R_t(A_j, u_j, \ell_j)},$$

and

$$\bar{\Delta}_\ell = \frac{\Delta_{\ell,j}}{N_j} = (1 - 2\varepsilon) \cdot \frac{\varepsilon}{q \cdot \sum_{t=1}^m R_t(A_j, u_j, \ell_j)}$$

to express the average change of the barrier values  $\Delta_{u,j}$  and  $\Delta_{\ell,j}$ . We further define for  $0 \leq j \leq N_j$  that  $\hat{u}_i = u_j + i \cdot \bar{\Delta}_u$ ,  $\hat{\ell}_i = \ell_j + i \cdot \bar{\Delta}_\ell$ .

Assuming  $\mathbf{0} \preceq W_j \preceq \frac{1}{2}(u_j \mathbf{I} - A_j)$ , we claim that

$$w_i w_i^T \preceq \frac{2\varepsilon}{q} \cdot (\hat{u}_i \mathbf{I} - B_{i-1}) \text{ and } w_i w_i^T \preceq \frac{2\varepsilon}{q} \cdot (B_{i-1} - \hat{\ell}_i \mathbf{I}), \quad (4.13)$$

for any  $1 \leq i \leq N_j$ . Based on this, we apply Lemma 4.3.2 and get that

$$\begin{aligned}
\Phi_{\hat{u}_i, \hat{\ell}_i}(B_{i-1} + w_i w_i^T) &\leq \Phi_{\hat{u}_i, \hat{\ell}_i}(B_{i-1}) + q(1 + 2\varepsilon) \mathrm{Tr} \left( (\hat{u}_i \mathbf{I} - B_{i-1})^{-(q+1)} \mathbf{E} [w_i w_i^T] \right) \\
&\quad - q(1 - 2\varepsilon) \mathrm{Tr} \left( (B_{i-1} - \hat{\ell}_i \mathbf{I})^{-(q+1)} \mathbf{E} [w_i w_i^T] \right) \\
&= \Phi_{\hat{u}_i, \hat{\ell}_i}(B_{i-1}) + q \cdot \bar{\Delta}_u \cdot \mathrm{Tr} \left( (\hat{u}_i \mathbf{I} - B_{i-1})^{-(q+1)} \right) \\
&\quad - q \cdot \bar{\Delta}_\ell \cdot \mathrm{Tr} \left( (B_{i-1} - \hat{\ell}_i \mathbf{I})^{-(q+1)} \right). \quad (4.14)
\end{aligned}$$

We define a function  $f_i : \mathbb{R} \rightarrow \mathbb{R}$  by

$$f_i(x) = \mathrm{Tr} \left( (\hat{u}_{i-1} + x \cdot \bar{\Delta}_u) \mathbf{I} - B_{i-1} \right)^{-q} + \mathrm{Tr} \left( B_{i-1} - (\hat{\ell}_{i-1} + x \cdot \bar{\Delta}_\ell) \mathbf{I} \right)^{-q}.$$

Notice that

$$\begin{aligned} \frac{df_i(x)}{dx} &= -q \cdot \bar{\Delta}_u \cdot \text{Tr} \left( (\hat{u}_{i-1} + x \cdot \bar{\Delta}_u) I - B_{i-1} \right)^{-(q+1)} \\ &\quad + q \cdot \bar{\Delta}_\ell \cdot \text{Tr} \left( B_{i-1} - (\hat{\ell}_{i-1} + x \cdot \bar{\Delta}_\ell) I \right)^{-(q+1)}. \end{aligned}$$

Since  $f$  is convex, we have that

$$\left. \frac{df_i(x)}{dx} \right|_{x=1} \geq f_i(1) - f_i(0) = \Phi_{\hat{u}_i, \hat{\ell}_i}(B_{i-1}) - \Phi_{\hat{u}_{i-1}, \hat{\ell}_{i-1}}(B_{i-1}). \quad (4.15)$$

Putting (4.14) and (4.15) together, we have that

$$\Phi_{\hat{u}_i, \hat{\ell}_i}(B_i) \leq \Phi_{\hat{u}_i, \hat{\ell}_i}(B_{i-1}) - \left. \frac{df_i(x)}{dx} \right|_{x=1} \leq \Phi_{\hat{u}_{i-1}, \hat{\ell}_{i-1}}(B_{i-1}).$$

Repeating this argument we have that

$$\Phi_{u_{j+1}, \ell_{j+1}}(A_{j+1}) = \Phi_{\hat{u}_{N_j}, \hat{\ell}_{N_j}}(B_{N_j}) \leq \Phi_{\hat{u}_0, \hat{\ell}_0}(B_0) = \Phi_{u_j, \ell_j}(A_j),$$

which proves the statement.

So, it suffices to prove (4.13). Since  $vv^T \preceq (v^T B^{-1}v)B$  for any vector  $v$  and PSD matrix  $B$ , we have that

$$\frac{v_i v_i^T}{R_i(A_j, u_j, \ell_j)} \preceq \frac{v_i v_i^T}{v_i^T (u_j I - A_j)^{-1} v_i} \preceq u_j I - A_j.$$

By the assumption of  $W_j \preceq \frac{1}{2}(u_j I - A_j)$ , it holds that

$$w_i w_i^T = \frac{\varepsilon}{q R_i(A_j, u_j, \ell_j)} v_i v_i^T \preceq \frac{\varepsilon}{q} (u_j I - A_j) \preceq \frac{2\varepsilon}{q} (\hat{u}_i I - B_{i-1}).$$

This proves the first statement of the claim.

For the second statement, notice that

$$\ell_{j+1} - \ell_j = \Delta_{\ell, j} \leq \frac{\varepsilon N_j}{q \sum_{t=1}^m R_t(A_j, u_j, \ell_j)} \leq \frac{1}{2} \cdot \lambda_{\min}(A_j - \ell_j I),$$

and hence

$$w_i w_i^T \preceq \frac{\varepsilon}{q} (A_j - \ell_j I) \preceq \frac{2\varepsilon}{q} (A_j - \hat{\ell}_i I) \preceq \frac{2\varepsilon}{q} (B_{i-1} - \hat{\ell}_i I).$$

□

### ■ 4.3.2 Analysis of the approximation guarantee

In this subsection we will prove that the algorithm produces a linear-sized  $O(\varepsilon)$ -spectral sparsifier. We assume that the algorithm finishes after  $k$  iterations, and will prove that the condition number of  $A_k$  is small, which follows from our setting of parameters.

**Lemma 4.3.4.** *The output matrix  $A_k$  has condition number at most  $1 + O(\varepsilon)$ .*

*Proof.* Since the condition number of  $A_k$  is at most  $\frac{u_k}{\ell_k} = \left(1 - \frac{u_k - \ell_k}{u_k}\right)^{-1}$ , it suffices to prove that  $(u_k - \ell_k)/u_k = O(\varepsilon)$ .

Since the increase rate of  $\Delta_{u, j} - \Delta_{\ell, j}$  with respect to  $\Delta_{u, j}$  for any iteration  $j$  is

$$\frac{\Delta_{u, j} - \Delta_{\ell, j}}{\Delta_{u, j}} = \frac{(1 + 2\varepsilon) - (1 - 2\varepsilon)}{1 + 2\varepsilon} = \frac{4\varepsilon}{1 + 2\varepsilon} \leq 4\varepsilon,$$

we have that

$$\begin{aligned} \frac{u_k - \ell_k}{u_k} &= \frac{2 \cdot (2n)^{1/q} + \sum_{j=0}^{k-1} (\Delta_{u,j} - \Delta_{\ell,j})}{(2n)^{1/q} + \sum_{j=0}^{k-1} \Delta_{u,j}} \\ &\leq \frac{2 \cdot (2n)^{1/q} + \sum_{j=0}^{k-1} (\Delta_{u,j} - \Delta_{\ell,j})}{(2n)^{1/q} + (4\varepsilon)^{-1} \sum_{j=0}^{k-1} (\Delta_{u,j} - \Delta_{\ell,j})}. \end{aligned}$$

By the ending condition of the algorithm, it holds that  $u_k - \ell_k \geq 4 \cdot (2n)^{1/q}$ , i.e.

$$\sum_{j=0}^{k-1} (\Delta_{u,j} - \Delta_{\ell,j}) \geq 2 \cdot (2n)^{1/q}.$$

Hence, it holds that

$$\frac{u_k - \ell_k}{u_k} \leq \frac{2 \cdot (2n)^{1/q} + 2 \cdot (2n)^{1/q}}{(2n)^{1/q} + (4\varepsilon)^{-1} 2 \cdot (2n)^{1/q}} \leq 8\varepsilon,$$

which finishes the proof.  $\square$

Now we prove that the algorithm finishes in  $O\left(\frac{qn^{3/q}}{\varepsilon^2}\right)$  iterations, and picks  $O(qn/\varepsilon^2)$  vectors in total.

**Lemma 4.3.5.** *The following statements hold:*

1. *With probability at least  $4/5$ , the algorithm finishes in  $\frac{10qn^{3/q}}{\varepsilon^2}$  iterations.*
2. *With probability at least  $4/5$ , the algorithm chooses at most  $10qn/\varepsilon^2$  vectors.*

*Proof.* Notice that after iteration  $j$  the barrier gap  $u_j - \ell_j$  is increased by

$$\begin{aligned} \Delta_{u,j} - \Delta_{\ell,j} &= \frac{4\varepsilon^2}{q} \frac{N_j}{\sum_{i=1}^m R_i(A_j, u_j, \ell_j)} \\ &= \frac{4\varepsilon^2}{q} \frac{1}{n^{2/q}} \cdot \min\{\lambda_{\min}(u_j I - A_j), \lambda_{\min}(A_j - \ell_j I)\} \\ &\geq \frac{4\varepsilon^2}{q} \frac{1}{n^{2/q}} \cdot (\Phi_{u_j, \ell_j}(A_j))^{-1/q}. \end{aligned}$$

Since the algorithm finishes within  $k$  iterations if  $\sum_{j=0}^{k-1} (\Delta_{u,j} - \Delta_{\ell,j}) \geq 2 \cdot (2n)^{1/q}$ , it holds that

$$\begin{aligned} \Pr[\text{algorithm finishes within } k \text{ iterations}] &\geq \Pr\left[\sum_{j=0}^{k-1} (\Delta_{u,j} - \Delta_{\ell,j}) \geq 2 \cdot (2n)^{1/q}\right] \\ &\geq \Pr\left[\sum_{j=0}^{k-1} \frac{4\varepsilon^2}{qn^{2/q}} \cdot (\Phi_{u_j, \ell_j}(A_j))^{-1/q} \geq 2 \cdot (2n)^{1/q}\right] \\ &= \Pr\left[\sum_{j=0}^{k-1} (\Phi_{u_j, \ell_j}(A_j))^{-1/q} \geq \frac{q}{2\varepsilon^2} \cdot (2n^3)^{1/q}\right] \\ &\geq \Pr\left[\sum_{j=0}^{k-1} (\Phi_{u_j, \ell_j}(A_j))^{1/q} \leq 2 \cdot \frac{k^2 \varepsilon^2}{q} \cdot \left(\frac{1}{2n^3}\right)^{1/q}\right], \end{aligned}$$

where the last inequality follows from the fact that

$$\left( \sum_{j=0}^{k-1} (\Phi_{u_j, \ell_j}(A_j))^{-1/q} \right) \cdot \left( \sum_{j=0}^{k-1} (\Phi_{u_j, \ell_j}(A_j))^{1/q} \right) \geq k^2.$$

By Lemma 4.3.1, every picked matrix  $W_j$  in iteration  $j$  satisfies  $0 \preceq W_j \preceq \frac{1}{2} \cdot (u_j I - A)$  with probability at least  $1 - \frac{\varepsilon^2}{100qn}$ , and with probability 9/10 all matrices picked in  $k = 10qn/\varepsilon^2$  iterations satisfy the condition above. Also, by Lemma 4.3.3 we have that

$$\sum_{j=0}^{k-1} (\Phi_{u_j, \ell_j}(A_j))^{1/q} = \sum_{j=0}^{k-1} (\Phi_{u_j, \ell_j}(A_j))^{1/q} \leq \sum_{j=0}^{k-1} (\Phi_{u_j, \ell_j}(A_j))^{1/q} \leq k, \quad (4.16)$$

since the initial value of the potential function is at most 1. Therefore, it holds that

$$\begin{aligned} & \Pr[\text{algorithm finishes in more than } k \text{ iterations}] \\ & \leq \Pr \left[ \sum_{j=0}^{k-1} (\Phi_{u_j, \ell_j}(A_j))^{1/q} \geq 2 \cdot \frac{k^2 \varepsilon^2}{q} \cdot \left( \frac{1}{2n^3} \right)^{1/q} \right] \\ & \leq \Pr \left[ \sum_{j=0}^{k-1} (\Phi_{u_j, \ell_j}(A_j))^{1/q} \geq 2 \cdot \frac{k^2 \varepsilon^2}{q} \cdot \left( \frac{1}{2n^3} \right)^{1/q} \text{ and } \forall j : W_j \preceq \frac{1}{2} (u_j I - A_j) \right] \\ & \quad + \Pr \left[ \exists j : W_j \not\preceq \frac{1}{2} (u_j I - A_j) \right] \leq \frac{q}{2 \cdot k \varepsilon^2} \cdot (2n^3)^{1/q} + 1/10 \leq 1/5, \end{aligned}$$

where the second last inequity follows from Markov's inequality and (4.16), and the last inequality follows by our choice of  $k$ . This proves the first statement.

Now for the second statement. Notice that for every vector chosen in iteration  $j$ , the barrier gap  $\Delta_{u,j} - \Delta_{\ell,j}$  is increased on average by

$$\frac{\Delta_{u,j} - \Delta_{\ell,j}}{N_j} = \frac{4\varepsilon^2}{q \sum_{i=1}^m R_i(A_j, u_j, \ell_j)}.$$

To bound  $R_i(A_j, u_j, \ell_j)$ , let the eigenvalues of matrix  $A_j$  be  $\lambda_1, \dots, \lambda_n$ . Then, it holds that

$$\begin{aligned} \sum_{i=1}^m R_i(A_j, u_j, \ell_j) &= \sum_{i=1}^m \vec{v}_i^T (u_j I - A_j)^{-1} \mathbf{v}_i + \sum_{i=1}^m \vec{v}_i^T (A_j - \ell_j I)^{-1} \vec{v}_i \\ &= \sum_{i=1}^n \frac{1}{u_j - \lambda_i} + \sum_{i=1}^n \frac{1}{\lambda_i - \ell_j} \\ &\leq \left( \sum_{i=1}^n (u_j - \lambda_i)^{-q} + \sum_{i=1}^n (\lambda_i - \ell_j)^{-q} \right)^{1/q} (2n)^{1-1/q} \\ &= (\Phi_{u_j, \ell_j}(A_j))^{1/q} \cdot (2n)^{1-1/q}. \end{aligned}$$

Therefore, we have that

$$\frac{\Delta_{u,j} - \Delta_{\ell,j}}{N_j} \geq \frac{4\varepsilon^2}{q} \cdot \frac{1}{(2n)^{1-1/q} \cdot (\Phi_{u_j, \ell_j}(A_j))^{1/q}}. \quad (4.17)$$

Let  $v_1, \dots, v_z$  be the vectors sampled by the algorithm, and  $v_j$  is picked in iteration  $\tau_j$ , where



$1 \leq j \leq z$ . We first assume that the algorithm could check the ending condition after adding every single vector. In such case, it holds that

$$\begin{aligned} & \Pr[\text{algorithm finishes after choosing } z \text{ vectors}] \\ & \geq \Pr \left[ \sum_{j=1}^z \frac{4\varepsilon^2}{q} \cdot \frac{1}{(2n)^{1-1/q} \cdot (\Phi_{u_{\tau_j}, \ell_{\tau_j}}(A_{\tau_j}))^{1/q}} \geq 2 \cdot (2n)^{1/q} \right] \\ & = \Pr \left[ \sum_{j=1}^z (\Phi_{u_{\tau_j}, \ell_{\tau_j}}(A_{\tau_j}))^{-1/q} \geq qn/\varepsilon^2 \right]. \end{aligned}$$

Following the same proof as the first part and noticing that in the final iteration the algorithm chooses at most  $O(n)$  extra vectors, we obtain the second statement.  $\square$

### ■ 4.3.3 Proof of the main results

Now we analyze the runtime of the algorithm, and prove the main results. We first analyze the algorithm for sparsifying the sum of rank-1 PSD matrices, and prove Theorem 4.1.1.

*Proof of Theorem 4.1.1.* By Theorem 3.7.3, we can first find a  $\varepsilon$ -approximation with  $O(n \log n / \varepsilon^2)$  vectors in time  $\tilde{O}(Z + n^\omega \varepsilon^{-2})$ . Note that if  $q < \log n$ , we are done.

Otherwise, by Lemma 4.3.5, with probability at least  $4/5$  the algorithm chooses at most  $10qn/\varepsilon^2$  vectors, and by Lemma 4.3.4 the condition number of  $A_k$  is at most  $1 + O(\varepsilon)$ , implying that the matrix  $A_k$  is a  $(1 + O(\varepsilon))$ -approximation of  $I$ . These two results together prove that  $A_k$  is a linear-sized spectral sparsifier.

For the runtime, Lemma 4.3.5 proves that the algorithm finishes in  $\frac{10qn^{3/q}}{\varepsilon^2}$  iterations, and it is easy to see that all the required quantities in each iteration can be approximately computed in  $\tilde{O}(m \cdot n^{\omega-1})$  time using fast matrix multiplication. Therefore, the total runtime of the algorithm is  $\tilde{O}(Z + n^\omega \varepsilon^{-2} + \frac{q \cdot m}{\varepsilon^2} \cdot n^{\omega-1+3/q})$ . Using  $q < \log n$  and  $m = O(n \log n / \varepsilon^2)$ , we have

$$\tilde{O}\left(Z + n^{\omega+3/q} \varepsilon^{-4}\right).$$

$\square$

Next we apply our algorithm in the graph setting, and prove Theorem 4.1.2.

*Proof of Theorem 4.1.2.* For any edge  $e = \{u, v\}$  in graph  $G = (V, E)$ , we define  $b_e \in \mathbb{R}^n$ , where  $b_e(w) = 1$  if  $w = u$ ,  $b_e(w) = -1$  if  $w = v$ , and  $b_e(w) = 0$  otherwise. Then the Laplacian matrix of  $G$  can be written as  $L = \sum_{e \in E[G]} b_e b_e^T$ . By setting  $v_e = L^{-1/2} b_e$  for  $e \in E[G]$ , it is easy to see that constructing a spectral sparsifier of  $G$  is equivalent to sparsifying the matrix  $\sum_{e \in E[G]} v_e v_e^T$ . We apply the same analysis as in the proof of Theorem 4.1.1, and hence the output is a linear-sized spectral sparsifier of graph  $G$ . So it suffices to analyze the runtime of the algorithm.

By Lemma 4.3.1 and the Union Bound, with probability at least  $9/10$  all the matrices picked in  $k = \frac{10qn^{3/q}}{\varepsilon^2}$  iterations satisfy  $W_j \preceq \frac{1}{2}(u_j I - A_j)$ . Conditioning on the event, with constant probability  $\mathbf{E}[\Phi_{u_j, \ell_j}(A_j)] \leq 2$  for all iterations  $j$ , and by Markov's inequality with high probability it holds that  $\Phi_{u_j, \ell_j}(A_j) = O(qn/\varepsilon^2)$  for all iterations  $j$ .

On the other hand, notice that it holds for any  $1 \leq j \leq n$  that

$$(u - \lambda_j)^{-q} \leq \sum_{i=1}^n (u - \lambda_i)^{-q} < \Phi_{u, \ell}(\mathbf{A}),$$

where  $\lambda_1, \dots, \lambda_n$  are the eigenvalues of  $A$ . This implies that  $\lambda_j < u - (\Phi_{u,\ell}(\mathbf{A}))^{-1/q}$ . Similarly, it holds that  $\lambda_j > \ell + (\Phi_{u,\ell}(\mathbf{A}))^{-1/q}$  for any  $1 \leq j \leq n$ . Therefore, we have that

$$\left( \ell_j + O\left(\left(\frac{\varepsilon^2}{qn}\right)^{1/q}\right) \right) \mathbf{I} \prec \mathbf{A}_j \prec \left( u_j - O\left(\left(\frac{\varepsilon^2}{qn}\right)^{1/q}\right) \right) \mathbf{I}.$$

Since both of  $u_j$  and  $\ell_j$  are of the order  $O(n^{1/q})$ , we set  $\eta = O((\varepsilon/n)^{2/q})$  and obtain that

$$(\ell_j + |\ell_j|\eta)\mathbf{I} \prec \mathbf{A}_j \prec (1 - \eta)u_j\mathbf{I}.$$

Hence, we apply Lemma 4.4.3 and Lemma 4.4.4 to compute all required quantities in each iteration up to constant approximation in time

$$\tilde{O}\left(\frac{m}{\varepsilon^2 \cdot \eta}\right) = \tilde{O}\left(\frac{m \cdot n^{2/q}}{\varepsilon^{2+2/q}}\right).$$

Since by Lemma 4.3.5 the algorithm finishes in  $\frac{10qn^{3/q}}{\varepsilon^2}$  iterations with probability at least  $4/5$ , the total runtime of the algorithm is

$$\tilde{O}\left(\frac{q \cdot m \cdot n^{5/q}}{\varepsilon^{4+4/q}}\right).$$

Again, we can assume  $m = \tilde{O}(n \log n / \varepsilon^2)$  and  $q < \log n$  by existing algorithm by using existing algorithm for graph sparsification. Hence, we have the promised runtime.  $\square$

## ■ 4.4 Appendix

In this section we study fast approximation of the required quantities for constructing graph sparsifiers, and show that the number of sampled vectors  $N_j$  and the relative leverage scores  $\{R_i(A_j, u_j, \ell_j)\}_{i=1}^m$  used in each iteration can be approximately computed in almost-linear time. For simplicity we drop the subscript  $j$  expressing the iterations in this subsection. We will assume that the following assumption holds on  $A$ , and show that the matrix under consideration always satisfies this condition.

**Assumption 4.4.1.** *Let  $L$  and  $\tilde{L}$  be the Laplacian matrices of graph  $G$  and its subgraph after reweighting. Let  $A = L^{-1/2}\tilde{L}L^{-1/2}$ , and assume that*

$$(\ell + |\ell|\eta) \cdot \mathbf{I} \prec A \prec (1 - \eta)u \cdot \mathbf{I}$$

holds for some  $0 < \eta < 1$ .

**Lemma 4.4.2.** *Under Assumption 4.4.1, the following statements hold:*

1. *We can construct a matrix  $S_u$  such that*

$$S_u \approx_{\varepsilon/10} (u\mathbf{I} - A)^{-1/2},$$

*and  $S_u = p(A)$  for a polynomial  $p$  of degree  $O\left(\frac{\log(1/\varepsilon\eta)}{\eta}\right)$ .*

2. *We can construct a matrix  $S_\ell$  such that*

$$S_\ell \approx_{\varepsilon/10} (A - \ell\mathbf{I})^{-1/2}.$$

*Moreover,  $S_\ell$  is of the form  $(A')^{-1/2}q((A')^{-1})$ , where  $q$  is a polynomial of degree  $O\left(\frac{\log(1/\varepsilon\eta)}{\eta}\right)$  and  $A' = L^{-1/2}L'L^{-1/2}$  for some Laplacian matrix  $L'$ .*

*Proof.* By Taylor expansion, it holds that

$$(1-x)^{-1/2} = 1 + \sum_{k=1}^{\infty} \prod_{j=0}^{k-1} \left(j + \frac{1}{2}\right) \frac{x^k}{k!}.$$

We define for any  $T \in \mathbb{N}$  that

$$p_T(x) = 1 + \sum_{k=1}^T \prod_{j=0}^{k-1} \left(j + \frac{1}{2}\right) \frac{x^k}{k!}.$$

Then, it holds for any  $0 < x < 1 - \eta$  that

$$\begin{aligned} p_T(x) &\leq (1-x)^{-1/2} = p_T(x) + \sum_{k=T+1}^{\infty} \prod_{j=0}^{k-1} \left(j + \frac{1}{2}\right) \frac{x^k}{k!} \\ &\leq p_T(x) + \sum_{k=T+1}^{\infty} x^k \\ &\leq p_T(x) + \frac{(1-\eta)^{T+1}}{\eta}. \end{aligned}$$

Hence, it holds that

$$(uI - A)^{-1/2} = u^{-1/2}(I - u^{-1}A)^{-1/2} \succeq u^{-1/2}p_T(u^{-1}A),$$

and

$$(uI - A)^{-1/2} \preceq u^{-1/2} \left( p_T(u^{-1}A) + \frac{(1-\eta)^{T+1}}{\eta} \cdot I \right),$$

since  $u^{-1}A \preceq (1-\eta)I$ . Notice that  $u^{-1/2}I \preceq (uI - A)^{-1/2}$ , and therefore

$$(uI - A)^{-1/2} \preceq u^{-1/2}p_T(u^{-1}A) + \frac{(1-\eta)^{T+1}}{\eta} \cdot (uI - A)^{-1/2}.$$

Setting  $T = \frac{c \log(1/(\varepsilon\eta))}{\eta}$  for some constant  $c$  and defining  $S_u = u^{-1/2}p_T(u^{-1}A)$  gives us that

$$S_u \approx_{\varepsilon/10} (uI - A)^{-1/2}.$$

Now for the second statement. Our construction of  $S_\ell$  is based on the case distinction ( $\ell > 0$ , and  $\ell \leq 0$ ).

Case 1 ( $\ell > 0$ ): Notice that

$$(A - \ell I)^{-1/2} = A^{-1/2}(I - \ell A^{-1})^{-1/2},$$

and

$$p_T(\ell A^{-1}) \preceq (I - \ell A^{-1})^{-1/2} \preceq p_T(\ell A^{-1}) + \frac{(1-\eta/2)^{T+1}}{\eta/2} \cdot I.$$

Using the same analysis as before, we have that

$$A^{-1/2}(I - \ell A^{-1})^{-1/2} \approx_{\varepsilon/10} A^{-1/2}p_T(\ell A^{-1}).$$

By defining  $S_\ell = A^{-1/2}p_T(\ell A^{-1})$ , i.e.,  $A' = A$  and  $q((A')^{-1}) = p_T(\ell A^{-1})$ , we have that

$$S_\ell \approx_{\varepsilon/10} (A - \ell I)^{-1/2}.$$

Case 2 ( $\ell \leq 0$ ): We look at the matrix

$$A - \ell I = L^{-1/2} \tilde{L} L^{-1/2} - \ell I = L^{-1/2} (\tilde{L} - \ell L) L^{-1/2}.$$

Notice that  $\tilde{L} - \ell L$  is a Laplacian matrix, and hence this reduces to the case of  $\ell = 0$ , for which we simply set  $S_\ell = (A - \ell I)^{-1/2}$ . Therefore, we can write  $S_\ell$  as a desired form, where  $A' = A - \ell I$  and polynomial  $q = 1$ .  $\square$

The following two lemmas present nearly-linear time algorithms for computing the required quantities within each iteration. We remark that an almost-linear time algorithm for computing similar quantities was shown in [9].

**Lemma 4.4.3.** *Let  $A = \sum_{i=1}^m v_i v_i^T$ , and suppose that  $A$  satisfies Assumption 4.4.1. Then, we can compute  $\{r_i\}_{i=1}^m$  and  $\{t_i\}_{i=1}^m$  in  $\tilde{O}\left(\frac{m}{\varepsilon^2 \eta}\right)$  time such that*

$$(1 - \varepsilon)r_i \leq v_i^T (uI - A)^{-1} v_i \leq (1 + \varepsilon)r_i,$$

and

$$(1 - \varepsilon)t_i \leq v_i^T (A - \ell I)^{-1} v_i \leq (1 + \varepsilon)t_i.$$

*Proof.* Define  $u_i = L^{1/2} v_i$  for any  $1 \leq i \leq m$ . By Lemma 4.4.2, we have that

$$\begin{aligned} v_i^T (uI - A)^{-1} v_i &\approx_{3\varepsilon/10} \|p(A)v_i\|^2 \\ &= \left\| p\left(L^{-1/2} \tilde{L} L^{-1/2}\right) L^{-1/2} u_i \right\|^2 \\ &= \left\| L^{1/2} p\left(L^{-1} \tilde{L}\right) L^{-1} u_i \right\|^2. \end{aligned}$$

Let  $L = B^T B$  for some  $B \in \mathbb{R}^{m \times n}$ . Then, it holds that

$$v_i^T (uI - A)^{-1} v_i \approx_{3\varepsilon/10} \left\| B p\left(L^{-1} \tilde{L}\right) L^{-1} u_i \right\|^2.$$

We invoke the Johnson-Lindenstrauss Lemma and find a random matrix  $Q \in \mathbb{R}^{O(\log n/\varepsilon^2) \times m}$ . With high probability, it holds that

$$v_i^T (uI - A)^{-1} v_i \approx_{4\varepsilon/10} \left\| Q B p\left(L^{-1} \tilde{L}\right) L^{-1} u_i \right\|^2.$$

We apply a nearly-linear time Laplacian solver to compute  $\left\| Q B p\left(L^{-1} \tilde{L}\right) L^{-1} u_i \right\|^2$  for all  $\{u_i\}_{i=1}^m$  up to  $(1 \pm \varepsilon/10)$ -multiplicative error in time  $\tilde{O}\left(\frac{m}{\varepsilon^2 \eta}\right)$ . This gives the desired  $\{r_i\}_{i=1}^m$ .

The computation for  $\{t_i\}_{i=1}^m$  is similar. By Lemma 4.4.2, it holds for any  $1 \leq i \leq m$  that

$$\begin{aligned} v_i^T (A - \ell I)^{-1} v_i &\approx_{3\varepsilon/10} \left\| (A')^{-1/2} q((A')^{-1}) v_i \right\|^2 \\ &= \left\| (A')^{-1/2} q\left(L^{1/2} (L')^{-1} L^{1/2}\right) L^{-1/2} u_i \right\|^2 \\ &= \left\| (A')^{-1/2} L^{-1/2} q(L(L')^{-1}) u_i \right\|^2. \end{aligned}$$

Let  $L' = (B')^T(B')$  for some  $B' \in \mathbb{R}^{m \times n}$ . Then, it holds that

$$\begin{aligned} v_i^T (A - \ell I)^{-1} v_i &\approx_{3\varepsilon/10} \left\| (L')^{-1/2} q(L(L')^{-1}) u_i \right\|^2 \\ &= \left\| (L')^{1/2} (L')^{-1} q(L(L')^{-1}) u_i \right\|^2 \\ &= \left\| B'(L')^{-1} q(L(L')^{-1}) u_i \right\|^2. \end{aligned}$$

We invoke the Johnson-Lindenstrauss Lemma and a nearly-linear time Laplacian solver as before to obtain required  $\{t_i\}_{i=1}^m$ . The total runtime is  $\tilde{O}\left(\frac{m}{\eta\varepsilon^2}\right)$ .  $\square$

**Lemma 4.4.4.** *Under Assumption 4.4.1, we can compute values  $\alpha, \beta$  in  $\tilde{O}\left(\frac{m}{\eta\varepsilon^3}\right)$  time such that*

$$(1 - \varepsilon)\alpha \leq \lambda_{\min}(uI - A) \leq (1 + \varepsilon)\alpha$$

and

$$(1 - \varepsilon)\beta \leq \lambda_{\min}(A - \ell I) \leq (1 + \varepsilon)\beta.$$

*Proof.* By Lemma 4.4.2, we have that  $S_u \approx_{\varepsilon/10} (uI - A)^{-1/2}$ . Hence,  $\lambda_{\max}(S_u)^{-2} \approx_{3\varepsilon/10} \lambda_{\min}(uI - A)$ , and it suffices to estimate  $\lambda_{\max}(S_u)$ . Since

$$\lambda_{\max}(S_u) \leq \left( \text{Tr}(S_u^{2k}) \right)^{1/2k} \leq n^{1/2k} \lambda_{\max}(S_u),$$

by picking  $k = \log n/\varepsilon$  we have that  $(\text{Tr}(S_u^{2k}))^{1/2k} \approx_{\varepsilon/2} \lambda_{\max}(S_u)$ . Notice that

$$\text{Tr}(S_u^{2k}) = \text{Tr}\left(p^{2k}\left(L^{-1/2}\tilde{L}L^{-1/2}\right)\right) = \text{Tr}\left(p^{2k}\left(L^{-1}\tilde{L}\right)\right).$$

Set  $\tilde{L} = \tilde{B}^T \tilde{B}$  for some matrix  $\tilde{B} \in \mathbb{R}^{m \times n}$ , and we have that

$$\text{Tr}(S_u^{2k}) = \text{Tr}\left(p^{2k}\left(\tilde{B}L^{-1}\tilde{B}^T\right)\right).$$

Since we can apply  $p^k\left(\tilde{B}L^{-1}\tilde{B}^T\right)$  to vectors in  $\tilde{O}\left(\frac{m}{\eta\varepsilon}\right)$  time, we invoke the Johnson-Lindenstrauss Lemma and approximate  $\text{Tr}(S_u^{2k})$  in  $\tilde{O}\left(\frac{m}{\eta\varepsilon^3}\right)$  time.

We approximate  $\lambda_{\min}(A - \ell I)$  in a similar way. Notice that

$$\begin{aligned} \text{Tr}(S_\ell^{4k}) &= \text{Tr}\left((A')^{-1/2} q((A')^{-1})\right)^{4k} \\ &= \text{Tr}\left(q((A')^{-1})(A')^{-1} q((A')^{-1})\right)^{2k}. \end{aligned}$$

Let  $z$  be a polynomial defined by  $z(x) = xq^2(x)$  and  $L' = (B')^T(B')$ . Then, we have that

$$\text{Tr}(S_\ell^{4k}) = \text{Tr}\left(z^{2k}((A')^{-1})\right) = \text{Tr}\left(z^{2k}\left(L^{1/2}(L')^{-1}L^{1/2}\right)\right).$$

Applying the same analysis as before, we can estimate the trace in  $\tilde{O}\left(\frac{m}{\eta\varepsilon^3}\right)$  time.  $\square$



## ■ 5.1 Introduction

In this chapter, we study sparsification for  $\ell_p$  regression problems.

For  $\ell_2$  regression, Lemma 2.3.3 shows that given a tall matrix, we can sample  $O(n \log n)$  rows of  $\mathbf{A}$  to form a matrix  $\tilde{\mathbf{A}}$  such that, for all vectors  $\mathbf{x}$ ,  $\|\tilde{\mathbf{A}}\mathbf{x}\|_2 \approx \|\mathbf{A}\mathbf{x}\|_2$ . Similarly, for  $\ell_p$  regression, it is known that an approximation can be obtained by sampling rows with a certain probability (called  $\ell_p$  Lewis weights). The main goal of this chapter is to show how to approximate  $\ell_p$  Lewis weights by solving  $\tilde{O}(1)$  linear systems (Theorem 5.3.4). Since we know how to solve linear systems in input sparsity time, this allows us to compute  $\ell_p$  Lewis Weights in nearly input sparsity time.

This chapter is organized as follows: In Section 5.2, we define  $\ell_p$  Lewis weights; in Section 5.3, we show how to approximate  $\ell_p$  Lewis weights efficiently; in Section 5.4, we discuss some geometric implications of  $\ell_p$  Lewis weights.

## ■ 5.2 Lewis Weights

Unlike leverage scores,  $\ell_p$  Lewis weights do not have closed formulas and are defined by the following non-linear equation.

**Definition 5.2.1.** For a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $1 \leq p < \infty$ , the  $\ell_p$  Lewis weight is the unique vector  $\mathbf{w} \in \mathbb{R}^m$  such that

$$\mathbf{w}_i = \sigma_i(\mathbf{W}^{1/2-1/p}\mathbf{A}) \quad \text{for all } i \in [m]$$

where  $\sigma_i(\mathbf{A}) = (\mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1}\mathbf{A})_{ii}$ .

The following result show how to use  $\ell_p$  Lewis weights to construct sparsifiers for  $\ell_p$  regression. They follow from a collection of results for Banach spaces (See [56] for the summary of each result).

**Theorem 5.2.2** ([169, 38, 245, 246, 56]). *Let  $\mathbf{w}$  be the  $\ell_p$  Lewis weight of some given matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $1 \leq p < \infty$ . Given an error parameter  $0 < \varepsilon < 1$  and sampling values  $\mathbf{p}$  that*

$$\mathbf{p}_i \geq f(p, n, \varepsilon, \delta)\mathbf{w}_i$$

$p$	Success probability: $1 - \delta$	Number of rows: $n \cdot f(p, n, \varepsilon, \delta)$
$p = 1$	$1 - 1/n^C$	$O(n \log(n/\varepsilon)/\varepsilon^2)$
$p = 1$	$1 - 1/C$	$O(n \log n/\varepsilon^2)$
$1 < p < 2$	$1 - 1/C$	$O(n \log(n/\varepsilon) \log(\log n/\varepsilon)^2/\varepsilon^2)$
$p > 2$	$1 - 1/n^C$	$O(n^{p/2} \log n \log(1/\varepsilon)/\varepsilon^5)$

**Table 5.1:**  $n \cdot f(p, n, \varepsilon, \delta)$ : Number of rows needed for  $\ell_p$  sampling. The constants in big O notation depend on  $C$ .

where  $f(p, n, \varepsilon, \delta)$  is shown in the Table 5.1. Let  $\mathbf{S}$  be the diagonal matrix with independently chosen entries.  $\mathbf{S}_{ii} = \frac{1}{p_i^{1/p}}$  with probability  $p_i$  and 0 otherwise. Then, with probability at least  $1 - \delta$ , we have that

$$(1 - \varepsilon)\|\mathbf{A}\mathbf{x}\|_p \leq \|\mathbf{S}\mathbf{A}\mathbf{x}\|_p \leq (1 + \varepsilon)\|\mathbf{A}\mathbf{x}\|_p$$

for all  $\mathbf{x} \in \mathbb{R}^n$ .

This theorem shows if we have a good upper estimate of  $\ell_p$  Lewis weights, then we can sample  $\tilde{O}(\max(n, n^{p/2}))$  rows to approximate the original matrix. Since  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_p = \|[\mathbf{A}, \mathbf{b}] \begin{pmatrix} \mathbf{x} \\ -1 \end{pmatrix}\|_p$ , Lewis weights allow us to reduce a  $\ell_p$  regression problem involving a tall matrix to a matrix with only  $\tilde{O}(\max(n, n^{p/2}))$  rows.

### ■ 5.3 Approximate Lewis Weights

In this section, we show how to compute  $\ell_p$  Lewis weights by solving  $\tilde{O}(1)$  linear systems for  $1 \leq p < \infty$ . For  $p < 4$ , this can be done by the following algorithm [56].

---

**Algorithm 5:** Approximate  $\ell_p$  Lewis weights for  $p < 4$

---

**Input:** a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , accuracy  $0 < \varepsilon < 1$ ,  $1 \leq p < 4$ .

Let  $\mathbf{v}_i^{(1)} = n/m$  for all  $i \in [m]$  and  $T = \Theta\left(\frac{\log(m/\varepsilon)}{4-p}\right)$ .

**for**  $k = 1, \dots, T$  **do**

Let  $\tilde{\sigma}$  be a  $c_1(4-p)\varepsilon$ -approximation of  $\sigma_i\left(\left(\mathbf{V}^{(k)}\right)^{1/2-1/p} \mathbf{A}\right)$  for some small enough constant  $c_1$ .

$\mathbf{v}_i^{(k+1)} = \left(\left(\mathbf{v}_i^{(k)}\right)^{2/p-1} \tilde{\sigma}_i\right)^{p/2}$ .

**end**

**Output:**  $\mathbf{v}^{(T+1)}$ .

---

**Theorem 5.3.1** ([56]). *For any matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , accuracy  $0 < \varepsilon < 1$  and  $1 \leq p < 4$ , the algorithm 5 outputs a vector  $\mathbf{v}$  in time  $\tilde{O}(\mathcal{T}\varepsilon^{-2})$  with high probability such that*

$$(1 - \varepsilon)\mathbf{w}_i \leq \mathbf{v}_i \leq (1 + \varepsilon)\mathbf{w}_i \tag{5.1}$$

where  $\mathbf{w}$  is the  $\ell_p$  Lewis weight for  $\mathbf{A}$  and  $\mathcal{T}$  be the time needed to apply  $(\mathbf{A}^\top \mathbf{D} \mathbf{A})^{-1}$  to a vector for any diagonal matrix  $\mathbf{D}$ .

This algorithm relies on the stability of  $\ell_p$  Lewis weights for  $1 \leq p < 4$ , namely, if we rescale each row of  $\mathbf{A}$  by a constant, the Lewis weight of  $\mathbf{A}$  changes only by a constant. However, this is not true for large  $p$ . For  $p \sim \infty$ , the Lewis weights can change arbitrarily even with a tiny rescaling and hence it is almost impossible to approximate Lewis weight in the sense of (5.1). To avoid this inability issue, we define the  $\varepsilon$ -approximate  $\ell^p$  Lewis weight as follows:

**Definition 5.3.2.** For a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $1 \leq p < \infty$ , we call a vector  $\mathbf{w} \in \mathbb{R}^m$  is an  $\varepsilon$ -approximate  $\ell^p$  Lewis weight for  $\mathbf{A}$  if

$$(1 - \varepsilon)\sigma_i(\mathbf{W}^{1/2-1/p} \mathbf{A}) \leq \mathbf{w}_i \leq (1 + \varepsilon)\sigma_i(\mathbf{W}^{1/2-1/p} \mathbf{A}).$$

For  $1 \leq p < 4$ , an  $\varepsilon$ -approximate  $\ell^p$  Lewis weight is indeed a  $O(\varepsilon)$  approximate of the true  $\ell^p$  Lewis weight in the sense of (5.1). Although this is not true for larger  $p$ , the following lemma shows that  $\varepsilon$ -approximate  $\ell^p$  Lewis weights suffice for the sampling purpose.



**Lemma 5.3.3.** *Given an error parameter  $0 < \varepsilon < \frac{1}{2}$ . Let  $\mathbf{w}$  be an  $\varepsilon$ -approximate  $\ell_p$  Lewis weight of some given matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $1 \leq p < \infty$ . Let  $\mathbf{p}$  be the sampling values such that*

$$\mathbf{p}_i \geq f(p, n, \varepsilon, \delta) \mathbf{w}_i$$

where  $f(p, n, \varepsilon, \delta)$  is shown in the Table 5.1. Let  $\mathbf{S}$  be the diagonal matrix with independently chosen entries.  $\mathbf{S}_{ii} = \frac{1}{\mathbf{p}_i^{1/p}}$  with probability  $\mathbf{p}_i$  and 0 otherwise. Then, with probability at least  $1 - \delta$ , we have that

$$(1 - 2\varepsilon) \|\mathbf{A}\mathbf{x}\|_p \leq \|\mathbf{S}\mathbf{A}\mathbf{x}\|_p \leq (1 + 2\varepsilon) \|\mathbf{A}\mathbf{x}\|_p$$

for all  $\mathbf{x} \in \mathbb{R}^n$ .

*Proof.* Let  $\mathbf{d}_i = \sigma_i(\mathbf{W}^{1/2-1/p}\mathbf{A})/\mathbf{w}_i$ ,  $\tilde{\mathbf{A}} = \mathbf{D}^{-1/2+1/p}\mathbf{A}$  and  $\tilde{\mathbf{W}} = \mathbf{D}\mathbf{W}$ . From the definition of  $\sigma_i$  (Definition 5.2.1), we see that

$$\begin{aligned} \sigma_i(\tilde{\mathbf{W}}^{1/2-1/p}\tilde{\mathbf{A}}) &= \sigma_i(\mathbf{W}^{1/2-1/p}\mathbf{A}) \\ &= \mathbf{d}_i \mathbf{w}_i = \tilde{\mathbf{w}}_i. \end{aligned}$$

Hence,  $\tilde{\mathbf{w}}$  is the  $\ell_p$  Lewis weight of  $\tilde{\mathbf{A}}$  and theorem 5.2.2 shows that

$$(1 - \varepsilon) \|\tilde{\mathbf{A}}\mathbf{x}\|_p \leq \|\mathbf{S}\tilde{\mathbf{A}}\mathbf{x}\|_p \leq (1 + \varepsilon) \|\tilde{\mathbf{A}}\mathbf{x}\|_p.$$

From the definition of  $\varepsilon$ -approximate  $\ell_p$  Lewis weight, we know that  $1 - \varepsilon \leq \mathbf{d}_i \leq 1 + \varepsilon$  and hence

$$(1 - 2\varepsilon) \|\mathbf{A}\mathbf{x}\|_p \leq \|\mathbf{S}\mathbf{A}\mathbf{x}\|_p \leq (1 + 2\varepsilon) \|\mathbf{A}\mathbf{x}\|_p.$$

□

Note that if we have a constant approximate of leverage scores, we can construct an arbitrarily good approximation by sampling more rows. However, this is not true for  $\ell_p$  Lewis weights; we need a good approximation of Lewis weight to get a good approximate of the original matrix.

Now, we have everything to describe our result, a simple multiplicative weight update algorithm.

---

**Algorithm 6:** Approximate  $\ell_p$  Lewis weights for  $p \geq 2$

---

**Input:** a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , accuracy  $0 < \varepsilon < 1$ ,  $p \geq 2$ .

Let  $\mathbf{v}_i^{(1)} = n/m$  for all  $i \in [m]$  and  $T = \Theta\left(\frac{\log(m/n)}{\varepsilon}\right)$ .

**for**  $k = 1, \dots, T$  **do**

Let  $\tilde{\sigma}$  be a  $\frac{\varepsilon}{4}$ -approximation of  $\sigma_i\left(\left(\mathbf{V}^{(k)}\right)^{1/2-1/p}\mathbf{A}\right)$ .

$\mathbf{v}_i^{(k+1)} = \tilde{\sigma}_i$ .

**end**

**Output:**  $\frac{1}{T} \sum_{k=1}^T \mathbf{v}^{(k)}$ .

---

**Theorem 5.3.4.** *For any matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , accuracy  $0 < \varepsilon < 1$  and  $p \geq 2$ , the algorithm 6 outputs an  $\varepsilon$ -approximate  $\ell_p$  Lewis weight in time  $\tilde{O}(\mathcal{T}\varepsilon^{-3})$  where  $\mathcal{T}$  be the time needed to apply  $(\mathbf{A}^\top \mathbf{D} \mathbf{A})^{-1}$  to a vector for any diagonal matrix  $\mathbf{D}$ .*

Since we know how to solve linear systems in input sparsity time (Chapter 3), this theorem gives an input sparsity time algorithm for approximating Lewis weights. This algorithm is based on the following convexity result.

**Lemma 5.3.5.** For any matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $i \in [m]$  and  $p \geq 2$ , the function

$$\phi_i(\mathbf{w}) = \log(\mathbf{a}_i^\top (\mathbf{A}^\top \mathbf{W}^{1-2/p} \mathbf{A})^{-1} \mathbf{a}_i \mathbf{w}_i^{-2/p})$$

is convex.

*Proof.* For any vector  $\mathbf{h} \in \mathbb{R}^m$ , we define  $\phi_i(t) = \phi_i(\mathbf{w} + t\mathbf{h})$ . Let  $\mathbf{M} = \left(\mathbf{A}^\top \mathbf{W}^{1-2/p} \mathbf{A}\right)^{-1}$ , we have that

$$\begin{aligned} \left. \frac{d}{dt} \right|_{t=0} \phi_i(t) &= -\frac{2 \mathbf{h}_i}{p \mathbf{w}_i} - \frac{(1-2/p) \mathbf{a}_i^\top \mathbf{M} \mathbf{A}^\top \mathbf{H} \mathbf{W}^{-2/p} \mathbf{A} \mathbf{M} \mathbf{a}_i}{\mathbf{a}_i^\top \mathbf{M} \mathbf{a}_i} \\ \left. \frac{d^2}{dt^2} \right|_{t=0} \phi_i(t) &= \frac{2 \mathbf{h}_i^2}{p \mathbf{w}_i^2} - \frac{(1-2/p)^2 \left(\mathbf{a}_i^\top \mathbf{M} \mathbf{A}^\top \mathbf{H} \mathbf{W}^{-2/p} \mathbf{A} \mathbf{M} \mathbf{a}_i\right)^2}{\left(\mathbf{a}_i^\top \mathbf{M} \mathbf{a}_i\right)^2} \\ &\quad - \frac{(-2/p)(1-2/p) \mathbf{a}_i^\top \mathbf{M} \mathbf{A}^\top \mathbf{H}^2 \mathbf{W}^{-1-2/p} \mathbf{A} \mathbf{M} \mathbf{a}_i}{\mathbf{a}_i^\top \mathbf{M} \mathbf{a}_i} \\ &\quad + \frac{2(1-2/p)^2 \mathbf{a}_i^\top \mathbf{M} \mathbf{A}^\top \mathbf{H} \mathbf{W}^{-2/p} \mathbf{A} \mathbf{M} \mathbf{A}^\top \mathbf{H} \mathbf{W}^{-2/p} \mathbf{A} \mathbf{M} \mathbf{a}_i}{\mathbf{a}_i^\top \mathbf{M} \mathbf{a}_i}. \end{aligned}$$

By Cauchy-Schwarz inequality, we have that

$$\left(\mathbf{a}_i^\top \mathbf{M} \mathbf{A}^\top \mathbf{H} \mathbf{W}^{-2/p} \mathbf{A} \mathbf{M} \mathbf{a}_i\right)^2 \leq \left(\mathbf{a}_i^\top \mathbf{M} \mathbf{A}^\top \mathbf{H} \mathbf{W}^{-2/p} \mathbf{A} \mathbf{M} \mathbf{A}^\top \mathbf{H} \mathbf{W}^{-2/p} \mathbf{A} \mathbf{M} \mathbf{a}_i\right) \left(\mathbf{a}_i^\top \mathbf{M} \mathbf{a}_i\right)$$

and hence

$$\begin{aligned} \left. \frac{d^2}{dt^2} \right|_{t=0} \phi_i(t) &\geq \frac{2}{p} \left( \frac{\mathbf{h}_i^2}{\mathbf{w}_i^2} + \frac{(1-2/p) \mathbf{a}_i^\top \mathbf{M} \mathbf{A}^\top \mathbf{H}^2 \mathbf{W}^{-1-2/p} \mathbf{A} \mathbf{M} \mathbf{a}_i}{\mathbf{a}_i^\top \mathbf{M} \mathbf{a}_i} \right) \\ &\quad + \frac{(1-2/p)^2 \mathbf{a}_i^\top \mathbf{M} \mathbf{A}^\top \mathbf{H} \mathbf{W}^{-2/p} \mathbf{A} \mathbf{M} \mathbf{A}^\top \mathbf{H} \mathbf{W}^{-2/p} \mathbf{A} \mathbf{M} \mathbf{a}_i}{\mathbf{a}_i^\top \mathbf{M} \mathbf{a}_i} \\ &\geq 0. \end{aligned}$$

Since this holds for any any vector  $\mathbf{h} \in \mathbb{R}^m$ , this shows that  $\phi_i$  is convex.  $\square$

Now, we prove Theorem 5.3.4.

*Proof of Theorem 5.3.4.* Let  $\bar{\mathbf{v}} = \frac{1}{T} \sum_{k=1}^T \mathbf{v}_i^{(k)}$ . Let  $\mathbf{d}_i^{(k)} = \tilde{\sigma}_i / \sigma_i \left( \left( \mathbf{V}^{(k)} \right)^{1/2-1/p} \mathbf{A} \right)$ . Then, we have that

$$\begin{aligned} \mathbf{v}_i^{(k+1)} &= \mathbf{d}_i^{(k)} \sigma_i \left( \left( \mathbf{V}^{(k)} \right)^{1/2-1/p} \mathbf{A} \right) \\ &= \mathbf{d}_i^{(k)} \left( \mathbf{v}_i^{(k)} \right)^{1-2/p} \mathbf{a}_i^\top \left( \mathbf{A}^\top \left( \mathbf{V}^{(k)} \right)^{1-2/p} \mathbf{A} \right)^{-1} \mathbf{a}_i. \end{aligned}$$

Hence, we have

$$\log \mathbf{v}_i^{(k+1)} = \log \mathbf{v}_i^{(k)} + \log \left( \mathbf{a}_i^\top \left( \mathbf{A}^\top \left( \mathbf{V}^{(k)} \right)^{1-2/p} \mathbf{A} \right)^{-1} \mathbf{a}_i \left( \mathbf{v}_i^{(k)} \right)^{-2/p} \right) + \log(\mathbf{d}_i^{(k)}).$$

Using  $\mathbf{v}_i^{(1)} = \frac{n}{m}$  and  $\mathbf{v}_i^{(T+1)} \leq 1$ , we have that

$$\begin{aligned} 0 &\geq \log \mathbf{v}_i^{(T+1)} \\ &= \log \mathbf{v}_i^{(1)} + \sum_{k=1}^T \log \left( \mathbf{a}_i^\top \left( \mathbf{A}^\top \left( \mathbf{V}^{(k)} \right)^{1-2/p} \mathbf{A} \right)^{-1} \mathbf{a}_i \left( \mathbf{v}_i^{(k)} \right)^{-2/p} \right) + \sum_{k=1}^T \log(\mathbf{d}_i^{(k)}). \end{aligned}$$

By Lemma 5.3.5, we have that

$$\log(\mathbf{a}_i^\top (\mathbf{A}^\top \bar{\mathbf{V}}^{1-2/p} \mathbf{A})^{-1} \mathbf{a}_i \bar{\mathbf{v}}_i^{-2/p}) \leq \frac{1}{T} \sum_{k=1}^T \log \left( \mathbf{a}_i^\top \left( \mathbf{A}^\top \left( \mathbf{V}^{(k)} \right)^{1-2/p} \mathbf{A} \right)^{-1} \mathbf{a}_i \left( \mathbf{v}_i^{(k)} \right)^{-2/p} \right).$$

Hence, we have that

$$\begin{aligned} \log \left( \frac{\sigma_i(\bar{\mathbf{V}}^{1/2-1/p} \mathbf{A})}{\bar{\mathbf{v}}_i} \right) &= \log(\mathbf{a}_i^\top (\mathbf{A}^\top \bar{\mathbf{V}}^{1-2/p} \mathbf{A})^{-1} \mathbf{a}_i \bar{\mathbf{v}}_i^{-2/p}) \\ &\leq \frac{\log \left( \frac{m}{n} \right)}{T} - \frac{1}{T} \sum_{i=1}^T \log(\mathbf{d}_i^{(k)}) \\ &\leq \frac{\log \left( \frac{m}{n} \right)}{T} + \frac{1}{2} \varepsilon. \end{aligned}$$

By setting  $T$  large enough, we have that

$$(1 - \varepsilon) \sigma_i(\bar{\mathbf{V}}^{1/2-1/p} \mathbf{A}) \leq \bar{\mathbf{v}}_i \leq (1 + \varepsilon) \sigma_i(\bar{\mathbf{V}}^{1/2-1/p} \mathbf{A}).$$

□

## ■ 5.4 Approximate John Ellipsoids

In this section, we discuss some geometric implications of  $\ell_p$  Lewis weight. [270] shows that the  $\ell_p$  Lewis weight is given by the following optimization problem.

**Lemma 5.4.1.** *Given a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Let  $\mathbf{M}$  be the maximizer of the following optimization problem*

$$\max_{\mathbf{M}} \log \det \mathbf{M} \text{ subject to } \sum_{i=1}^m \left( \mathbf{a}_i^\top \mathbf{M} \mathbf{a}_i \right)^{p/2} \leq n \text{ and } \mathbf{M} \succeq \mathbf{0}. \quad (5.2)$$

Then, the  $\ell_p$  Lewis weight of  $\mathbf{A}$  is given by  $\mathbf{w}_i = \left( \mathbf{a}_i^\top \mathbf{M} \mathbf{a}_i \right)^{p/2}$ . Furthermore, we have that  $\mathbf{M}^{-1} = \sum_i \mathbf{w}_i^{1-2/p} \mathbf{a}_i \mathbf{a}_i^\top$ .

To understand this optimization problem, we define the  $\ell_p$  Lewis Ellipsoids and John ellipsoids as follows:

**Definition 5.4.2.** For any matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , we define the  $\ell_p$  Lewis Ellipsoid of  $\mathbf{A}$  by  $E = \{ \mathbf{x} \in \mathbb{R}^n \text{ such that } \mathbf{x}^\top \mathbf{A}^\top \mathbf{W}^{1-2/p} \mathbf{A} \mathbf{x} \leq 1 \}$  where  $\mathbf{w}$  is the  $\ell_p$  Lewis weight of  $\mathbf{A}$ . For any convex set  $K$ , we define the John ellipsoid of  $K$  is the maximum volume ellipsoid contained inside  $K$

We note that  $\sum_{i=1}^m \left( \mathbf{a}_i^\top \mathbf{M} \mathbf{a}_i \right)^{p/2}$  is proportional to the average value of  $\| \mathbf{A} \mathbf{x} \|_p^p$  inside the ellipsoid  $\{ \mathbf{x}^\top \mathbf{M}^{-1} \mathbf{x} \leq 1 \}$ . Therefore, the  $\ell_p$  Lewis Ellipsoids is the maximum volume ellipsoid  $E$  such that the average value of  $\| \mathbf{A} \mathbf{x} \|_p^p$  inside that ellipsoid is less than certain value. In particular, the  $\ell_\infty$  Lewis

Ellipsoid of  $\mathbf{A}$  exactly corresponds to the John ellipsoid of the symmetric polytope  $\{\|\mathbf{A}\mathbf{x}\|_\infty \leq 1\}$ . Therefore, our algorithm for computing Lewis weights gives an efficient algorithm for approximating John ellipsoids for symmetric polytopes.

**Lemma 5.4.3.** *Given a symmetric polytope  $P = \{\mathbf{x} \text{ such that } \|\mathbf{A}\mathbf{x}\|_\infty \leq 1\}$ . Then, we can find  $\mathbf{v}$  such that*

$$\sigma_i(\mathbf{V}^{1/2}\mathbf{A}) \leq \mathbf{v}_i \leq (1 + \varepsilon)\sigma_i(\mathbf{V}^{1/2}\mathbf{A})$$

in time  $\tilde{O}(\mathcal{T}\varepsilon^{-3})$  where  $\mathcal{T}$  be the time needed to apply  $(\mathbf{A}^\top \mathbf{D} \mathbf{A})^{-1}$  to a vector for any diagonal matrix  $\mathbf{D}$ . Furthermore, the ellipsoid  $E = \{\mathbf{x} \in \mathbb{R}^n \text{ such that } \mathbf{x}^\top \mathbf{A}^\top \mathbf{V} \mathbf{A} \mathbf{x} \leq 1\}$  satisfies the rounding condition

$$E \subset P \subset \sqrt{(1 + \varepsilon)n}E. \quad (5.3)$$

*Proof.* By Lemma 5.3.4, we can find  $\mathbf{v}$  such that  $(1 - \frac{\varepsilon}{3})\sigma_i(\mathbf{V}^{1/2}\mathbf{A}) \leq \mathbf{v}_i \leq (1 + \frac{\varepsilon}{3})\sigma_i(\mathbf{V}^{1/2}\mathbf{A})$  in time  $\tilde{O}(\mathcal{T}\varepsilon^{-3})$ . By scaling  $\mathbf{v}$  up  $1/(1 - \varepsilon/3)$ , we get the first guarantee.

By  $\sigma_i(\mathbf{V}^{1/2}\mathbf{A}) \leq \mathbf{v}_i$ , we know that for all  $i$

$$\mathbf{a}_i^\top (\mathbf{A}^\top \mathbf{V} \mathbf{A})^{-1} \mathbf{a}_i \leq 1.$$

Since  $(\mathbf{a}_i^\top \mathbf{x})^2 \leq (\mathbf{a}_i^\top (\mathbf{A}^\top \mathbf{V} \mathbf{A})^{-1} \mathbf{a}_i) (\mathbf{x}^\top \mathbf{A}^\top \mathbf{V} \mathbf{A} \mathbf{x})$ , this shows that  $E \subset P$ .

On the other hand, for any  $\mathbf{x} \in P$ , we have that

$$\begin{aligned} \mathbf{x}^\top \mathbf{A}^\top \mathbf{V} \mathbf{A} \mathbf{x} &\leq \sum_{i=1}^m \mathbf{v}_i \|\mathbf{A}\mathbf{x}\|_\infty^2 \\ &\leq \sum_{i=1}^m (1 + \varepsilon)\sigma_i(\mathbf{V}^{1/2}\mathbf{A}) \\ &= (1 + \varepsilon)n. \end{aligned}$$

Therefore, we have that  $P \subset \sqrt{(1 + \varepsilon)n}E$ . □

*Remark 5.4.4.* We call the ellipsoid found in this theorem is an approximate John ellipsoid because it almost satisfies the optimality condition of John ellipsoid, i.e.  $\mathbf{v}_i = \sigma_i(\mathbf{V}^{1/2}\mathbf{A})$ . Although this does not show that the ellipsoid is actually close to the true John ellipsoid, this has similar rounding guarantee (5.3) as John ellipsoid and hence satisfies the requirement for most of the applications such as volume computation [262] and solving linear systems (next chapter).

## Faster Linear Programming

### ■ 6.1 Introduction

In this chapter, we present a new algorithm for the following linear programs

$$\begin{aligned} \min \quad & \mathbf{c}^T \mathbf{x}. \\ \mathbf{x} \in \mathbb{R}^m \quad & : \mathbf{A}^T \mathbf{x} = \mathbf{b} \\ \forall i \in [m] \quad & : l_i \leq x_i \leq u_i \end{aligned} \tag{6.1}$$

Our algorithm is based on interior point methods, a framework for solving linear programming. Roughly speaking, interior point methods approximate the domain by a sequence of ellipsoids and the number of iterations depends on how good are these approximations. In Chapter 5, we showed a family of ellipsoids that is easy to compute and gives a good approximations for arbitrary symmetric polytopes. Using these ellipsoids, we develop an interior point method that takes  $\tilde{O}(\sqrt{\text{rank}(\mathbf{A})}L)^1$  iterations (Theorem 6.8.1). Furthermore, we show how to achieve this convergence rate while only solving  $\tilde{O}(1)$  linear systems and performing additional  $\tilde{O}(\text{nnz}(\mathbf{A}))$  work in each iteration. This is the first polynomial improvement on the running time of solving linear programming over 25 years.

Applying our techniques to the linear programming formulation of the maximum flow problem we achieve a running time of  $\tilde{O}(|E|\sqrt{|V|}\log(U))$  for solving the problem on a graph with  $|E|$  edges,  $|V|$  vertices, and capacity ratio  $U$  (Theorem 6.9.3). This improves upon the previous fastest running time of  $\tilde{O}(|E|\min\{|E|^{1/2}, |V|^{2/3}\}\log(U))$  achieved over 15 years ago by Goldberg and Rao [104]<sup>2</sup> and the previous fastest running times for solving dense directed unit capacity graphs of  $O(|E|\min\{|E|^{1/2}, |V|^{2/3}\})$  achieved by Even and Tarjan [85] over 35 years ago and of  $\tilde{O}(|E|^{10/7})$  achieved recently by Mařdry [183].

#### ■ 6.1.1 Organization

The rest of this chapter is structured as follows. In Section 6.3, we provide some background information on linear programming and interior point methods and show how to obtain the “first”  $\tilde{O}(\sqrt{\text{rank}(\mathbf{A})}L)$  iterations algorithm for the standard dual linear program (6.2). In Section 6.4, we explain why this algorithm is not sufficient for the maximum flow problem. In Section 6.5 we introduce our framework to use our alternative barrier function more efficiently and in Section 6.6 we present the key lemmas used to analyze progress along paths and in Section 6.7 we introduce the weight function we use to find paths. In Section 6.8 provide our linear programming algorithm for the linear program (6.1). In Section 6.9 we use these results to achieve our desired running times for the maximum flow

<sup>1</sup>Throughout this thesis  $L$  denotes the standard “bit complexity” of the linear program, a quantity less than the number of bits needed to represent (6.2). For integral  $\mathbf{A}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$  this quantity is often defined as  $L \stackrel{\text{def}}{=} \log(1 + d_{\max}) + \log(1 + \max\{\|\mathbf{c}\|_\infty, \|\mathbf{b}\|_\infty\})$  where  $d_{\max}$  is the largest absolute value of the determinant of a square sub-matrix of  $\mathbf{A}$  [136].

<sup>2</sup>In this thesis we are primarily concerned with “weakly” polynomial time algorithms. The current fastest “strongly” polynomial running time for solving this problem is  $O(nm)$  [218].

problem and its generalizations.

## ■ 6.2 Previous Work

Linear programming and maximum flow are extremely well studied problems with long histories. Here we focus on recent history that particularly influenced our work.

### ■ 6.2.1 Convergence Rate of Interior Point Methods

Given a matrix,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , and vectors,  $\mathbf{b} \in \mathbb{R}^m$  and  $\mathbf{c} \in \mathbb{R}^n$ , solving the linear program<sup>3</sup>

$$\min_{\mathbf{x} \in \mathbb{R}^n : \mathbf{Ax} \geq \mathbf{b}} \mathbf{c}^T \mathbf{x} \quad (6.2)$$

is a core algorithmic task for both the theory and practice of computer science.

In 1984, Karmarkar [136] provided the first proof of an interior point method running in polynomial time. This method required  $O(mL)$  iterations where the running time of each iteration was dominated by the time needed to solve a linear system of the form  $(\mathbf{A}^T \mathbf{D} \mathbf{A}) \mathbf{x} = \mathbf{y}$  for some positive diagonal matrix  $\mathbf{D} \in \mathbb{R}^{m \times m}$  and some  $\mathbf{y} \in \mathbb{R}^n$ . Using low rank matrix updates and preconditioning Karmarkar achieved a running time of  $O(m^{3.5}L)$ .

In 1988, Renegar improved the convergence rate of interior point methods to  $O(\sqrt{m}L)$  iterations [227]. He presented a path following method where the iterations had comparable complexity to Karmarkar's method. Using a combination of techniques involving low rank updates, preconditioning and fast matrix multiplication, the amortized complexity of each iteration was improved [254, 107, 211] yielding the current state of the art running time of  $\tilde{O}(m^{1.5}nL)$  [257].

In 1989, Vaidya [258] proposed two barrier functions which were shown to yield  $O((m \text{rank}(\mathbf{A}))^{1/4}L)$  and  $O(\text{rank}(\mathbf{A})L)$  iteration linear programming algorithms [255, 258, 253]. Unfortunately each iteration of these methods required explicit computation of the projection matrix  $\mathbf{D}^{1/2} \mathbf{A} (\mathbf{A}^T \mathbf{D} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{D}^{1/2}$  for a positive diagonal matrix  $\mathbf{D} \in \mathbb{R}^{m \times m}$ . This was improved by Anstreicher [15] who showed it sufficed to compute the diagonal of this projection matrix. Unfortunately both methods do not yield faster running times unless  $m \gg n$ .

In 1994 Nesterov and Nemirovski [211] showed that path-following methods can in principle be applied to any solve convex optimization problem, given a suitable barrier function. The number of iterations of their method depends on a square root of a quantity known as the *self-concordance* of the barrier. They showed that for any convex set in  $\mathbb{R}^n$ , there exists a barrier function, called the *universal barrier* function, with self-concordance  $O(n)$  and therefore in principle any convex optimization problem with  $n$  variables can be solved in  $O(\sqrt{n}L)$  iterations. However, this result is generally considered to be only of theoretical interest as the universal barrier function is defined as the volume of certain polytopes, a problem which in full-generality is NP-hard and can only be computed approximately in  $O(n^c)$  for some large constant  $c$  [180].

These results suggested that you can solve linear programs closer to the  $\tilde{O}(\sqrt{\text{rank}(\mathbf{A})}L)$  bound achieved by the universal barrier only if you pay more in each iteration. In this chapter we show that this is not the case and up to polylogarithmic factors we achieve the convergence rate of the universal barrier function while only having iterations of cost comparable to that of Renegar's algorithm.

---

<sup>3</sup>This expression is the dual of a linear program written in standard form. It is well known that all linear programs can be written as (6.2). Note that this notation of  $m$  and  $n$  differs from that in some papers. Here  $m$  denotes the number of constraints and  $n$  denotes the number of variables. To avoid confusion we state many of our results in terms of  $\sqrt{\text{rank}(\mathbf{A})}$  instead of  $\sqrt{n}$ .

Year	Author	Number of Iterations	Nature of iterations
1984	Karmarkar [136]	$O(mL)$	Linear system solve
1986	Renegar [227]	$O(\sqrt{m}L)$	Linear system solve
1989	Vaidya [253]	$O((m\text{rank}(\mathbf{A}))^{1/4}L)$	Expensive linear algebra
1994	Nesterov and Nemirovskii [211]	$O(\sqrt{\text{rank}(\mathbf{A})}L)$	Volume computation
-	This chapter	$\tilde{O}(\sqrt{\text{rank}(\mathbf{A})}L)$	$\tilde{O}(1)$ Linear system solves

**Figure 6-1:** Interior Point Iteration Improvements

## ■ 6.2.2 Recent Breakthroughs on the Maximum Flow Problem

A beautiful line of work on solving the maximum flow on undirected graphs began with a result of Benzur and Karger in which they showed how to reduce approximately computing minimum s-t cuts on arbitrary undirected graphs to the same problem on sparse graphs, i.e. those with  $|E| = \tilde{O}(|V|)$  [33]. Pushing these ideas further Karger showed how to perform a similar reduction for the maximum flow problem [134] and Karger and Levine showed how to compute the exact maximum flow in an unweighted undirected graph with maximum flow value  $F$  in  $\tilde{O}(|E| + |V|F)$  time [135].

In a breakthrough result of Spielman and Teng [243] they showed that a particular class of linear systems, Laplacians, can be solved in nearly linear time. Leveraging these fast Laplacian system solvers Christiano, Kelner, Mądry, and Spielman [51] showed how to compute  $(1 - \varepsilon)$  approximations to the maximum flow problem on undirected graphs in  $\tilde{O}(|E| \cdot |V|^{1/3} \varepsilon^{-11/3})$ . Later Lee, Rao and Srivastava [161] improved the running time to  $\tilde{O}(|E| \cdot |V|^{1/3} \varepsilon^{-2/3})$  for unweighted undirected graphs. This line of work culminated in recent results of Sherman [233] and Kelner, Lee, Orecchia, and Sidford (Chapter 12) who showed how to solve the problem in  $O(|E|^{1+o(1)} \varepsilon^{-2})$ , using congestion-approximators, oblivious routings, and techniques developed by Mądry [182].

In 2008, Daitch and Spielman [63] showed that, by careful application of interior point techniques, fast Laplacian system solvers [243], and a novel method for solving M-matrices, they could match (up to polylogarithmic factors) the running time of Goldberg Rao and achieve a running time of  $\tilde{O}(|E|^{3/2} \log^2(U))$  not just for maximum flow but also for the minimum cost flow and lossy generalized minimum cost flow problems (see Fig 6-2).

Very recently Mądry [183] achieved an astounding running time of  $\tilde{O}(|E|^{10/7})$  for solving the maximum flow problem on un-capacitated directed graphs by a novel application and modification of interior point methods. This shattered numerous barriers providing the first general improvement over the running time of  $O(|E| \min\{|E|^{1/2}, |V|^{2/3}\})$  for solving unit capacity graphs proven over 35 years ago by Even and Tarjan [85] in 1975.

While our algorithm for solving the maximum flow problem is new, we make extensive use of these breakthroughs. We use sampling techniques first discovered in the context of graph sparsification (previous two chapters) to re-weight the graph so that we make progress at a rate commensurate with the number of vertices and not the number of edges. We use fast Laplacian system solvers as in [51, 161] to make the cost of interior point iterations cheap when applied to the linear program formulations analyzed by Daitch and Spielman [63]. Furthermore, as in Mądry [183] we use weights to change the central path (albeit for a slightly different purpose). We believe this further emphasizes the power of these tools as general purpose techniques for algorithm design.

Year	Author	Running Time
1972	Edmonds, Karp [84]	$\tilde{O}( E ^2 \log(U))$
1984	Tardos [247]	$O( E ^4)$
1984	Orlin [217]	$\tilde{O}( E ^3)$
1986	Galil, Tardos [98]	$\tilde{O}( E  V ^2)$
1987	Goldberg, Tarjan [106]	$\tilde{O}( E  V  \log(U))$
1988	Orlin [215]	$\tilde{O}( E ^2)$
2008	Daitch, Spielman [63]	$\tilde{O}( E ^{3/2} \log^2(U))$
-	This chapter	$\tilde{O}( E  \sqrt{ V } \log^2(U))$

Figure 6-2: Minimum Cost Flow Running Times Improvements

## ■ 6.3 Simple Case

In this section, we introduce the notation of  $\nu$ -self-concordant barrier functions and explain how such barrier functions can be used to develop an  $O(\sqrt{\nu}L)$  iterations linear programming algorithm. Then, we construct an  $\tilde{O}(n)$ -self-concordant barrier function which can be computed in polynomial time. This barrier function will only be used to give the intuition on how we are getting our main algorithm. The reader well versed in path following methods can likely skip to subsection 6.3.5 and to the more curious reader we encourage them to consider some of the many wonderful expositions on this subject [209, 273, 107] for further reading.

### ■ 6.3.1 The Setup

Given a matrix,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , and vectors,  $\mathbf{b} \in \mathbb{R}^m$  and  $\mathbf{c} \in \mathbb{R}^n$ , the goal of this section is to solve the following linear program

$$\min_{\mathbf{x} \in \mathbb{R}^n : \mathbf{Ax} \geq \mathbf{b}} \mathbf{c}^T \mathbf{x} \quad (6.3)$$

This is called the dual of the *standard form* of a linear program and all linear programs can be expressed in this form. We call a vector  $\mathbf{x} \in \mathbb{R}^m$  *feasible* if  $\mathbf{Ax} \geq \mathbf{b}$ , we call  $\mathbf{c}^T \mathbf{x}$  the *cost* of such a vector. therefore our goal is to either compute a minimum cost feasible vector or determine that none exists.

We assume that  $\mathbf{A}$  is full rank, i.e.  $\text{rank}(\mathbf{A}) = n$ , and that  $m \geq n$ . Nevertheless, we still write many of our results using  $\text{rank}(\mathbf{A})$  rather than  $n$  for two reasons. First, this notation makes clear that  $\text{rank}(\mathbf{A})$  is referring to the smaller of the two quantities  $m$  and  $n$ . Second, if  $\text{rank}(\mathbf{A}) < n$ , then we can reduce the number of variables to  $\text{rank}(\mathbf{A})$  by a change of basis. Hence, we only need to solve linear programs in the full rank version.

### ■ 6.3.2 Central Path

Interior point methods solve (6.3) by maintaining a point  $\mathbf{x}$  that is in the *interior* of the feasible region, i.e.  $\mathbf{x} \in \Omega$  where

$$\Omega \stackrel{\text{def}}{=} \{\mathbf{x} \in \mathbb{R}^n : \mathbf{Ax} > \mathbf{b}\}.$$

These methods attempt to iteratively decrease the cost of  $\mathbf{x}$  while maintaining strict feasibility. This is often done by considering some measurement of the distance to feasibility such as  $\mathbf{s} \stackrel{\text{def}}{=} \mathbf{Ax} - \mathbf{b}$ , called the *slacks*, and creating some penalty for these distances approaching 0. By carefully balancing penalties for small  $\mathbf{s}(\mathbf{x})$  and penalties for large  $\mathbf{c}^T \mathbf{x}$ , these methods eventually compute a point close



enough to the optimum solution that it can be computed exactly.

Path following methods fix ratios between the the penalty for large  $\mathbf{c}^T \mathbf{x}$  and the penalty for small  $\mathbf{s}$  and alternate between steps of optimizing with respect to this ratio and changing the ratio. These methods typically encode the penalties through a *barrier function*  $\phi : \Omega \rightarrow \mathbb{R}$  such that  $\phi(\mathbf{x}) \rightarrow \infty$  as  $\mathbf{s}(\mathbf{x})_i \rightarrow 0$  for any  $i \in [m]$  and they encode the ratio through some parameter  $t > 0$ . Formally, they attempt to solve optimization problems of the following form for increasing values of  $t$

$$\min_{\mathbf{x} \in \mathbb{R}^m} f_t(\mathbf{x}) \quad \text{where} \quad f_t(\mathbf{x}) \stackrel{\text{def}}{=} t \cdot \mathbf{c}^T \mathbf{x} + \phi(\mathbf{x}) \quad (6.4)$$

Since  $\phi(\mathbf{x}) \rightarrow \infty$  as  $\mathbf{s}(\mathbf{x})_i \rightarrow 0$ , the minimizer of  $f_t(\mathbf{x})$ , denoted  $\mathbf{x}^*(t)$ , is in  $\Omega$  for all  $t$ . As  $t$  increases the effect of the cost vector on  $\mathbf{x}^*(t)$  increases and the distance from the boundary of the feasible region as measured by  $\mathbf{s}(\mathbf{x})$  decreases. One can think of the points  $\{\mathbf{x}^*(t) \mid t > 0\}$  as a path in  $\mathbb{R}^n$ , called *the central path*, where  $\mathbf{x}^*(t)$  approaches a solution to (6.3) as  $t \rightarrow \infty$ . A standard choice of barrier is the *standard log barrier*,  $\phi(\mathbf{x}) \stackrel{\text{def}}{=} -\sum_{i=1}^m \log(\mathbf{s}(\mathbf{x})_i)$  and for this choice of barrier we refer to  $\{\mathbf{x}^*(t) \mid t > 0\}$  as the *standard central path*.

Path following methods typically follow the following framework:

- (1) *Compute Initial Point*: Compute an approximation  $\mathbf{x}^*(t)$  for some  $t$ .
- (2) *Follow the central path*: Repeatedly increase  $t$  and compute an approximation to  $\mathbf{x}^*(t)$ .
- (3) *Round to optimal solution*: Use the approximation to  $\mathbf{x}^*(t)$  to compute the solution to (6.3).

Usually, steps (1) and (3) are not the bottleneck of interior point methods (including our algorithm) and can be carried out by standard interior point techniques. In this chapter, we ignore the details for step (1) and (3) and refer the interested reader to the paper version of this chapter [165, 166]. In the following subsection, we provide a simple technique for performing (2) and serves as the foundation for the algorithms considered in the remainder of this section.

### ■ 6.3.3 Following the Path

Typically, path following methods compute an approximation to  $\mathbf{x}^*(t)$  using Newton's method. While for an arbitrary current point  $\mathbf{x} \in \Omega$  and  $t > 0$ , the function  $f_t(\mathbf{x})$  can be ill-behaved, in a region near  $\mathbf{x}^*(t)$ , the Hessian  $\nabla^2 f_t(\mathbf{x})$  changes fairly slowly. More precisely, if one considers the second order approximation of  $f_t(\mathbf{z})$  around some point  $\mathbf{x} \in \Omega$  "close enough" to  $\mathbf{x}^*(t)$

$$f_t(\mathbf{z}) \approx f_t(\mathbf{x}) + \langle \nabla f_t(\mathbf{x}), \mathbf{z} - \mathbf{x} \rangle + \frac{1}{2} (\mathbf{z} - \mathbf{x})^T (\nabla^2 f_t(\mathbf{x})) (\mathbf{z} - \mathbf{x}),$$

and applies one step of Newton's method, i.e. minimizes this quadratic approximation to compute

$$\mathbf{x}^{(new)} := \mathbf{x} - (\nabla^2 f_t(\mathbf{x}))^{-1} \nabla f_t(\mathbf{x})$$

for  $\mathbf{s} \stackrel{\text{def}}{=} \mathbf{s}(\mathbf{x})$ , then this procedure rapidly converges to  $\mathbf{x}^*(t)$ .

To quantify this, we measure *centrality*, i.e. how close the current point  $\mathbf{x} \in \Omega$  is to  $\mathbf{x}^*(t)$ , by the size of this *Newton step* in the Hessian induced norm. For  $\mathbf{x} \in \Omega$  and Newton step  $\mathbf{h}_t(\mathbf{x}) \stackrel{\text{def}}{=} (\nabla^2 f_t(\mathbf{x}))^{-1} \nabla f_t(\mathbf{x})$ , we denote centrality by  $\delta_t(\mathbf{x}) \stackrel{\text{def}}{=} \|\mathbf{h}_t(\mathbf{x})\|_{\nabla^2 f_t(\mathbf{x})}$ .

For the standard log barrier function,  $\phi(\mathbf{x}) \stackrel{\text{def}}{=} -\sum_{i \in [m]} \log(\mathbf{s}(\mathbf{x})_i)$ , a standard analysis of Newton's method shows that if  $\delta_t(\mathbf{x}) \leq \frac{1}{2}$ , then for  $\mathbf{x}^{(new)} := \mathbf{x} - \mathbf{h}_t(\mathbf{x})$ , we have  $\delta_t(\mathbf{x}^{(new)}) \leq 2\delta_t(\mathbf{x})^2$ . Furthermore, if  $\delta_t(\mathbf{x}) \leq \frac{1}{4}$ , it is not hard to show that for  $t' = t(1 + m^{-1/2}/4)$  we have  $\delta_{t'}(\mathbf{x}) \leq 1/2$ . Combining these facts yields that in  $O(\sqrt{m})$  iterations we can double  $t$  while maintaining a *nearly centered*  $\mathbf{x}$ ,

i.e.  $\delta_t(\mathbf{x})$  at most a constant. With some additional work, it can be shown that by maintaining a nearly centered  $\mathbf{x}$  and increasing  $t$  at most  $O(\sqrt{m}L)$  times one can compute a solution to (6.3) exactly. Therefore, this barrier function gives an  $O(\sqrt{m}L)$  iterations algorithm for solving (6.3).

### ■ 6.3.4 Self-Concordant Barrier Functions

In a seminal result of Nesterov and Nemirovski [211], they studied the sufficient conditions (called *self-concordance*) on  $\phi$  that allows us to follow the path efficiently. They showed that given the ability to construct a  $\nu$ -self-concordant barrier for a convex set, one can minimize linear functions over that convex set with a convergence rate of  $O(\sqrt{\nu})$ .

**Definition 6.3.1.** Given a open convex set  $K \subset \mathbb{R}^n$ . We call  $\phi$  is a  $\nu$ -self-concordant barrier function for  $K$  if the following conditions are satisfies

- $\phi$  is convex and thrice continuously differentiable,
- $\phi(\mathbf{x}_i) \rightarrow \infty$  for any sequence  $\mathbf{x}_i \in K$  converging to boundary of  $K$ .
- $|D^3\phi(\mathbf{x})[\mathbf{h}, \mathbf{h}, \mathbf{h}]| \leq 2|D^2\phi(\mathbf{x})[\mathbf{h}, \mathbf{h}]|^{3/2}$  for any  $\mathbf{x} \in K$ ,  $\mathbf{h} \in \mathbb{R}^n$ ,
- $|D\phi(\mathbf{x})[\mathbf{h}]| \leq \sqrt{\nu}|D^2\phi(\mathbf{x})[\mathbf{h}, \mathbf{h}]|^{1/2}$  for any  $\mathbf{x} \in K$ ,  $\mathbf{h} \in \mathbb{R}^n$ .

For example, the standard log barrier function is a  $m$ -self concordant barrier function for polytopes. Nesterov and Nemirovski further showed that, for any convex set in  $\mathbb{R}^n$ , there exists a barrier function, called the *universal barrier* function, with self-concordance  $O(n)$  [211]. However, the universal barrier function is defined via the volume of certain polytopes, which is very difficult to compute.

### ■ 6.3.5 A $\tilde{O}(n)$ -Self-Concordant Barrier Function

In this section, we showed that for any polytope, there is an easy-to-compute barrier function with self-concordance  $\tilde{O}(n)$ . This barrier function is based on the following intuition on interior point methods. In high level, interior point methods replace the key difficulty of linear programming, the hard constraints, by barrier functions and solve the linear programs by a sequence of linear systems induced by barrier functions. Since these linear system corresponds to minimizing linear functions over ellipsoids and these ellipsoids comes from the second order approximation of barrier functions, interior point methods, to a certain degree, use a sequence of ellipsoids to approximate polytopes. Self-concordance can be thought as a geometric condition that relates how good those ellipsoids approximate the domain. In particular, the following lemma shows that the second order approximation of the barrier function at the minimum point indeed has a geometric implication, it approximates the domain.

**Theorem 6.3.2** ([206, Thm 4.2.6]). *Given a  $\nu$ -self-concordant barrier function  $\phi$  for some convex set  $\Omega \subset \mathbb{R}^n$ . Let  $x_\phi$  be the minimizer of  $\phi$  and let the Dikin ellipse  $E = \{x \in \mathbb{R}^n : (x - x_\phi)^T \nabla^2 \phi(x_\phi)(x - x_\phi) \leq 1\}$ . Then, we have that*

$$E \subset \Omega \subset (\nu + 2\sqrt{\nu})E. \quad (6.5)$$

This shows that in order to construct a  $\tilde{O}(n)$ -self-concordant barrier function, we must able to find an ellipse that gives a  $\tilde{O}(n)$  approximation of the convex set  $\Omega$ . There are quite a few ellipsoids satisfying the condition (6.5). In contrast to other known ellipsoids that gives an approximation guarantee such as second moment, John ellipsoid is defined by a convex optimization problem and therefore can be computed in weakly polynomial time.

**Lemma 6.3.3** ([128]). *Given a convex  $\Omega \subset \mathbb{R}^n$ . We define the John ellipsoid  $J(\Omega)$  as the largest volume ellipsoid contained inside  $\Omega$ . Then, we have that*

$$J(\Omega) \subset K \subset nJ(\Omega).$$

There are different ways to write down the John ellipsoid as a convex problem. Our barrier function is motivated by the following formulation, called *D-optimal design*.

**Lemma 6.3.4** ([143]). *For any polytope  $\Omega = \{\mathbf{A}\mathbf{x} > \mathbf{b}\}$ , the John ellipsoid  $J(K)$  is given by  $\{\mathbf{y} \in \mathbb{R}^n : (\mathbf{y} - \mathbf{x})^T \mathbf{A}^T \mathbf{W} \mathbf{A} (\mathbf{y} - \mathbf{x}) \leq 1\}$  where  $(\mathbf{x}, \mathbf{w})$  is the saddle point of the following convex concave problem*

$$\min_{\mathbf{x}} \max_{\sum w_i = n, w_i \geq 0} \ln \det (\mathbf{A}^T \mathbf{S}^{-1} \mathbf{W} \mathbf{S}^{-1} \mathbf{A}) \quad (6.6)$$

where  $\mathbf{S}$  and  $\mathbf{W}$  are diagonal  $m \times m$  matrices with  $S_{ii} = a_i^T \mathbf{x} - b_i$  and  $W_{ii} = w_i$ .

Motivated by Theorem 6.3.2, we want to find a barrier function whose minimizer is exactly the center of John ellipsoid. From (6.6), we see that the function

$$\phi_{\infty}(\mathbf{x}) \stackrel{\text{def}}{=} \max_{\sum w_i = n, w_i \geq 0} \ln \det (\mathbf{A}^T \mathbf{S}^{-1} \mathbf{W} \mathbf{S}^{-1} \mathbf{A})$$

satisfies this requirement exactly.

However,  $\phi_{\infty}(\mathbf{x})$  is not even continuously differentiable. To see this, we define  $J(\Omega, x)$  be the maximum volume ellipsoid inside  $K$  and centered at  $x$ . We note that  $\phi_{\infty}(\mathbf{x}) = c \log(\text{vol}(J(\Omega, x)))$  for some universal constant  $c$ . Consider the convex set  $\Omega = [0, 1]$ . For any  $x < 1/2$ ,  $J(\Omega, x)$  is the “ellipsoid”  $[0, 2x]$  and in some sense this ellipsoid does not see the constraint  $x \leq 1$ . Therefore, the function  $\phi_{\infty}(x)$  is not continuously differentiable at  $1/2$ .

To make  $\phi_{\infty}(\mathbf{x})$  smooth, we employ the standard approach, adding a strongly concave term into the objective function  $\ln \det (\mathbf{A}^T \mathbf{S}^{-1} \mathbf{W} \mathbf{S}^{-1} \mathbf{A})$ . In general, if a smooth function  $f(x, y)$  is strongly concave in  $y$ , then  $\max_y f(x, y)$  is smooth. Since the domain of  $w$  is a simplex, a natural choice for the strongly concave term is the entropy function. Therefore, we consider the following function

$$\max_{\sum w_i = n, w_i \geq 0} \ln \det (\mathbf{A}^T \mathbf{S}^{-1} \mathbf{W} \mathbf{S}^{-1} \mathbf{A}) + \rho(\mathbf{w})$$

with  $\rho(\mathbf{w}) = -\sum_{i=1}^m w_i \ln w_i$ . By some direct calculations, one can show this is indeed a  $\tilde{O}(n)$  self-concordant barrier function. In fact, there are many other choices of  $\rho$  that give similar result. For example, we considered  $\rho(\mathbf{w}) = -\sum_{i=1}^m w_i^{1+1/\log(m/n)}$  in our original paper. In the appendix, we proved the following:

**Theorem 6.3.5.** *For any polytope  $\Omega = \{\mathbf{A}\mathbf{x} > \mathbf{b}\}$ , we define the barrier  $p(x)$  for  $\Omega$  as  $\max_{w_i \geq 0} f(x, w)$  where*

$$f(x, w) = \ln \det (\mathbf{A}^T \mathbf{S}^{-1} \mathbf{W} \mathbf{S}^{-1} \mathbf{A}) - \frac{n}{m} \sum_{i=1}^m w_i \ln w_i - \frac{n}{m} \sum_{i=1}^m \ln s_i \quad (6.7)$$

and  $\mathbf{S}$  and  $\mathbf{W}$  are diagonal  $m \times m$  matrices with  $S_{ii} = a_i^T \mathbf{x} - b_i$  and  $W_{ii} = w_i$ . Then,  $p$  is an  $O(n \ln^6(\frac{me}{n}))$ -self concordant barrier function (after rescaling). More precisely, we have

1. For any  $x \in \Omega^\circ$ ,  $(\nabla p(x))^T (\nabla^2 p(x))^{-1} (\nabla p(x)) \leq O(n)$ .
2. For any  $x \in \Omega^\circ$ ,  $h \in \mathbb{R}^n$ , we have  $|\nabla^3 p(x)[h, h, h]| \leq O(\ln^3(\frac{me}{n})) |\nabla^2 p(x)[h, h]|^{3/2}$ .

We remark that both the scaling  $n/m$  and the extra term  $\sum_{i=1}^m \ln s_i$  is not essential. In fact, one consider the  $\ell_p$  Lewis ellipsoids instead and use the following barrier function

$$\phi_p(x) \stackrel{\text{def}}{=} \min_{\mathbf{M}} -\log \det \mathbf{M} \text{ subject to } \sum_{i=1}^m \left( \frac{\mathbf{a}_i^\top \mathbf{M} \mathbf{a}_i}{s_i^2(x)} \right)^{p/2} \leq n \text{ and } \mathbf{M} \succeq \mathbf{0}.$$

One can show that  $\phi_{\log(m)}$  is an  $\tilde{O}(n)$ -self concordant barrier function. However, this barrier function is more difficult to analyze and hence we left it as an 15 pages calculus exercise.

## ■ 6.4 Difficulties

In this section, we discuss two issues arises from using the barrier function defined in the previous section and outline our approaches used in the rest of this chapter for fixing them.

### ■ 6.4.1 Avoiding Matrix Multiplication

The minimax formula involves the maximizer  $\mathbf{w}$  in (6.7), which we call weight function  $g(\mathbf{x})$ . Computing  $g(\mathbf{x})$  requires computing the diagonal of some projection matrices as in Vaidya and Anstreicher's work [257, 15] and is therefore not efficient enough for our purposes. Although we can compute  $g(\mathbf{x})$  approximately up to certain multiplicative coordinate-wise error using dimension reduction techniques (Section 6.7.1), this error is still too much for path following to handle the directly as multiplicatively changing weights can hurt our measures of progress too much.

To avoid the minimax formulation, we note that the barrier function  $p(\mathbf{x})$  has similar behaviors as  $-\sum \mathbf{w}_i \ln s_i(\mathbf{x})$  around  $\mathbf{x}$  if  $\mathbf{w}$  is close to the maximizer in (6.7) for that given  $\mathbf{x}$ . Therefore, rather than using the minimax formula, we maintain separate weights  $\mathbf{w}$  and current point  $\mathbf{x}$  and use the barrier

$$\phi(\mathbf{x}, \mathbf{w}) = - \sum_{i \in [m]} \mathbf{w}_i \ln s_i(\mathbf{x}).$$

We then maintain two invariants, (1)  $\mathbf{x}$  is centered, i.e.  $\mathbf{x}$  close to the minimum point of  $t \cdot \mathbf{c}^T \mathbf{x} + \phi(\mathbf{x}, \mathbf{w})$  and (2)  $\mathbf{w}$  close to  $g(\mathbf{x})$  multiplicatively.

We separate the problem of maintaining these invariants into two steps. First, we show that a Newton step for  $\mathbf{x}$  improves centrality without moving  $\mathbf{w}$  too far away from  $g(\mathbf{x})$ . Second, we show that given a multiplicative approximation to  $g(\mathbf{x})$  and bounds for how much  $g(\mathbf{x})$  may have changed, we can maintain the invariant that  $g(\mathbf{x})$  is close to  $\mathbf{w}$  multiplicatively without moving  $\mathbf{w}$  too much. We formulate this as a general two player game and prove that there is an efficient strategy to maintain our desired invariant. Combining these and standard techniques in path-following methods, we obtain an  $\tilde{O}(\sqrt{\text{rank}(\mathbf{A})L})$  iterations path-following algorithm for solving (6.2) where each iterations consists of  $\tilde{O}(1)$  linear system solves.

### ■ 6.4.2 Generalizing The Result

Unfortunately, this result is insufficient to produce faster algorithms for the maximum flow problem and its generalizations. Given an arbitrary minimum cost maximum flow instance there is a natural linear program that one can use to express the problem:

$$\begin{aligned} & \min && \mathbf{c}^T \mathbf{x} \\ & \mathbf{x} \in \mathbb{R}^m & : & \mathbf{A}^T \mathbf{x} = \mathbf{b} \\ & \forall i \in [m] & : & l_i \leq x_i \leq u_i \end{aligned} \tag{6.8}$$

where the variables  $x_i$  denote the flow on an edge, the  $l_i$  and  $u_i$  denote lower and upper bounds on how much flow we can put on the edge, and  $\mathbf{A}$  is the incidence matrix associated with the graph [63]. In this formulation,  $\text{rank}(\mathbf{A})$  is less than the number of vertices in the graph and using fast Laplacian system solvers [243, 147, 146, 141, 162, 170, 223] we can solve linear systems involving  $\mathbf{A}$  in time nearly linear in the number of edges in the graph. Thus, if we could solve (6.8) in time comparable to that we achieve for solving (6.2) this would immediately yield a  $\tilde{O}(m\sqrt{n}\log^{O(1)}(U))$  algorithm for the maximum flow problem. Unfortunately, it is not clear how to apply our previous results in this more general setting and naive attempts to write (6.8) in the form of (6.2) without increasing  $\text{rank}(\mathbf{A})$  fail.

Even more troubling, achieving a faster than  $\tilde{O}(\sqrt{m}L)$  iterations interior point method for solving general linear programs in this form would break a long-standing barrier for the convergence rate of interior point methods. To the best of the authors knowledge, there is no general purpose interior point method that achieves a convergence rate faster than the self-concordance of the best barrier of the feasible region. Furthermore, using lower bounds results of Nesterov and Nemirovski [211, Proposition 2.3.6], it is not hard to see that any general barrier for (6.8) must have self-concordance  $\Omega(m)$ .

After all, linear programs of the form (6.2) are easier because each step of Newton steps lies in some ellipsoids that is strictly contained inside the domain  $\mathbf{A}\mathbf{x} \geq \mathbf{b}$ . However, such restriction implies that it takes at least  $O(\sqrt{m})$  iterations to go from one point in the box  $\{l_i \leq x_i \leq u_i\}$  to another point and that is too slow. Therefore, we need to explicitly bound the size of the Newton step in both the  $\ell_\infty$  norm and a weighted  $\ell_2$  norm. Hence, we measure the progress by the size of the Newton step in a mixed norm of the form  $\|\cdot\| = \|\cdot\|_\infty + C_{\text{norm}}\|\cdot\|_{\mathbf{W}}$  to keep track of these two quantities simultaneously.

Measuring of Newton step size both with respect to the mixed norm helps to explain how our method outperforms the self-concordance of the best barrier for the space. Self-concordance is based on  $\ell_2$  analysis and the lower bounds for self-concordance are precisely the failure of the sphere to approximate a box. While ideally we would just perform optimization over the  $\ell_\infty$  box directly,  $\ell_\infty$  is ripe with degeneracies that makes this analysis difficult. Nevertheless, unconstrained minimization over a box is quite simple and by working with this mixed norm and choosing weights to improve the conditioning, we are taking advantage of the simplicity of minimizing  $\ell_\infty$  over most of the domain and only paying for the  $n$ -self-concordance of a barrier for the smaller subspace induce by the requirement that  $\mathbf{A}^T \mathbf{x} = \mathbf{b}$ .

## ■ 6.5 Weighted Path Finding

In this section, we formally define the formulation of linear program we solve (Subsection 6.5.1), introduce the weighted central path we follow (Subsection 6.5.3), define key properties of the path (Subsection 6.5.4) and the weights (Subsection 6.5.5) that we will use to produce an efficient path finding scheme.

### ■ 6.5.1 The Problem

As we explained in subsection 6.4.2, the standard dual linear program is not suitable for the maximum flow application. From now on, our goal is to efficiently solve the following linear program

$$\begin{aligned} & \min && \mathbf{c}^T \mathbf{x} \\ & \mathbf{x} \in \mathbb{R}^m & : & \mathbf{A}^T \mathbf{x} = \mathbf{b} \\ & \forall i \in [m] & : & l_i \leq x_i \leq u_i \end{aligned} \tag{6.9}$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^n$ ,  $\mathbf{c} \in \mathbb{R}^m$ ,  $l_i \in \mathbb{R} \cup \{-\infty\}$ , and  $u_i \in \mathbb{R} \cup \{+\infty\}$ .<sup>4</sup> We assume that for all  $i \in [m]$  the domain of variable  $x_i$ ,  $\text{dom}(x_i) \stackrel{\text{def}}{=} \{x : l_i \leq x \leq u_i\}$ , is non-degenerate. In particular we assume that  $\text{dom}(x_i)$  is not the empty set, a singleton, or the entire real line, i.e.  $l_i < u_i$  and either  $l_i \neq -\infty$  or  $u_i \neq +\infty$ . Furthermore we make the standard assumptions that  $\mathbf{A}$  has full column rank, and therefore  $m \geq n$ , and we assume that the interior of the polytope,  $\Omega \stackrel{\text{def}}{=} \{\mathbf{x} \in \mathbb{R}^m : \mathbf{A}^T \mathbf{x} = \mathbf{b}, l_i < x_i < u_i\}$ , is non-empty.

The linear program (6.9) is a generalization of standard form, the case where for all  $i \in [m]$  we have  $l_i = 0$  and  $u_i = +\infty$ . While it is well known that all linear programs can be written in standard form, the transformations to rewrite (6.9) in standard form may increase the rank of  $\mathbf{A}$  and therefore we solve (6.9) directly.

## ■ 6.5.2 Coordinate Barrier Functions

Rather than working directly with the different domain of the  $x_i$  we take a slightly more general approach and for the remainder of the chapter assume that for all  $i \in [m]$  we have a *barrier function*,  $\phi_i : \text{dom}(x_i) \rightarrow \mathbb{R}$ , such that

$$\lim_{x \rightarrow l_i} \phi_i(x) = \lim_{x \rightarrow u_i} \phi_i(x) = +\infty.$$

More precisely, we assume that each  $\phi_i$  is a *1-self-concordant barrier function*. Recall from definition 6.3.1,  $\phi$  is a *1-self-concordant barrier function* if

$$|\phi'''(x)| \leq 2(\phi''(x))^{3/2} \text{ for all } x \in \text{dom}(\phi) \quad (6.10)$$

and

$$|\phi'(x)| \leq \sqrt{\phi''(x)} \text{ for all } x \in \text{dom}(\phi). \quad (6.11)$$

The first condition (6.10) bounds how quickly the second order approximation to the function can change and the second condition (6.11) bounds how much force the barrier can exert.

The existence of a self-concordant barrier for the domain is a standard assumption for interior point methods [211]. However, for completeness, here we show how for each possible setting of the  $l_i$  and  $u_i$  there is an explicit 1-self-concordant barrier function we can use:

- *Case (1):  $l_i$  finite and  $u_i = \infty$ :* Here we use a *log barrier* defined as  $\phi_i(x) \stackrel{\text{def}}{=} -\log(x - l_i)$ . For this barrier we have

$$\phi_i'(x) = -\frac{1}{x - l_i}, \quad \phi_i''(x) = \frac{1}{(x - l_i)^2}, \quad \text{and} \quad \phi_i'''(x) = -\frac{2}{(x - l_i)^3}$$

and therefore clearly  $|\phi_i'''(x)| = 2(\phi_i''(x))^{3/2}$ ,  $|\phi_i'(x)| = \sqrt{\phi_i''(x)}$ , and  $\lim_{x \rightarrow l_i} \phi_i(x) = \infty$ .

- *Case (2):  $l_i = -\infty$  and  $u_i$  finite:* Here we use a *log barrier* defined as  $\phi_i(x) \stackrel{\text{def}}{=} -\log(u_i - x)$ . For this barrier we have

$$\phi_i'(x) = \frac{1}{u_i - x}, \quad \phi_i''(x) = \frac{1}{(u_i - x)^2}, \quad \text{and} \quad \phi_i'''(x) = -\frac{2}{(u_i - x)^3}$$

and therefore clearly  $|\phi_i'''(x)| = 2(\phi_i''(x))^{3/2}$ ,  $|\phi_i'(x)| = \sqrt{\phi_i''(x)}$ , and  $\lim_{x \rightarrow u_i} \phi_i(x) = \infty$ .

<sup>4</sup>Typically (6.9) is written as  $\mathbf{Ax} = \mathbf{b}$  rather than  $\mathbf{A}^T \mathbf{x} = \mathbf{b}$ . We chose this formulation to be consistent with the previous section and to be consistent with the standard use of  $n$  to denote the number of vertices and  $m$  to denote the number of edges in a graph in the linear program formulation of flow problems.

- *Case (3):  $l_i$  finite and  $u_i$  finite:* Here we use a *trigonometric barrier*<sup>5</sup> defined as  $\phi_i(x) \stackrel{\text{def}}{=} -\log \cos(a_i x + b_i)$  for  $a_i = \frac{\pi}{u_i - l_i}$  and  $b_i = -\frac{\pi}{2} \frac{u_i + l_i}{u_i - l_i}$ . Note for this choice as  $x \rightarrow u_i$  we have  $a_i x + b_i \rightarrow \frac{\pi}{2}$  and as  $x \rightarrow l_i$  we have  $a_i x + b_i \rightarrow -\frac{\pi}{2}$  and in both cases  $\phi_i(x) \rightarrow \infty$ . Furthermore,

$$\phi_i'(x) = a_i \tan(a_i x + b_i) \quad , \quad \phi_i''(x) = \frac{a_i^2}{\cos^2(a_i x + b_i)} \quad , \quad \text{and} \quad \phi_i'''(x) = \frac{2a_i^3 \sin(a_i x + b_i)}{\cos^3(a_i x + b_i)}.$$

Therefore, we have  $|\phi_i'''(x)| \leq 2(\phi_i''(x))^{3/2}$  and  $|\phi_i'(x)| \leq \sqrt{\phi_i''(x)}$ .

For the remainder of this chapter we will simply assume that we have a 1 self-concordant barrier  $\phi_i$  for each of the  $\text{dom}(\phi_i)$  and will not use any more structure about the barriers.

While there is much theory regarding properties of self-concordant barrier functions we will primarily use two common properties about self-concordant barrier functions. The first property, Lemma 6.5.1, shows that the Hessian of the barrier cannot change too quickly, and the second property, Lemma 6.5.2 we use to reason about how the force exerted by the barrier changes over the domain.

**Lemma 6.5.1** ([206, Theorem 4.1.6]). *Suppose  $\phi$  is a 1-self-concordant barrier function. For all  $s \in \text{dom}(\phi)$  if  $r \stackrel{\text{def}}{=} \sqrt{\phi''(s)} |s - t| < 1$  then  $t \in \text{dom}(\phi)$  and*

$$(1 - r) \sqrt{\phi''(s)} \leq \sqrt{\phi''(t)} \leq \frac{\sqrt{\phi''(s)}}{1 - r}.$$

**Lemma 6.5.2** ([206, Theorem 4.2.4]). *Suppose  $\phi$  is a 1-self-concordant barrier function. For all  $x, y \in \text{dom}(\phi)$ , we have*

$$\phi'(x) \cdot (y - x) \leq 1.$$

### ■ 6.5.3 The Weighted Central Path

Our linear programming algorithm maintains a feasible point  $\mathbf{x} \in \Omega$ , weights  $\mathbf{w} \in \mathbb{R}_{>0}^m$ , and minimizes the following *penalized objective function*

$$\min_{\mathbf{A}^T \mathbf{x} = \mathbf{b}} f_t(\mathbf{x}, \mathbf{w}) \quad \text{where} \quad f_t(\mathbf{x}, \mathbf{w}) \stackrel{\text{def}}{=} t \cdot \mathbf{c}^T \mathbf{x} + \sum_{i \in [m]} w_i \phi_i(\mathbf{x}_i) \quad (6.12)$$

for increasing  $t$  and small  $\mathbf{w}$ . For every fixed set of *weights*,  $\mathbf{w} \in \mathbb{R}_{>0}^m$  the set of points  $\mathbf{x}_{\mathbf{w}}(t) = \arg \min_{\mathbf{x} \in \Omega} f_t(\mathbf{x}, \mathbf{w})$  for  $t \in [0, \infty)$  forms a path through the interior of the polytope that we call the *weighted central path*. We call  $\mathbf{x}_{\mathbf{w}}(0)$  a *weighted center* of the polytope and note that  $\lim_{t \rightarrow \infty} \mathbf{x}_{\mathbf{w}}(0)$  is a solution to the linear program.

Our weighted path finding algorithm follows a simple iterative scheme. We assume we have a feasible point  $\{\mathbf{x}, \mathbf{w}\} \in \{\Omega \times \mathbb{R}_{>0}^m\}$  and a weight function  $\mathbf{g}(\mathbf{x}) : \Omega \rightarrow \mathbb{R}_{>0}^m$ , such that for any point  $\mathbf{x} \in \mathbb{R}_{>0}^m$  the function  $\mathbf{g}(\mathbf{x})$  returns a good set of weights that suggest a possibly better weighted path. Although the weight function we use is motivated by the barrier function introduced in Theorem 6.3.5, we do not assume anything on  $\mathbf{g}$  in this section. Our algorithm then repeats the following.

1. If  $\mathbf{x}$  close to  $\arg \min_{\mathbf{y} \in \Omega} f_t(\mathbf{y}, \mathbf{w})$ , then increase  $t$ .
2. Otherwise, use projected Newton step to update  $\mathbf{x}$  and move  $\mathbf{w}$  closer to  $\mathbf{g}(\mathbf{x})$ .
3. Repeat.

---

<sup>5</sup>The “trigonometric barrier” we use arises as the (unique) solution of the ODE  $\phi''' = 2(\phi'')^{3/2}$  such that the function value goes to infinity up at  $u_i$  and  $l_i$ .

In the remainder of this section we present how we measure both the quality of a current feasible point  $\{\mathbf{x}, \mathbf{w}\} \in \{\Omega \times \mathbb{R}_{>0}^m\}$  and the quality of the weight function. In Section 6.5.4 we derive and present both how we measure how close  $\{\mathbf{x}, \mathbf{w}\}$  is to the weighted central path and the step we take to improve this *centrality*. Then in Section 6.5.5 we present how we measure the quality of a weight function, i.e. how good the weighted paths it finds are.

### ■ 6.5.4 Measuring Centrality.

Here we explain how we measure the distance from  $\mathbf{x}$  to the minimum of  $f_t(\mathbf{x}, \mathbf{w})$  for fixed  $\mathbf{w}$ . This distance is a measure of how close  $\mathbf{x}$  is to the weighted central path and we refer to it as the *centrality* of  $\mathbf{x}$ , denote  $\delta_t(\mathbf{x}, \mathbf{w})$ .

To motivate our centrality measure, we first compute a projected Newton step for  $\mathbf{x}$ . For all  $\mathbf{x} \in \Omega$ , we define  $\phi(\mathbf{x}) \in \mathbb{R}^m$  by  $\phi(\mathbf{x})_i = \phi_i(\mathbf{x}_i)$  for  $i \in [m]$ . We define  $\phi'(\mathbf{x})$ ,  $\phi''(\mathbf{x})$ , and  $\phi'''(\mathbf{x})$  similarly and let  $\Phi, \Phi', \Phi'', \Phi'''$  denote the diagonal matrices corresponding to these matrices. Using this, we have<sup>6</sup>

$$\nabla_{\mathbf{x}} f_t(\mathbf{x}, \mathbf{w}) = t \cdot \mathbf{c} + \mathbf{w} \phi'(\mathbf{x}) \quad \text{and} \quad \nabla_{\mathbf{x}\mathbf{x}} f_t(\mathbf{x}, \mathbf{w}) = \mathbf{W} \Phi''(\mathbf{x}).$$

Therefore, a Newton step for  $\mathbf{x}$  is given by

$$\begin{aligned} \mathbf{h}_t(\mathbf{x}, \mathbf{w}) &= -(\mathbf{W} \Phi''(\mathbf{x}))^{-1/2} \mathbf{P}_{\mathbf{A}^T(\mathbf{W} \Phi''(\mathbf{x}))^{-1/2}} (\mathbf{W} \Phi''(\mathbf{x}))^{-1/2} \nabla_{\mathbf{x}} f_t(\mathbf{x}, \mathbf{w}) \\ &= -\Phi''(\mathbf{x})^{-1/2} \mathbf{P}_{\mathbf{x}, \mathbf{w}} \mathbf{W}^{-1} \Phi''(\mathbf{x})^{-1/2} \nabla_{\mathbf{x}} f_t(\mathbf{x}, \mathbf{w}) \end{aligned} \quad (6.13)$$

where  $\mathbf{P}_{\mathbf{A}^T(\mathbf{W} \Phi''(\mathbf{x}))^{-1/2}}$  is the orthogonal projection onto the kernel of  $\mathbf{A}^T(\mathbf{W} \Phi''(\mathbf{x}))^{-1/2}$  and  $\mathbf{P}_{\mathbf{x}, \mathbf{w}}$  is the orthogonal projection onto the kernel of  $\mathbf{A}^T(\Phi''(\mathbf{x}))^{-1/2}$  with respect to the norm  $\|\cdot\|_{\mathbf{W}}$ , i.e.

$$\mathbf{P}_{\mathbf{x}, \mathbf{w}} \stackrel{\text{def}}{=} \mathbf{I} - \mathbf{W}^{-1} \mathbf{A}_x (\mathbf{A}_x^T \mathbf{W}^{-1} \mathbf{A}_x)^{-1} \mathbf{A}_x^T \quad \text{for} \quad \mathbf{A}_x \stackrel{\text{def}}{=} \Phi''(\mathbf{x})^{-1/2} \mathbf{A}. \quad (6.14)$$

As with standard convergence analysis of Newton's method, we wish to keep the Newton step size in the Hessian norm, i.e.  $\|\mathbf{h}_t(\mathbf{x}, \mathbf{w})\|_{\mathbf{W} \Phi''(\mathbf{x})} = \|\sqrt{\phi''(\mathbf{x})} \mathbf{h}_t(\mathbf{x}, \mathbf{w})\|_{\mathbf{W}}$ , small and the multiplicative change in the Hessian,  $\|\sqrt{\phi''(\mathbf{x})} \mathbf{h}_t(\mathbf{x}, \mathbf{w})\|_{\infty}$ , small (See Lemma 6.5.2). While in the unweighted case we can bound the multiplicative change by the change in the hessian norm (since  $\|\cdot\|_{\infty} \leq \|\cdot\|_2$ ), here we would like to use small weights and this comparison would be insufficient.

To track both these quantities simultaneously, we define the *mixed norm* for all  $\mathbf{y} \in \mathbb{R}^m$  by

$$\|\mathbf{y}\|_{\mathbf{w}+\infty} \stackrel{\text{def}}{=} \|\mathbf{y}\|_{\infty} + C_{\text{norm}} \|\mathbf{y}\|_{\mathbf{W}} \quad (6.15)$$

for some  $C_{\text{norm}} > 0$  that we define later. Note that  $\|\cdot\|_{\mathbf{w}+\infty}$  is indeed a norm for  $\mathbf{w} \in \mathbb{R}_{>0}^m$  as in this case both  $\|\cdot\|_{\infty}$  and  $\|\cdot\|_{\mathbf{W}}$  are norms. However, rather than measuring centrality by the quantity  $\|\sqrt{\phi''(\mathbf{x})} \mathbf{h}_t(\mathbf{x}, \mathbf{w})\|_{\mathbf{w}+\infty} = \left\| \mathbf{P}_{\mathbf{x}, \mathbf{w}} \left( \frac{\nabla_{\mathbf{x}} f_t(\mathbf{x}, \mathbf{w})}{\mathbf{w} \sqrt{\phi''(\mathbf{x})}} \right) \right\|_{\mathbf{w}+\infty}$ , we instead find it more convenient to use the following idealized form

$$\delta_t(\mathbf{x}, \mathbf{w}) \stackrel{\text{def}}{=} \min_{\eta \in \mathbb{R}^n} \left\| \frac{\nabla_{\mathbf{x}} f_t(\mathbf{x}, \mathbf{w}) - \mathbf{A} \eta}{\mathbf{w} \sqrt{\phi''(\mathbf{x})}} \right\|_{\mathbf{w}+\infty}.$$

We justify this definition by showing these two quantities differ by at most a multiplicative factor of  $\|\mathbf{P}_{\mathbf{x}, \mathbf{w}}\|_{\mathbf{w}+\infty}$  as follows

$$\delta_t(\mathbf{x}, \mathbf{w}) \leq \left\| \sqrt{\phi''(\mathbf{x})} \mathbf{h}_t(\mathbf{x}, \mathbf{w}) \right\|_{\mathbf{w}+\infty} \leq \|\mathbf{P}_{\mathbf{x}, \mathbf{w}}\|_{\mathbf{w}+\infty} \cdot \delta_t(\mathbf{x}, \mathbf{w}). \quad (6.16)$$

<sup>6</sup>Recall that  $\mathbf{w} \phi'(\mathbf{x})$  denotes the entry-wise multiplication of the vectors  $\mathbf{w}$  and  $\phi'(\mathbf{x})$ .



This a direct consequence of the more general Lemma 6.11.1 that we prove in the appendix.

We summarize this section with the following definition.

**Definition 6.5.3** (Centrality Measure). For  $\{\mathbf{x}, \mathbf{w}\} \in \{\Omega \times \mathbb{R}_{>0}^m\}$  and  $t \geq 0$ , we let  $\mathbf{h}_t(\mathbf{x}, \mathbf{w})$  denote the *projected newton step* for  $\mathbf{x}$  on the penalized objective  $f_t$  given by

$$\mathbf{h}_t(\mathbf{x}, \mathbf{w}) \stackrel{\text{def}}{=} -\frac{1}{\sqrt{\vec{\phi}''(\mathbf{x})}} \mathbf{P}_{\mathbf{x}, \mathbf{w}} \left( \frac{\nabla_x f_t(\mathbf{x}, \mathbf{w})}{\mathbf{w} \sqrt{\vec{\phi}''(\mathbf{x})}} \right)$$

where  $\mathbf{P}_{\mathbf{x}, \mathbf{w}}$  is the orthogonal projection onto the kernel of  $\mathbf{A}^T (\Phi'')^{-1/2}$  with respect to the norm  $\|\cdot\|_{\mathbf{W}}$  (see (6.15)). We measure the *centrality* of  $\{\mathbf{x}, \mathbf{w}\}$  by

$$\delta_t(\mathbf{x}, \mathbf{w}) \stackrel{\text{def}}{=} \min_{\boldsymbol{\eta} \in \mathbb{R}^n} \left\| \frac{\nabla_x f_t(\mathbf{x}, \mathbf{w}) - \mathbf{A}\boldsymbol{\eta}}{\mathbf{w} \sqrt{\vec{\phi}''(\mathbf{x})}} \right\|_{\mathbf{w} + \infty} \quad (6.17)$$

where for all  $\mathbf{y} \in \mathbb{R}^m$  we let  $\|\mathbf{y}\|_{\mathbf{w} + \infty} \stackrel{\text{def}}{=} \|\mathbf{y}\|_{\infty} + C_{\text{norm}} \|\mathbf{y}\|_{\mathbf{W}}$  for some  $C_{\text{norm}} > 0$  we define later.

### ■ 6.5.5 The Weight Function

With the Newton step and centrality conditions defined, the specification of our algorithm becomes more clear. Our algorithm is as follows

1. If  $\delta_t(\mathbf{x}, \mathbf{w})$  is small, then increase  $t$ .
2. Set  $\mathbf{x}^{(\text{new})} \leftarrow \mathbf{x} + \mathbf{h}_t(\mathbf{x}, \mathbf{w})$  and move  $\mathbf{w}^{(\text{new})}$  towards  $\mathbf{g}(\mathbf{x}^{(\text{new})})$ .
3. Repeat.

To prove this algorithm converges, we need to show what happens to  $\delta_t(\mathbf{x}, \mathbf{w})$  when we change  $t$ ,  $\mathbf{x}$  and  $\mathbf{w}$ . At the heart of this chapter is understanding what conditions we need to impose on the weight function  $\mathbf{g}(\mathbf{x}) : \Omega \rightarrow \mathbb{R}_{>0}^m$  so that we can bound this change in  $\delta_t(\mathbf{x}, \mathbf{w})$  and hence achieve fast converge rates. In Lemma 6.6.1 we show that the effect of changing  $t$  on  $\delta_t$  is bounded by  $C_{\text{norm}}$  and  $\|\mathbf{g}(\mathbf{x})\|_1$ , in Lemma 6.6.2 we show that the effect that a Newton Step on  $\mathbf{x}$  has on  $\delta_t$  is bounded by  $\|\mathbf{P}_{\mathbf{x}, \mathbf{g}(\mathbf{x})}\|_{\mathbf{g}(\mathbf{x}) + \infty}$ , and in Lemma 6.6.3 and 6.6.4 we show the change of  $\mathbf{w}$  as  $\mathbf{g}(\mathbf{x})$  changes is bounded by  $\|\mathbf{G}(\mathbf{x})^{-1} \mathbf{G}'(\mathbf{x}) (\Phi''(\mathbf{x}))^{-1/2}\|_{\mathbf{g}(\mathbf{x}) + \infty}$ .

Hence for the remainder of the chapter we assume we have a weight function  $\mathbf{g}(\mathbf{x}) : \Omega \rightarrow \mathbb{R}_{>0}^m$  and make the following assumptions regarding our weight function. In Section 6.7 we prove that such weight function exists.

**Definition 6.5.4** (Weight Function). A *weight function* is a differentiable function from  $\mathbf{g} : \Omega \rightarrow \mathbb{R}_{>0}^m$  such that for constants  $c_1(\mathbf{g})$ ,  $c_\gamma(\mathbf{g})$ , and  $c_\delta(\mathbf{g})$ , we have the following for all  $\mathbf{x} \in \Omega$ :

- *Size* : The size  $c_1(\mathbf{g}) = \|\mathbf{g}(\mathbf{x})\|_1$ .
- *Slack Sensitivity*: The slack sensitivity  $c_\gamma(\mathbf{g})$  satisfies  $1 \leq c_\gamma(\mathbf{g}) \leq \frac{5}{4}$  and  $\|\mathbf{P}_{\mathbf{x}, \mathbf{w}}\|_{\mathbf{w} + \infty} \leq c_\gamma(\mathbf{g})$  for any  $\mathbf{w}$  such that  $\frac{4}{5}\mathbf{g}(\mathbf{x}) \leq \mathbf{w} \leq \frac{5}{4}\mathbf{g}(\mathbf{x})$ .
- *Step Consistency* : The step consistency  $c_\delta(\mathbf{g})$  satisfies  $c_\delta(\mathbf{g}) \cdot c_\gamma(\mathbf{g}) < 1$  and

$$\left\| \mathbf{G}(\mathbf{x})^{-1} \mathbf{G}'(\mathbf{x}) (\Phi''(\mathbf{x}))^{-1/2} \right\|_{\mathbf{g}(\mathbf{x}) + \infty} \leq c_\delta \leq 1.$$

- *Uniformity* : The weight function satisfies  $\|\mathbf{g}(\mathbf{x})\|_{\infty} \leq 2$ .

### ■ 6.5.6 Notation

Here we introduce various notation that we will use throughout the chapter (excluding Section 6.13). This should be used primarily for reference as we reintroduce notation as needed later in the chapter. For a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and vector  $\mathbf{x} \in \mathbb{R}_{>0}^m$  we define the following matrices and vectors

- Projection matrix  $\mathbf{P}_\mathbf{A}(\mathbf{x}) \in \mathbb{R}^{m \times m}$ :  $\mathbf{P}_\mathbf{A}(\mathbf{x}) \stackrel{\text{def}}{=} \mathbf{X}^{1/2} \mathbf{A} (\mathbf{A}^T \mathbf{X} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{X}^{1/2}$ .
- Leverage scores  $\boldsymbol{\sigma}_\mathbf{A}(\mathbf{x}) \in \mathbb{R}^m$ :  $\boldsymbol{\sigma}_\mathbf{A}(\mathbf{x}) \stackrel{\text{def}}{=} \text{diag}(\mathbf{P}_\mathbf{A}(\mathbf{x}))$ .
- Leverage matrix  $\boldsymbol{\Sigma}_\mathbf{A}(\mathbf{x}) \in \mathbb{R}^{m \times m}$ :  $\boldsymbol{\Sigma}_\mathbf{A}(\mathbf{x}) \stackrel{\text{def}}{=} \text{Diag}(\mathbf{P}_\mathbf{A}(\mathbf{x}))$ .
- Projection Laplacian  $\boldsymbol{\Lambda}_\mathbf{A}(\mathbf{x}) \in \mathbb{R}^{m \times m}$ :  $\boldsymbol{\Lambda}_\mathbf{A}(\mathbf{x}) \stackrel{\text{def}}{=} \boldsymbol{\Sigma}_\mathbf{A}(\mathbf{x}) - \mathbf{P}_\mathbf{A}(\mathbf{x})^{(2)}$ .

## ■ 6.6 Progressing Along Weighted Paths

In this section, we provide the main lemmas we need for an  $\tilde{O}(\sqrt{\text{rank}(\mathbf{A})} \log(U/\varepsilon))$  iterations weighted path following algorithm for (6.9) assuming a weight function satisfying Definition 6.5.5. In Section 6.6.1, 6.6.2, and 6.6.3 we show how centrality,  $\delta_t(\mathbf{x}, \mathbf{w})$ , is affected by changing  $t$ ,  $\mathbf{x} \in \Omega$ , and  $\mathbf{w} \in \mathbb{R}_{>0}^m$  respectively. In Section 6.6.5 we then show how to use these Lemmas to improve centrality using approximate computations of the weight function,  $\mathbf{g} : \Omega \rightarrow \mathbb{R}_{>0}^m$ .

### ■ 6.6.1 Changing $t$

Here we bound how much centrality increases as we increase  $t$ . We show that this rate of increase is governed by  $C_{\text{norm}}$  and  $\|\mathbf{w}\|_1$ .

**Lemma 6.6.1.** *For all  $\{\mathbf{x}, \mathbf{w}\} \in \{\Omega \times \mathbb{R}_{>0}^m\}$ ,  $t > 0$  and  $\alpha \geq 0$ , we have*

$$\delta_{(1+\alpha)t}(\mathbf{x}, \mathbf{w}) \leq (1 + \alpha)\delta_t(\mathbf{x}, \mathbf{w}) + \alpha \left(1 + C_{\text{norm}} \sqrt{\|\mathbf{w}\|_1}\right).$$

*Proof.* Let  $\boldsymbol{\eta}_t \in \mathbb{R}^n$  be such that

$$\delta_t(\mathbf{x}, \mathbf{w}) = \left\| \frac{\nabla_x f_t(\mathbf{x}, \mathbf{w}) + \mathbf{A}\boldsymbol{\eta}_t}{\mathbf{w}\sqrt{\boldsymbol{\phi}''(\mathbf{x})}} \right\|_{\mathbf{w}+\infty} = \left\| \frac{t \cdot \mathbf{c} + \mathbf{w}\boldsymbol{\phi}'(\mathbf{x}) + \mathbf{A}\boldsymbol{\eta}_t}{\mathbf{w}\sqrt{\boldsymbol{\phi}''(\mathbf{x})}} \right\|_{\mathbf{w}+\infty}.$$

Applying this to the definition of  $\delta_{(1+\alpha)t}$  and using that  $\|\cdot\|_{\mathbf{w}+\infty}$  is a norm then yields

$$\begin{aligned} \delta_{(1+\alpha)t}(\mathbf{x}, \mathbf{w}) &= \min_{\boldsymbol{\eta} \in \mathbb{R}^n} \left\| \frac{(1 + \alpha)t \cdot \mathbf{c} + \mathbf{w}\boldsymbol{\phi}'(\mathbf{x}) + \mathbf{A}\boldsymbol{\eta}}{\mathbf{w}\sqrt{\boldsymbol{\phi}''(\mathbf{x})}} \right\|_{\mathbf{w}+\infty} \\ &\leq \left\| \frac{(1 + \alpha)t \cdot \mathbf{c} + \mathbf{w}\boldsymbol{\phi}'(\mathbf{x}) + \mathbf{A}(1 + \alpha)\boldsymbol{\eta}_t}{\mathbf{w}\sqrt{\boldsymbol{\phi}''(\mathbf{x})}} \right\|_{\mathbf{w}+\infty} \\ &\leq (1 + \alpha) \left\| \frac{t \cdot \mathbf{c} + \mathbf{w}\boldsymbol{\phi}'(\mathbf{x}) + \mathbf{A}\boldsymbol{\eta}_t}{\mathbf{w}\sqrt{\boldsymbol{\phi}''(\mathbf{x})}} \right\|_{\mathbf{w}+\infty} + \alpha \left\| \frac{\mathbf{w}\boldsymbol{\phi}'(\mathbf{x})}{\mathbf{w}\sqrt{\boldsymbol{\phi}''(\mathbf{x})}} \right\|_{\mathbf{w}+\infty} \\ &= (1 + \alpha)\delta_t(\mathbf{x}, \mathbf{w}) + \alpha \left( \left\| \frac{\boldsymbol{\phi}'(\mathbf{x})}{\sqrt{\boldsymbol{\phi}''(\mathbf{x})}} \right\|_{\infty} + C_{\text{norm}} \left\| \frac{\boldsymbol{\phi}'(\mathbf{x})}{\sqrt{\boldsymbol{\phi}''(\mathbf{x})}} \right\|_{\mathbf{w}} \right) \end{aligned}$$

Using that  $|\phi'_i(\mathbf{x})| \leq \sqrt{\phi''_i(\mathbf{x})}$  for all  $i \in [m]$  and  $\mathbf{x} \in \mathbb{R}^m$  by Equation (6.11) yields the result.  $\square$

### ■ 6.6.2 Changing $x$

Here we analyze the effect of a Newton step on  $x$  on centrality. We show for sufficiently central  $\{x, w\} \in \{\Omega \times \mathbb{R}_{>0}^m\}$  and  $w$  sufficiently close to  $g(x)$  Newton steps converge quadratically.

**Lemma 6.6.2.** *Let  $\{x_0, w\} \in \{\Omega \times \mathbb{R}_{>0}^m\}$  such that  $\delta_t(x_0, w) \leq \frac{1}{10}$  and  $\frac{4}{5}g(x) \leq w \leq \frac{5}{4}g(x)$  and consider a Newton step  $x_1 = x_0 + h_t(x, w)$ . Then,  $\delta_t(x_1, w) \leq 4(\delta_t(x_0, w))^2$ .*

*Proof.* Let  $\phi_0 \stackrel{\text{def}}{=} \phi(x_0)$  and let  $\phi_1 \stackrel{\text{def}}{=} \phi(x_1)$ . By the definition of  $h_t(x_0, w)$  and the formula of  $P_{x_0, w}$  we know that there is some  $\eta_0 \in \mathbb{R}^n$  such that

$$-\sqrt{\vec{\phi}_0''} h_t(x_0, w) = \frac{t \cdot c + w\phi_0' - A\eta_0}{w\sqrt{\vec{\phi}_0''}}.$$

Therefore,  $A\eta_0 = c + w\phi_0' + w\phi_0'' h_t(x_0, w)$ . Recalling the definition of  $\delta_t$  this implies that

$$\begin{aligned} \delta_t(x_1, w) &= \min_{\eta \in \mathbb{R}^n} \left\| \frac{t \cdot c + w\vec{\phi}_1' - A\eta}{w\sqrt{\vec{\phi}_1''}} \right\|_{w+\infty} \leq \left\| \frac{t \cdot c + w\vec{\phi}_1' - A\eta_0}{w\sqrt{\vec{\phi}_1''}} \right\|_{w+\infty} \\ &\leq \left\| \frac{w(\vec{\phi}_1' - \phi_0') - w\phi_0'' h_t(x_0, w)}{w\sqrt{\vec{\phi}_1''}} \right\|_{w+\infty} = \left\| \frac{(\vec{\phi}_1' - \phi_0') - \phi_0'' h_t(x_0, w)}{\sqrt{\vec{\phi}_1''}} \right\|_{w+\infty} \end{aligned}$$

By the mean value theorem, we have  $\vec{\phi}_1' - \phi_0' = \vec{\phi}''(\vec{\theta}) h_t(x_0, w)$  for some  $\vec{\theta}$  between  $x_0$  and  $x_1$  coordinate-wise. Hence,

$$\begin{aligned} \delta_t(x_1, w) &\leq \left\| \frac{\vec{\phi}''(\vec{\theta}) h_t(x_0, w) - \phi_0'' h_t(x_0, w)}{\sqrt{\vec{\phi}_1''}} \right\|_{w+\infty} = \left\| \frac{(\vec{\phi}''(\vec{\theta}) - \phi_0'')}{\sqrt{\vec{\phi}_1''} \sqrt{\vec{\phi}_0''}} (\sqrt{\vec{\phi}_0''} h_t(x_0, w)) \right\|_{w+\infty} \\ &\leq \left\| \frac{\vec{\phi}''(\vec{\theta}) - \phi_0''}{\sqrt{\vec{\phi}_1''} \sqrt{\vec{\phi}_0''}} \right\|_{\infty} \cdot \left\| \sqrt{\vec{\phi}_0''} h_t(x_0, w) \right\|_{w+\infty}. \end{aligned}$$

To bound the first term, we use Lemma 6.5.1 as follows

$$\begin{aligned} \left\| \frac{(\vec{\phi}''(\vec{\theta}) - \phi_0'')}{\sqrt{\vec{\phi}_1''} \sqrt{\vec{\phi}_0''}} \right\|_{\infty} &\leq \left\| \frac{\phi''(\vec{\theta})}{\phi_0''} - \mathbf{1} \right\|_{\infty} \cdot \left\| \frac{\sqrt{\vec{\phi}_0''}}{\sqrt{\vec{\phi}_1''}} \right\|_{\infty} \\ &\leq \left| \left(1 - \left\| \sqrt{\vec{\phi}_0''} h_t(x_0, w) \right\|_{\infty}\right)^{-2} - 1 \right| \cdot \left(1 - \left\| \sqrt{\vec{\phi}_0''} h_t(x_0, w) \right\|_{\infty}\right)^{-1}. \end{aligned}$$

Using (6.16), i.e. Lemma 6.11.1, the bound  $c_\gamma \leq 2$ , and the assumption on  $\delta_t(x_0, w)$ , we have

$$\left\| \sqrt{\vec{\phi}_0''} h_t(x_0, w) \right\|_{\infty} \leq \left\| \sqrt{\vec{\phi}_0''} h_t(x_0, w) \right\|_{w+\infty} \leq c_\gamma \cdot \delta_t(x_0, w) \leq \frac{1}{5}.$$

Using  $((1-t)^{-2} - 1) \cdot (1-t)^{-1} \leq 4t$  for  $t \leq 1/5$ , we have

$$\left\| \frac{(\vec{\phi}''(\vec{\theta}) - \phi_0'')}{\sqrt{\vec{\phi}_1''} \sqrt{\vec{\phi}_0''}} \right\|_{\infty} \leq 4 \left\| \sqrt{\vec{\phi}_0''} h_t(x_0, w) \right\|_{\infty}.$$

Combining the above formulas yields that  $\delta_t(\mathbf{x}_1, \mathbf{w}) \leq 4(\delta_t(\mathbf{x}_0, \mathbf{w}))^2$  as desired.  $\square$

### ■ 6.6.3 Changing $\mathbf{w}$

In the previous subsection we used the assumption that the weights,  $\mathbf{w}$ , were multiplicatively close to the output of the weight function,  $\mathbf{g}(\mathbf{x})$ , for the current point  $\mathbf{x} \in \Omega$ . In order to maintain this invariant when we change  $\mathbf{x}$  we will need to change  $\mathbf{w}$  to move it closer to  $\mathbf{g}(\mathbf{x})$ . Here we bound how much  $\mathbf{g}(\mathbf{x})$  can move as we move  $\mathbf{x}$  (Lemma 6.6.3) and we bound how much changing  $\mathbf{w}$  can hurt centrality (Lemma 6.6.4). Together these lemmas will allow us to show that we can keep  $\mathbf{w}$  close to  $\mathbf{g}(\mathbf{x})$  while still improving centrality (Section 6.6.5).

**Lemma 6.6.3.** *For all  $t \in [0, 1]$ , let  $\mathbf{x}_t \stackrel{\text{def}}{=} \mathbf{x}_0 + t\Delta_x$  for  $\Delta_x \in \mathbb{R}^m$ ,  $\mathbf{x}_t \in \Omega$ ,  $\mathbf{g}_t = \mathbf{g}(\mathbf{x}_t)$  and  $\varepsilon = \left\| \sqrt{\vec{\phi}_0''} \Delta_x \right\|_{\mathbf{g}_0+\infty} \leq 0.1$ . Then*

$$\|\log(\mathbf{g}_1) - \log(\mathbf{g}_0)\|_{\mathbf{g}_0+\infty} \leq c_\delta \varepsilon (1 + 4\varepsilon) \leq 0.2$$

and for all  $s, t \in [0, 1]$  and for all  $\mathbf{y} \in \mathbb{R}^m$  we have

$$\|\mathbf{y}\|_{\mathbf{g}_s+\infty} \leq (1 + 2\varepsilon) \|\mathbf{y}\|_{\mathbf{g}_t+\infty}. \quad (6.18)$$

*Proof.* Let  $\mathbf{q} : [0, 1] \rightarrow \mathbb{R}^m$  be given by  $\mathbf{q}(t) \stackrel{\text{def}}{=} \log(\mathbf{g}_t)$  for all  $t \in [0, 1]$ . Then, we have

$$\mathbf{q}'(t) = \mathbf{G}_t^{-1} \mathbf{G}'_t \Delta_x.$$

Let  $Q(t) \stackrel{\text{def}}{=} \|\vec{q}(t) - \vec{q}(0)\|_{\mathbf{g}_0+\infty}$ . Using Jensen's inequality we have that for all  $u \in [0, 1]$ ,

$$Q(u) \leq \bar{Q}(u) \stackrel{\text{def}}{=} \int_0^u \left\| \mathbf{G}_t^{-1} \mathbf{G}'_t (\vec{\phi}_t'')^{-1/2} \right\|_{\mathbf{g}_0+\infty} \left\| \sqrt{\vec{\phi}_t''} \Delta_x \right\|_{\mathbf{g}_0+\infty} dt.$$

Using Lemma 6.5.1 and  $\varepsilon \leq \frac{1}{10}$ , we have for all  $t \in [0, 1]$ ,

$$\begin{aligned} \left\| \sqrt{\vec{\phi}_t''} \Delta_x \right\|_{\mathbf{g}_0+\infty} &\leq \left\| \sqrt{\vec{\phi}_t''} / \sqrt{\vec{\phi}_0''} \right\|_{\infty} \left\| \sqrt{\vec{\phi}_0''} \Delta_x \right\|_{\mathbf{g}_0+\infty} \\ &\leq \left( 1 - \left\| \sqrt{\vec{\phi}_0''} \Delta_x \right\|_{\infty} \right)^{-1} \left\| \sqrt{\vec{\phi}_0''} \Delta_x \right\|_{\mathbf{g}_0+\infty} \leq \frac{\varepsilon}{1 - \varepsilon}. \end{aligned}$$

Thus, we have

$$\bar{Q}(u) \leq \frac{\varepsilon}{1 - \varepsilon} \int_0^u \left\| \mathbf{G}_t^{-1} \mathbf{G}'_t (\vec{\phi}_t'')^{-1/2} \right\|_{\mathbf{g}_0+\infty} dt. \quad (6.19)$$

Note that  $\bar{Q}$  is monotonically increasing. Let  $\theta = \sup_{u \in [0, 1]} \{\bar{Q}(u) \leq c_\delta \varepsilon (1 + 4\varepsilon)\}$ . Since  $\bar{Q}(\theta) \leq \frac{1}{2}$ , we know that for all  $s, t \in [0, \theta]$ , we have

$$\left\| \frac{\mathbf{g}(\mathbf{x}_s) - \mathbf{g}(\mathbf{x}_t)}{\mathbf{g}(\mathbf{x}_t)} \right\|_{\infty} \leq \|\vec{q}(s) - \vec{q}(t)\|_{\infty} + \|\vec{q}(s) - \vec{q}(t)\|_{\infty}^2$$

and therefore

$$\|\mathbf{g}_s / \mathbf{g}_t\|_{\infty} \leq (1 + \|\vec{q}(s) - \vec{q}(t)\|_{\infty} + \|\vec{q}(s) - \vec{q}(t)\|_{\infty}^2)^2 \leq (1 + c_\delta \varepsilon (1 + 4\varepsilon))^2$$

Consequently,

$$\|\mathbf{y}\|_{\mathbf{g}_s+\infty} \leq (1 + c_\delta \varepsilon (1 + 4\varepsilon)) \|\mathbf{y}\|_{\mathbf{g}_t+\infty} \leq (1 + 2\varepsilon) \|\mathbf{y}\|_{\mathbf{g}_t+\infty}.$$

Using (6.19), we have for all  $u \in [0, \theta]$ ,

$$\begin{aligned} Q(u) \leq \bar{Q}(u) &\leq \frac{\varepsilon}{1-\varepsilon} \int_0^u \left\| \mathbf{G}_t^{-1} \mathbf{G}'_t \left( \vec{\phi}_t'' \right)^{-1/2} \right\|_{\mathbf{g}_0+\infty} dt \\ &\leq \frac{\varepsilon}{1-\varepsilon} \int_0^u (1+2\varepsilon) \left\| \mathbf{G}_t^{-1} \mathbf{G}'_t \left( \vec{\phi}_t'' \right)^{-1/2} \right\|_{\mathbf{g}_t+\infty} dt \\ &\leq \frac{\varepsilon}{1-\varepsilon} (1+2\varepsilon) c_\delta \theta \leq c_\delta \varepsilon (1+4\varepsilon). \end{aligned}$$

Consequently, we have that  $\theta = 1$  and we have the desired result with  $Q(1) \leq c_\delta \varepsilon (1+4\varepsilon) < \frac{1}{5}$ .  $\square$

**Lemma 6.6.4.** *Let  $\mathbf{v}, \mathbf{w} \in \mathbb{R}_{>0}^m$  such that  $\varepsilon = \|\log(\mathbf{w}) - \log(\mathbf{v})\|_{\mathbf{w}+\infty} \leq 0.1$ . Then for  $\mathbf{x} \in \Omega$  we have*

$$\delta_t(\mathbf{x}, \mathbf{v}) \leq (1+4\varepsilon)(\delta_t(\mathbf{x}, \mathbf{w}) + \varepsilon).$$

*Proof.* Let  $\boldsymbol{\eta}_w$  be such that

$$\delta_t(\mathbf{x}, \mathbf{w}) = \left\| \frac{\mathbf{c} + \mathbf{w} \vec{\phi}'(\mathbf{x}) - \mathbf{A} \boldsymbol{\eta}_w}{\mathbf{w} \sqrt{\vec{\phi}''(\mathbf{x})}} \right\|_{\mathbf{w}+\infty} \quad (6.20)$$

Furthermore, the assumption shows that  $(1+\varepsilon)^{-2} \mathbf{w}_i \leq \mathbf{v}_i \leq (1+\varepsilon)^2 \mathbf{w}_i$  for all  $i$ . Using these, we bound the energy with the new weights as follows

$$\begin{aligned} \delta_t(\mathbf{x}, \mathbf{v}) &= \min_{\boldsymbol{\eta}} \left\| \frac{\mathbf{c} + \mathbf{v} \vec{\phi}'(\mathbf{x}) - \mathbf{A} \boldsymbol{\eta}}{\mathbf{v} \sqrt{\vec{\phi}''(\mathbf{x})}} \right\|_{\mathbf{v}+\infty} \leq \left\| \frac{\mathbf{c} + \mathbf{v} \vec{\phi}'(\mathbf{x}) - \mathbf{A} \boldsymbol{\eta}_w}{\mathbf{v} \sqrt{\vec{\phi}''(\mathbf{x})}} \right\|_{\mathbf{v}+\infty} \\ &\leq (1+\varepsilon) \left\| \frac{\mathbf{c} + \mathbf{v} \vec{\phi}'(\mathbf{x}) - \mathbf{A} \boldsymbol{\eta}_w}{\mathbf{v} \sqrt{\vec{\phi}''(\mathbf{x})}} \right\|_{\mathbf{w}+\infty} \\ &\leq (1+\varepsilon) \cdot \left( \left\| \frac{\mathbf{c} + \mathbf{w} \vec{\phi}'(\mathbf{x}) - \mathbf{A} \boldsymbol{\eta}_w}{\mathbf{v} \sqrt{\vec{\phi}''(\mathbf{x})}} \right\|_{\mathbf{w}+\infty} + \left\| \frac{(\mathbf{v} - \mathbf{w}) \vec{\phi}'(\mathbf{x})}{\mathbf{v} \sqrt{\vec{\phi}''(\mathbf{x})}} \right\|_{\mathbf{w}+\infty} \right) \\ &\leq (1+\varepsilon)^3 \delta_t(\mathbf{x}, \mathbf{w}) + (1+\varepsilon) \cdot \left\| \frac{\phi'(\mathbf{x})}{\sqrt{\phi''(\mathbf{x})}} \right\|_{\infty} \cdot \left\| \frac{(\mathbf{v} - \mathbf{w})}{\mathbf{v}} \right\|_{\mathbf{w}+\infty} \end{aligned}$$

Using that  $|\phi'_i(\mathbf{x})| \leq \sqrt{\phi''_i(\mathbf{x})}$  for all  $i \in [m]$  by Equation (6.11), we have that

$$\begin{aligned} \delta_t(\mathbf{x}, \mathbf{v}) &\leq (1+\varepsilon)^3 \delta_t(\mathbf{x}, \mathbf{w}) + (1+\varepsilon)^2 \varepsilon \\ &\leq (1+4\varepsilon) (\delta_t(\mathbf{x}, \mathbf{w}) + \varepsilon). \end{aligned}$$

$\square$

### ■ 6.6.4 The Chasing 0 Game

In the previous subsection, we saw how much the weight function,  $\mathbf{g}(\mathbf{x})$ , can change after a Newton step on  $\mathbf{x}$  and we bounded how much we can move the weights without affecting centrality too much.

Here we study how to correctly move the weights even when we cannot compute the weight function exactly. To effectively use multiplicative approximations to the weight function, we need to smooth out changes to the weights by using some slowly changing approximation to the weight function. In this subsection, we present the smoothing problem in a general form that we call the *chasing 0 game* and we provide an effective strategy for playing this game. In next subsection, we show how to use this strategy to produce an efficient weighted path following scheme.

The *chasing 0 game* is as follows. There is a player, an adversary, and a point  $\mathbf{x} \in \mathbb{R}^m$ . The goal of the player is to keep the point close to  $\mathbf{0}$  in  $\ell_\infty$  norm and the goal of the adversary tries to move  $\mathbf{x}$  away from  $\mathbf{0} \in \mathbb{R}^m$ . The game proceeds for an infinite number of iterations where in each iteration the adversary moves the current point  $\mathbf{x}^{(k)} \in \mathbb{R}^m$  to some new point  $\mathbf{y}^{(k)} \in \mathbb{R}^m$  and the player needs to respond. The player does not know  $\mathbf{x}^{(k)}$ ,  $\mathbf{y}^{(k)}$ , or the move of the adversary. All the player knows is that the adversary moved the point within some convex set  $U^{(k)}$  and the player knows some  $\mathbf{z}^{(k)} \in \mathbb{R}^n$  that is close to  $\mathbf{y}^{(k)}$  in  $\ell_\infty$  norm. With this information the player is allowed to move the point a little more than the adversary. Formally, the player is allowed to set the next point to  $\mathbf{x}^{(k+1)} \in \mathbb{R}^m$  such that  $\Delta^{(k)} \stackrel{\text{def}}{=} \mathbf{x}^{(k+1)} - \mathbf{y}^{(k)} \in (1 + \varepsilon)U$  for some fixed  $\varepsilon > 0$ .

The question we would like to address is, how close the player can keep  $\mathbf{x}^{(k+1)}$  to  $\mathbf{0}$  in  $\ell_\infty$  norm? In particular, we would like an efficient strategy for computing  $\Delta^{(k)}$  such that  $\|\mathbf{x}^{(k)}\|_\infty$  is bounded for all  $k \geq 0$ .

---

**Algorithm 7:** Chasing 0 Game

---

Given  $R > 0, \varepsilon > 0, \mathbf{x}^{(0)} \in \mathbb{R}^m$ .

**for**  $k = 1, 2, \dots$  **do**

The adversary announces symmetric convex set  $U^{(k)} \subseteq \mathbb{R}^n$  and  $\mathbf{u}^{(k)} \in U^{(k)}$ .

The adversary sets  $\mathbf{y}^{(k)} := \mathbf{x}^{(k)} + \mathbf{u}^{(k)}$ .

The adversary announces  $\mathbf{z}^{(k)}$  such that  $\|\mathbf{z}^{(k)} - \mathbf{y}^{(k)}\|_\infty \leq R$ .

The player chooses  $\Delta^{(k)} \in (1 + \varepsilon)U^{(k)}$ .

The player sets  $\mathbf{x}^{(k+1)} = \mathbf{y}^{(k)} + \Delta^{(k)}$ .

**end**

---

We show that assuming that the  $U^{(k)}$  are sufficiently bounded then there is strategy that the player can follow to ensure that  $\|\mathbf{x}^{(k)}\|_\infty$  is never too large. Our strategy simply consists of taking “gradient steps” using the following potential function.

**Definition 6.6.5.** For any  $\mu \geq 0$  let  $p_\mu : \mathbb{R} \rightarrow \mathbb{R}$  and  $\Phi_\mu : \mathbb{R}^m \rightarrow \mathbb{R}$  be given by

$$\forall x \in \mathbb{R} \quad : \quad p_\mu(x) \stackrel{\text{def}}{=} e^{\mu x} + e^{-\mu x} \quad \text{and} \quad \Phi_\mu(\mathbf{x}) \stackrel{\text{def}}{=} \sum_{i \in [m]} p_\mu(x_i).$$

In other words, for all  $k$  we simply set  $\Delta^{(k)}$  to be the vector in  $(1 + \varepsilon)U^{(k)}$  that best minimizes the potential function of the observed position, i.e.  $\Phi_\mu(\mathbf{z}^{(k)})$  for an appropriate choice of  $\mu$ . In the following theorem we show that this suffices to keep  $\Phi_\mu(\mathbf{x}^{(k)})$  small and that small  $\Phi_\mu(\mathbf{x}^{(k)})$  implies small  $\|\mathbf{x}^{(k)}\|_\infty$  and hence has the desired properties.

**Theorem 6.6.6.** *Suppose that each  $U^{(k)}$  is a symmetric convex set that contains an  $\ell_\infty$  ball of radius  $r_k$  and is contained in a  $\ell_\infty$  ball of radius  $R_k \leq R$ . Let  $0 < \varepsilon < \frac{1}{5}$  and consider the strategy*

$$\Delta^{(k)} = (1 + \varepsilon) \arg \min_{\Delta \in U^{(k)}} \left\langle \nabla \Phi_\mu(\mathbf{z}^{(k)}), \Delta \right\rangle \quad \text{where} \quad \mu = \frac{\varepsilon}{12R}.$$

Let  $\tau \stackrel{\text{def}}{=} \max_k \frac{R_k}{r_k}$  and suppose  $\Phi_\mu(\mathbf{x}^{(0)}) \leq \frac{12m\tau}{\varepsilon}$  (or more specifically  $\|\mathbf{x}^{(0)}\|_\infty \leq \frac{12R}{\varepsilon} \log\left(\frac{6\tau}{\varepsilon}\right)$ ) then

$$\forall k \geq 0 \quad : \quad \Phi_\mu(\mathbf{x}^{(k+1)}) \leq \left(1 - \frac{\varepsilon^2 r_k}{24R}\right) \Phi_\mu(\mathbf{x}^{(k)}) + \varepsilon m \frac{R_k}{2R} \leq \frac{12m\tau}{\varepsilon}.$$

In particular, we have  $\|\mathbf{x}^{(k)}\|_\infty \leq \frac{12R}{\varepsilon} \log\left(\frac{12m\tau}{\varepsilon}\right)$ .

To prove Theorem 6.6.6 we first provide the following lemma regarding properties of the potential function  $\Phi_\mu$ .

**Lemma 6.6.7.** *For all  $\mathbf{x} \in \mathbb{R}^m$ , we have*

$$e^{\mu\|\mathbf{x}\|_\infty} \leq \Phi_\mu(\mathbf{x}) \leq 2me^{\mu\|\mathbf{x}\|_\infty} \quad \text{and} \quad \mu\Phi_\mu(\mathbf{x}) - 2\mu m \leq \|\nabla\Phi_\mu(\mathbf{x})\|_1 \quad (6.21)$$

Furthermore, for any symmetric convex set  $U \subseteq \mathbb{R}^m$  and any  $\mathbf{x} \in \mathbb{R}^m$ , let  $\mathbf{x}^b \stackrel{\text{def}}{=} \arg \max_{\mathbf{y} \in U} \langle \mathbf{x}, \mathbf{y} \rangle$  and  $\|\mathbf{x}\|_U \stackrel{\text{def}}{=} \max_{\mathbf{y} \in U} \langle \mathbf{x}, \mathbf{y} \rangle$ . Then for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^m$  with  $\|\mathbf{x} - \mathbf{y}\|_\infty \leq \delta \leq \frac{1}{5\mu}$  we have

$$e^{-\mu\delta} \|\nabla\Phi_\mu(\mathbf{y})\|_U - \mu \|\nabla\Phi_\mu(\mathbf{y})^b\|_1 \leq \left\langle \nabla\Phi_\mu(\mathbf{x}), \nabla\Phi_\mu(\mathbf{y})^b \right\rangle \leq e^{\mu\delta} \|\nabla\Phi_\mu(\mathbf{y})\|_U + \mu e^{\mu\delta} \|\nabla\Phi_\mu(\mathbf{y})^b\|_1. \quad (6.22)$$

If additionally  $U$  is contained in a  $\ell_\infty$  ball of radius  $R$  then

$$e^{-\mu\delta} \|\nabla\Phi_\mu(\mathbf{y})\|_U - \mu m R \leq \|\nabla\Phi_\mu(\mathbf{x})\|_U \leq e^{\mu\delta} \|\nabla\Phi_\mu(\mathbf{y})\|_U + \mu e^{\mu\delta} m R. \quad (6.23)$$

*Proof.* First we note that for all  $x \in \mathbb{R}$  we have

$$e^{\mu|x|} \leq p_\mu(x) \leq 2e^{\mu|x|} \quad \text{and} \quad p'_\mu(x) = \mu \operatorname{sign}(x) \left( e^{\mu|x|} - e^{-\mu|x|} \right)$$

and therefore we have (6.21).

Next let  $x, y \in \mathbb{R}$  such that  $|x - y| \leq \delta$ . Note that  $|p'_\mu(x)| = p'_\mu(|x|) = \mu \left( e^{\mu|x|} - e^{-\mu|x|} \right)$  and since  $|x - y| \leq \delta$  we have that  $|x| = |y| + z$  for some  $z \in (-\delta, \delta)$ . Using that  $p'(|x|)$  is monotonic in  $|x|$  we then have

$$\begin{aligned} |p'_\mu(x)| &= p'_\mu(|x|) = p'_\mu(|y| + z) \leq p'_\mu(|y| + \delta) \\ &= \mu \left( e^{\mu|y| + \mu\delta} - e^{-\mu|y| - \mu\delta} \right) = e^{\mu\delta} p'(|y|) + \mu \left( e^{\mu\delta - \mu|y|} - e^{-\mu|y| - \mu\delta} \right) \\ &\leq e^{\mu\delta} |p'(y)| + \mu e^{\mu\delta}. \end{aligned} \quad (6.24)$$

By symmetry (i.e. replacing  $x$  and  $y$ ) this implies that

$$|p'_\mu(x)| \geq e^{-\mu\delta} |p'(y)| - \mu \quad (6.25)$$

Since  $U$  is symmetric this implies that for all  $i \in [m]$  we have  $\operatorname{sign}(\nabla\Phi_\mu(\mathbf{y})^b)_i = \nabla\Phi_\mu(\mathbf{y})_i = \operatorname{sign}(y_i)$ . Therefore, if for all  $i \in [n]$  we have  $\operatorname{sign}(x_i) = \operatorname{sign}(y_i)$ , by (6.24), we see that

$$\begin{aligned} \left\langle \nabla\Phi_\mu(\mathbf{x}), \nabla\Phi_\mu(\mathbf{y})^b \right\rangle &= \sum_i p'_\mu(x_i) \nabla\Phi_\mu(\mathbf{y})_i^b \\ &\leq \sum_i \left( e^{\mu\delta} p'_\mu(y_i) + \mu e^{\mu\delta} \right) \nabla\Phi_\mu(\mathbf{y})_i^b \\ &\leq e^{\mu\delta} \left\langle \nabla\Phi_\mu(\mathbf{y}), \nabla\Phi_\mu(\mathbf{y})^b \right\rangle + \mu e^{\mu\delta} \|\nabla\Phi_\mu(\mathbf{y})^b\|_1 \\ &= e^{\mu\delta} \|\nabla\Phi_\mu(\mathbf{y})\|_U + \mu e^{\mu\delta} \|\nabla\Phi_\mu(\mathbf{y})^b\|_1. \end{aligned}$$

Similarly, using (6.25), we have  $e^{-\mu\delta} \|\nabla\Phi_\mu(\mathbf{y})\|_U - \mu \|\nabla\Phi_\mu(\mathbf{y})^b\|_1 \leq \left\langle \nabla\Phi_\mu(\mathbf{x}), \nabla\Phi_\mu(\mathbf{y})^b \right\rangle$  and hence

(6.22) holds. On the other hand if  $\text{sign}(x_i) \neq \text{sign}(y_i)$  then we know that  $|x_i| \leq \delta$  and consequently  $|p'_\mu(x_i)| \leq \mu(e^{\mu\delta} - e^{-\mu\delta}) \leq \frac{\mu}{2}$  since  $\delta \leq \frac{1}{5\mu}$ . Thus, we have

$$e^{-\mu\delta} |p'_\mu(y_i)| - \mu \leq -\frac{\mu}{2} \leq \text{sign}(y_i) p'_\mu(x_i) \leq 0 \leq e^{\mu\delta} |p'_\mu(y_i)| + \mu e^{\mu\delta}.$$

Taking inner product on both sides with  $\nabla\Phi_\mu(\mathbf{y})^\flat$  and using definition of  $\|\cdot\|_U$  and  $\cdot^\flat$ , we get (6.22). Thus, (6.22) holds in general.

Finally we note that since  $U$  is contained in a  $\ell_\infty$  ball of radius  $R$ , we have  $\|\mathbf{y}^\flat\|_1 \leq mR$  for all  $\mathbf{y}$ . Using this fact, (6.22), and the definition of  $\|\cdot\|_U$ , we obtain

$$e^{-\mu\delta} \|\nabla\Phi_\mu(\mathbf{y})\|_U - \mu mR \leq \left\langle \nabla\Phi_\mu(\mathbf{x}), \nabla\Phi_\mu(\mathbf{y})^\flat \right\rangle \leq \|\nabla\Phi_\mu(\mathbf{x})\|_U$$

where the last line comes from the fact  $\nabla\Phi_\mu(\mathbf{y})^\flat \in U$  and the definition of  $\|\cdot\|_U$ . By symmetry (6.23) follows.  $\square$

Using Lemma 6.6.7 we prove Theorem 6.6.6.

*Theorem 6.6.6.* Let  $\|\mathbf{x}\|_{U^{(k)}} = \max_{\mathbf{y} \in U^{(k)}} \langle \mathbf{x}, \mathbf{y} \rangle$  and  $\mathbf{x}^{\flat(k)} = \arg \max_{\mathbf{y} \in U^{(k)}} \langle \mathbf{x}, \mathbf{y} \rangle$ . Since  $U^{(k)}$  is symmetric, we know that  $\Delta^{(k)} = -(1 + \varepsilon) (\nabla\Phi_\mu(\mathbf{z}^{(k)}))^\flat$  and therefore by applying the mean value theorem twice we have that

$$\begin{aligned} \Phi_\mu(\mathbf{x}^{(k+1)}) &= \Phi_\mu(\mathbf{y}^{(k)}) + \left\langle \nabla\Phi_\mu(\mathbf{z}), \mathbf{x}^{(k+1)} - \mathbf{y}^{(k)} \right\rangle \\ &= \Phi_\mu(\mathbf{x}^{(k)}) + \left\langle \nabla\Phi_\mu(\mathbf{y}), \mathbf{y}^{(k)} - \mathbf{x}^{(k)} \right\rangle + \left\langle \nabla\Phi_\mu(\mathbf{z}), \mathbf{x}^{(k+1)} - \mathbf{y}^{(k)} \right\rangle \end{aligned}$$

for some  $\mathbf{y}$  between  $\mathbf{y}^{(k)}$  and  $\mathbf{x}^{(k)}$  and some  $\mathbf{z}$  between  $\mathbf{x}^{(k+1)}$  and  $\mathbf{y}^{(k)}$ . Now, using that  $\mathbf{y}^{(k)} - \mathbf{x}^{(k)} \in U^{(k)}$  and that  $\mathbf{x}^{(k+1)} - \mathbf{y}^{(k)} = \Delta^{(k)}$  we have

$$\Phi_\mu(\mathbf{x}^{(k+1)}) \leq \Phi_\mu(\mathbf{x}^{(k)}) + \|\nabla\Phi_\mu(\mathbf{y})\|_{U^{(k)}} - (1 + \varepsilon) \left\langle \nabla\Phi_\mu(\mathbf{z}), \left(\nabla\Phi_\mu(\mathbf{z}^{(k)})\right)^\flat \right\rangle. \quad (6.26)$$

Since  $U^k$  is contained within the  $\ell_\infty$  ball of radius  $R_k$  Lemma 6.6.7 shows that

$$\|\nabla\Phi_\mu(\mathbf{y})\|_{U^{(k)}} \leq e^{\mu R_k} \|\nabla\Phi_\mu(\mathbf{x}^{(k)})\|_{U^{(k)}} + m\mu R_k e^{\mu R_k}. \quad (6.27)$$

Furthermore, since  $\varepsilon < \frac{1}{5}$  and  $R_k \leq R$ , by triangle inequality we have  $\|\mathbf{z} - \mathbf{z}^{(k)}\|_\infty \leq (1 + \varepsilon)R_k + R \leq 3R$  and  $\|\mathbf{z}^{(k)} - \mathbf{x}^{(k)}\|_\infty \leq 2R$ . Therefore, applying Lemma 6.6.7 twice yields that

$$\begin{aligned} \left\langle \nabla\Phi_\mu(\mathbf{z}), \nabla\Phi_\mu(\mathbf{z}^{(k)})^\flat \right\rangle &\geq e^{-3\mu R} \|\nabla\Phi_\mu(\mathbf{z}^{(k)})\|_{U^{(k)}} - \mu mR_k \\ &\geq e^{-5\mu R} \|\nabla\Phi_\mu(\mathbf{x}^{(k)})\|_{U^{(k)}} - 2\mu mR_k. \end{aligned} \quad (6.28)$$

Combining (6.26), (6.27), and (6.28) then yields that

$$\Phi_\mu(\mathbf{x}^{(k+1)}) \leq \Phi_\mu(\mathbf{x}^{(k)}) - ((1 + \varepsilon)e^{-5\mu R} - e^{\mu R}) \|\nabla\Phi_\mu(\mathbf{x}^{(k)})\|_{U^{(k)}} + m\mu R_k e^{\mu R} + 2(1 + \varepsilon)m\mu R_k.$$

Since we chose  $\mu = \frac{\varepsilon}{12R}$ , we have  $1 + \varepsilon \leq \frac{\varepsilon}{2} + (1 + 6\mu R) \leq \frac{\varepsilon}{2} e^{5\mu R} + e^{6\mu R}$ . Hence, we have  $(1 + \varepsilon)e^{-5\mu R} - e^{\mu R} \leq \frac{\varepsilon}{2}$ . Also, since  $0 < \varepsilon < \frac{1}{5}$  we have

$$m\mu R_k e^{\mu R} + 2(1 + \varepsilon)m\mu R_k \leq (e^{\mu R} + 2(1 + \varepsilon)) m\mu R_k \leq \varepsilon m \frac{7R_k}{24R}.$$



Thus, we have

$$\Phi_\mu(\mathbf{x}^{(k+1)}) \leq \Phi_\mu(\mathbf{x}^{(k)}) - \frac{\varepsilon}{2} \|\nabla \Phi_\mu(\mathbf{x}^{(k)})\|_{U^{(k)}} + \varepsilon m \frac{7R_k}{24R}.$$

Using Lemma 6.6.7 and the fact that  $U_k$  contains a  $\ell_\infty$  ball of radius  $r_k$ , we have

$$\|\nabla \Phi_\mu(\mathbf{x}^{(k)})\|_{U^{(k)}} \geq r_k \|\nabla \Phi_\mu(\mathbf{x}^{(k)})\|_1 \geq \frac{\varepsilon r_k}{12R} \left( \Phi_\mu(\mathbf{x}^{(k)}) - 2m \right).$$

Therefore, we have that

$$\begin{aligned} \Phi_\mu(\mathbf{x}^{(k+1)}) &\leq \left(1 - \frac{\varepsilon^2 r_k}{24R}\right) \Phi_\mu(\mathbf{x}^{(k)}) + \frac{\varepsilon r_k}{12R} m + \varepsilon m \frac{7R_k}{24R} \\ &\leq \left(1 - \frac{\varepsilon^2 r_k}{24R}\right) \Phi_\mu(\mathbf{x}^{(k)}) + \varepsilon m \frac{R_k}{2R}. \end{aligned}$$

Hence, if  $\Phi_\mu(\mathbf{x}^{(k)}) \leq \frac{12m\tau}{\varepsilon}$ , we have  $\Phi_\mu(\mathbf{x}^{(k+1)}) \leq \frac{12m\tau}{\varepsilon}$ . Since  $\Phi_\mu(\mathbf{x}^{(0)}) \leq \frac{12m\tau}{\varepsilon}$  by assumption we have by induction that  $\Phi_\mu(\mathbf{x}^{(k)}) \leq \frac{12m\tau}{\varepsilon}$  for all  $k$ . The necessary bound on  $\|\mathbf{x}^{(k)}\|_\infty$  then follows immediately from Lemma 6.6.7.  $\square$

### ■ 6.6.5 Centering

Here, we show how to use the results of the previous subsection to perform weighted path following given access only to a multiplicative approximation of the weight function. To apply the results of the previous subsection, we can think updating weight is playing this game, we want to make sure the error between  $\mathbf{w}$  and  $\mathbf{g}(\mathbf{x})$  is close to 0 while the adversary control the next point  $\mathbf{g}(\mathbf{x})$  and the noise in the approximate  $\mathbf{g}(\mathbf{x})$ . Theorem 6.6.6 shows that we can control the error to be small in  $\ell_\infty$  if we can approximate  $\mathbf{g}(\mathbf{x})$  with small  $\ell_\infty$  error.

Formally, we will measure its distance from the optimal weights in log scale by

$$\Psi(\mathbf{x}, \mathbf{w}) \stackrel{\text{def}}{=} \log(\mathbf{g}(\mathbf{x})) - \log(\mathbf{w}). \quad (6.29)$$

Our goal will be to keep  $\|\Psi(\mathbf{x}, \mathbf{w})\|_{\mathbf{w}+\infty} \leq K$  for some error  $K$  that is just small enough to not impair our ability to decrease  $\delta_t$  linearly and not to impair our ability to approximate  $\mathbf{g}$ . We will attempt to do this without moving  $\mathbf{w}$  too much in  $\|\cdot\|_{\mathbf{w}+\infty}$ .

---

**Algorithm 8:**  $(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})}) = \text{centeringInexact}(\mathbf{x}, \mathbf{w}, K)$

---

Let  $c_k = \frac{1}{1-c_\delta c_\gamma}$ ,  $R = \frac{K}{48c_k \log(400m)}$ ,  $\delta_t = \delta_t(\mathbf{x}, \mathbf{w})$  and  $\varepsilon = \frac{1}{2c_k}$ .

Set  $\mathbf{x}^{(\text{new})} = \mathbf{x} - \frac{1}{\sqrt{\tilde{\phi}''(\mathbf{x})}} \mathbf{P}_{\mathbf{x}, \mathbf{w}} \left( \frac{t\mathbf{c} - \mathbf{w}\tilde{\phi}'(\mathbf{x})}{\mathbf{w}\sqrt{\tilde{\phi}''(\mathbf{x})}} \right)$ .

Let  $U = \{\mathbf{x} \in \mathbb{R}^m \mid \|\mathbf{x}\|_{\mathbf{w}+\infty} \leq \left(1 - \frac{7}{8c_k}\right) \delta_t\}$ .

Find  $\mathbf{z}$  such that  $\|\mathbf{z} - \log(\mathbf{g}(\mathbf{x}^{(\text{new})}))\|_\infty \leq R$ .

Set  $\mathbf{w}^{(\text{new})} = \exp\left(\log(\mathbf{w}) + (1 + \varepsilon) \arg \min_{\mathbf{u} \in U} \left\langle \nabla \Phi_{\frac{\varepsilon}{12R}}(\mathbf{z} - \log(\mathbf{w})), \mathbf{u} \right\rangle\right)$ .

**Output**  $(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})})$ .

---

The minimization problem in last step in `centeringInexact` is simply a projection onto the convex set  $U$  and it can be done in  $\tilde{O}(1)$  depth and  $\tilde{O}(m)$  work. See section 6.12 for details.

**Theorem 6.6.8.** Assume that  $24m^{1/4} \geq c_k \stackrel{\text{def}}{=} \frac{1}{1-c_\delta c_\gamma} \geq 5$ ,  $1 \leq C_{\text{norm}} \leq 2c_k$  and  $K \leq \frac{1}{20c_k}$ . Let

$\Psi(\mathbf{x}, \mathbf{w}) \stackrel{\text{def}}{=} \log(\mathbf{g}(\mathbf{x})) - \log(\mathbf{w})$ . Suppose that

$$\delta \stackrel{\text{def}}{=} \delta_t(\mathbf{x}, \mathbf{w}) \leq \frac{K}{48c_k \log(400m)} \quad \text{and} \quad \Phi_\mu(\Psi(\mathbf{x}, \mathbf{w})) \leq (400m)^2$$

where  $\mu = \frac{\varepsilon}{12R} = 2 \log(400m) / K$ . Let  $(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})}) = \text{centeringInexact}(\mathbf{x}, \mathbf{w}, K)$ , then

$$\delta_t(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})}) \leq \left(1 - \frac{1}{4c_k}\right) \delta \quad \text{and} \quad \Phi_\mu(\Psi(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})})) \leq (400m)^2.$$

Also, we have  $\|\log(\mathbf{g}(\mathbf{x}^{(\text{new})})) - \log(\mathbf{w})\|_\infty \leq K$ .

*Proof.* By Lemma 6.6.3, inequality (6.16),  $c_\delta c_\gamma \leq 1$  and  $c_\gamma \leq \frac{5}{4}$  (see Def 6.5.4), we have

$$\begin{aligned} \left\| \log(\mathbf{g}(\mathbf{x}^{(\text{new})})) - \log(\mathbf{g}(\mathbf{x})) \right\|_{\mathbf{g}(\mathbf{x})+\infty} &\leq c_\delta c_\gamma \delta (1 + 4c_\gamma \delta) \\ &\leq c_\delta c_\gamma \delta + 5\delta^2 \\ &\leq \left(1 - \frac{15}{16c_k}\right) \delta. \end{aligned}$$

Using  $K \leq \frac{1}{20c_k}$ , we have

$$\left\| \frac{\mathbf{w} - \mathbf{g}(\mathbf{x})}{\mathbf{g}(\mathbf{x})} \right\|_\infty \leq \|\log(\mathbf{w}) - \log(\mathbf{g}(\mathbf{x}))\|_\infty + \|\log(\mathbf{w}) - \log(\mathbf{g}(\mathbf{x}))\|_\infty^2 \leq \frac{21}{20}K \leq \frac{1}{16c_k}.$$

Hence, we have

$$\begin{aligned} \left\| \log(\mathbf{g}(\mathbf{x}^{(\text{new})})) - \log(\mathbf{g}(\mathbf{x})) \right\|_{\mathbf{w}+\infty} &\leq \left(1 + \frac{1}{16c_k}\right) \left(1 - \frac{15}{16c_k}\right) \delta \\ &\leq \left(1 - \frac{7}{8c_k}\right) \delta. \end{aligned}$$

Therefore, we know that for the Newton step, we have  $\Psi(\mathbf{x}^{(\text{new})}, \mathbf{w}) - \Psi(\mathbf{x}, \mathbf{w}) \in U$  where  $U$  is the symmetric convex set given by

$$U \stackrel{\text{def}}{=} \{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\|_{\mathbf{w}+\infty} \leq C\}$$

where  $C = \left(1 - \frac{7}{8c_k}\right) \delta$ . Note that from our assumption on  $\delta$ , we have

$$C \leq \delta \leq \frac{K}{48c_k \log(400m)} = R.$$

It ensures that  $U$  are contained in some  $\ell_\infty$  ball of radius  $R$ . Therefore, we can play the chasing 0 game on  $\Psi(\mathbf{x}, \mathbf{w})$  attempting to maintain the invariant that  $\|\Psi(\mathbf{x}, \mathbf{w})\|_\infty \leq K$  without taking steps that are more than  $1 + \varepsilon$  times the size of  $U$  where we pick  $\varepsilon = \frac{1}{2c_k}$  so to not interfere with our ability to decrease  $\delta_t$  linearly.

However, to do this with the chasing 0 game, we need to ensure that  $R$  satisfying the following

$$\frac{12R}{\varepsilon} \log\left(\frac{12m\tau}{\varepsilon}\right) \leq K$$

where here  $\tau$  is as defined in Theorem 6.6.6.

To bound  $\tau$ , we need to lower bound the radius of  $\ell_\infty$  ball it contains. Since by assumption

$\|\mathbf{g}(\mathbf{x})\|_\infty \leq 2$  and  $\|\Psi(\mathbf{x}, \mathbf{w})\|_\infty \leq \frac{1}{8}$ , we have that  $\|\mathbf{w}\|_\infty \leq 3$ . Hence, we have

$$\forall \mathbf{u} \in \mathbb{R}^m \quad : \quad \|\mathbf{u}\|_\infty^2 \geq \frac{1}{3m} \|\mathbf{u}\|_2^2.$$

Consequently, if  $\|\mathbf{u}\|_\infty \leq \frac{\delta}{5C_{\text{norm}}\sqrt{m}}$ , then  $\mathbf{u} \in U$ . So, we have that  $U$  contains a box of radius  $\frac{\delta}{5C_{\text{norm}}\sqrt{m}}$  and since  $U$  is contained in a box of radius  $\delta$ , we have that

$$\tau \leq 5C_{\text{norm}}\sqrt{m} \leq 10c_k\sqrt{m}.$$

Using  $c_k \leq 24m^{1/4}$ , we have

$$\begin{aligned} \frac{12R}{\varepsilon} \log\left(\frac{12m\tau}{\varepsilon}\right) &\leq 24c_k R \log\left(240m^{3/2}c_k^2\right) \\ &\leq 48c_k R \log(400m) = K. \end{aligned}$$

and

$$\frac{12m\tau}{\varepsilon} \leq 240m^{3/2}c_k^2 \leq (400m)^2.$$

This proves that we meet the conditions of Theorem 6.6.6. Consequently,  $\|\Psi(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})})\|_\infty \leq K$  and  $\Phi_\alpha(\Psi(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})})) \leq (400m)^2$ .

Since  $K \leq \frac{1}{4}$ , Lemma 6.6.2 shows that

$$\delta_t(\mathbf{x}^{(\text{new})}, \mathbf{w}) \leq 4(\delta_t(\mathbf{x}, \mathbf{w}))^2.$$

Our choice of  $\mathbf{w}^{(\text{new})}$  shows that

$$\begin{aligned} \left\| \log(\mathbf{w}) - \log(\mathbf{w}^{(\text{new})}) \right\|_{\mathbf{w}+\infty} &\leq \left(1 + \frac{1}{2c_k}\right) \left(1 - \frac{7}{8c_k}\right) \delta \\ &\leq \left(1 - \frac{3}{8c_k}\right) \delta. \end{aligned}$$

Using  $\delta \leq \frac{1}{80c_k}$ , the Lemma 6.6.4 shows that

$$\begin{aligned} \delta_t(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})}) &\leq \left(1 + 4\left(1 - \frac{3}{8c_k}\right)\delta\right) \left(\delta_t(\mathbf{x}^{(\text{new})}, \mathbf{w}) + \left(1 - \frac{3}{8c_k}\right)\delta\right) \\ &\leq (1 + 4\delta) \left(4\delta^2 + \left(1 - \frac{3}{8c_k}\right)\delta\right) \\ &\leq \left(1 - \frac{3}{8c_k}\right)\delta + 4\delta^2 + 16\delta^3 + 4\delta^2 \\ &\leq \left(1 - \frac{1}{4c_k}\right)\delta. \end{aligned}$$

□

## ■ 6.7 Weight Function

In this section we present the weight function that we use to achieve our  $\tilde{O}(\sqrt{\text{rank}(\mathbf{A})} \log(U/\varepsilon))$  iteration linear program solver. This weight function is motivated from the barrier introduced in

Subsection 6.3.5. We define the weight function  $\mathbf{g} : \Omega \rightarrow \mathbb{R}_{>0}^m$  for all  $\mathbf{x} \in \mathbb{R}_{>0}^m$  as follows

$$\mathbf{g}(\mathbf{x}) \stackrel{\text{def}}{=} \arg \min_{\mathbf{w} \in \mathbb{R}_{>0}^m} \hat{f}(\mathbf{x}, \mathbf{w}) \quad \text{where} \quad \hat{f}(\mathbf{x}, \mathbf{w}) \stackrel{\text{def}}{=} \mathbf{1}^T \mathbf{w} + \frac{1}{\alpha} \log \det (\mathbf{A}_x^T \mathbf{W}^{-\alpha} \mathbf{A}_x) - \beta \sum_{i \in [m]} \log w_i. \quad (6.30)$$

where here and in the remainder of the subsection we let  $\mathbf{A}_x \stackrel{\text{def}}{=} (\Phi''(\mathbf{x}))^{-1/2} \mathbf{A}$  and the parameters  $\alpha, \beta$  are chosen later such that the following hold

$$\alpha \in [1, 2), \quad \beta \in (0, 1), \quad \text{and} \quad \beta^{1-\alpha} \leq 2. \quad (6.31)$$

Here we choose  $\beta$  small and  $\alpha$  just slightly larger than 1.

We start by computing the gradient and Hessian of  $\hat{f}(\mathbf{x}, \mathbf{w})$  with respect to  $\mathbf{w}$ .

**Lemma 6.7.1.** *For all  $\mathbf{x} \in \Omega$  and  $\mathbf{w} \in \mathbb{R}_{>0}^m$ , we have*

$$\nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w}) = (\mathbf{I} - \Sigma \mathbf{W}^{-1} - \beta \mathbf{W}^{-1}) \mathbf{1} \quad \text{and} \quad \nabla_{\mathbf{w}\mathbf{w}}^2 \hat{f}(\mathbf{x}, \mathbf{w}) = \mathbf{W}^{-1} (\Sigma + \beta \mathbf{I} + \alpha \Lambda) \mathbf{W}^{-1}$$

where  $\Sigma \stackrel{\text{def}}{=} \Sigma_{\mathbf{A}_x}(\mathbf{w}^{-\alpha})$  and  $\Lambda \stackrel{\text{def}}{=} \Lambda_{\mathbf{A}_x}(\mathbf{w}^{-\alpha})$ .

*Proof.* Using Lemma 6.11.2 and the chain rule we compute the gradient of  $\nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w})$  as follows

$$\begin{aligned} \nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w}) &= \mathbf{1} + \frac{1}{\alpha} \Sigma \mathbf{W}^{\alpha} (-\alpha \mathbf{W}^{-\alpha-1}) - \beta \mathbf{W}^{-1} \mathbf{1} \\ &= (\mathbf{I} - \Sigma \mathbf{W}^{-1} - \beta \mathbf{W}^{-1}) \mathbf{1}. \end{aligned}$$

Next, using Lemma 6.11.2 and chain rule, we compute the following for all  $i, j \in [m]$

$$\begin{aligned} \frac{\partial (\nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w}))_i}{\partial w_j} &= - \frac{w_i \Lambda_{ij} w_j^{\alpha} (-\alpha w_j^{-\alpha-1}) - \Sigma_{ij} \mathbf{1}_{i=j}}{w_i^2} + \beta \mathbf{1}_{i=j} \{w_i^{-2}\} \\ &= \frac{\Sigma_{ij}}{w_i w_j} + \alpha \frac{\Lambda_{ij}}{w_i w_j} + \frac{\beta \mathbf{1}_{i=j}}{w_i^2}. \quad (\text{Using that } \Sigma \text{ is diagonal}) \end{aligned}$$

Consequently,  $\nabla_{\mathbf{w}\mathbf{w}}^2 \hat{f}(\mathbf{x}, \mathbf{w}) = \mathbf{W}^{-1} (\Sigma + \beta \mathbf{I} + \alpha \Lambda) \mathbf{W}^{-1}$  as desired.  $\square$

**Lemma 6.7.2.** *For all  $\mathbf{x} \in \Omega$ , the weight function  $\mathbf{g}(\mathbf{x})$  is a well defined with*

$$\beta \leq g_i(\mathbf{s}) \leq 1 + \beta \quad \text{and} \quad \|\mathbf{g}(\mathbf{x})\|_1 = \text{rank}(\mathbf{A}) + \beta \cdot m.$$

*Furthermore, for all  $\mathbf{x} \in \Omega$ , the weight function obeys the following equations*

$$\mathbf{G}(\mathbf{x}) = (\Sigma + \beta \mathbf{I}) \mathbf{1} \quad , \quad \text{and} \quad \mathbf{G}'(\mathbf{x}) = -\mathbf{G}(\mathbf{x}) (\mathbf{G}(\mathbf{x}) + \alpha \Lambda)^{-1} \Lambda (\Phi''(\mathbf{x}))^{-1} \Phi'''(\mathbf{x})$$

where  $\Sigma \stackrel{\text{def}}{=} \Sigma_{\mathbf{A}_x}(\mathbf{g}^{-\alpha}(\mathbf{x}))$ ,  $\Lambda \stackrel{\text{def}}{=} \Lambda_{\mathbf{A}_x}(\mathbf{g}^{-\alpha}(\mathbf{x}))$ , and  $\mathbf{G}'(\mathbf{x})$  is the Jacobian matrix of  $\mathbf{g}$  at  $\mathbf{x}$ .

*Proof.* By Lemma 6.11.2 we have that  $\Sigma \succeq \Lambda \succeq \mathbf{0}$ . Therefore, by Lemma 6.7.1, we have that  $\nabla_{\mathbf{w}\mathbf{w}}^2 \hat{f}(\mathbf{x}, \mathbf{w}) \succeq \beta \mathbf{W}^{-2}$  and  $\hat{f}(\mathbf{x}, \mathbf{w})$  is convex. Using the formula for the gradient in Lemma 6.7.1, we see that that for all  $i \in [m]$  it is the case that

$$\left[ \nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w}) \right]_i = \frac{1}{w_i} (w_i - \Sigma_{ii} - \beta).$$

Using that  $0 \leq \sigma_i \leq 1$  for all  $i$  by Lemma 6.11.2 and  $\beta \in (0, 1)$  by (6.31), we see that if  $w_i \in (0, \beta)$  then  $\left[ \nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w}) \right]_i$  is strictly negative and if  $w_i \in (1 + \beta, \infty)$  then  $\left[ \nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w}) \right]_i$  is strictly positive.

Therefore, for any  $\mathbf{x} \in \Omega$ , the  $\mathbf{w}$  that minimizes this convex function  $\hat{f}(\mathbf{x}, \mathbf{w})$  lies between the box between  $\beta$  to  $1 + \beta$ . Since  $\hat{f}$  is strongly convex in this region, the minimizer is unique.

The formula for  $\mathbf{G}(\mathbf{x})$  follows by setting  $\nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w}) = \mathbf{0}$  and the size of  $\mathbf{g}(\mathbf{x})$  follows from the fact that  $\|\boldsymbol{\sigma}\|_1 = \text{tr}(\mathbf{P}_{\mathbf{A}_x}(\mathbf{g}^{-\alpha}(\mathbf{x})))$ . Since  $\mathbf{P}_{\mathbf{A}_x}(\mathbf{g}^{-\alpha}(\mathbf{x}))$  is a projection onto the image of  $\mathbf{G}(\mathbf{x})^{-\alpha/2} \mathbf{A}_x$  and since  $\mathbf{g}(\mathbf{x}) > \mathbf{0}$  and  $\bar{\phi}''(\mathbf{x}) > \mathbf{0}$ , we have that the dimension of the image of  $\mathbf{G}(\mathbf{x})^{-\alpha/2} \mathbf{A}_x$  is the rank of  $\mathbf{A}$ . Hence, we have that  $\|\mathbf{g}(\mathbf{x})\|_1 = \text{rank}(\mathbf{A}) + \beta \cdot m$ .

By Lemma 6.11.2 and chain rule, we get the following for all  $i, j \in [m]$

$$\frac{\partial(\nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w}))_i}{\partial x_j} = -\mathbf{w}_i^{-1} \boldsymbol{\Lambda}_{ij} \bar{\phi}_j''(\mathbf{x}) \left( -(\bar{\phi}_j''(\mathbf{x}))^{-2} \bar{\phi}_j'''(\mathbf{x}) \right) = \mathbf{w}_i^{-1} \boldsymbol{\Lambda}_{ij} (\bar{\phi}_j''(\mathbf{x}))^{-1} \bar{\phi}_j'''(\mathbf{x}).$$

Consequently,  $\mathbf{J}_x(\nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w})) = \mathbf{W}^{-1} \boldsymbol{\Lambda} (\boldsymbol{\Phi}''(\mathbf{x}))^{-1} \boldsymbol{\Phi}'''(\mathbf{x})$  where  $\mathbf{J}_x$  denotes the Jacobian matrix of the function  $\nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w})$  with respect to  $\mathbf{x}$ . Since we have already know that  $\mathbf{J}_w(\nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w})) = \nabla_{\mathbf{w}\mathbf{w}}^2 \hat{f}_t(\mathbf{x}, \mathbf{w}) = \mathbf{W}^{-1} (\boldsymbol{\Sigma} + \beta \mathbf{I} + \alpha \boldsymbol{\Lambda}) \mathbf{W}^{-1}$  is positive definite (and hence invertible), by applying the implicit function theorem (Lemma 2.3.9) to the specification of  $\mathbf{g}(\mathbf{x})$  as the solution to  $\nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w}) = \mathbf{0}$ , we have

$$\mathbf{G}'(\mathbf{x}) = - \left( \mathbf{J}_w(\nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w})) \right)^{-1} \left( \mathbf{J}_x(\nabla_{\mathbf{w}} \hat{f}(\mathbf{x}, \mathbf{w})) \right) = -\mathbf{G}(\mathbf{x}) (\mathbf{G}(\mathbf{x}) + \alpha \boldsymbol{\Lambda})^{-1} \boldsymbol{\Lambda} (\boldsymbol{\Phi}''(\mathbf{x}))^{-1} \boldsymbol{\Phi}'''(\mathbf{x}).$$

□

Now we show the step consistency of  $\mathbf{g}$ .

**Lemma 6.7.3** (Step Consistency). *For all  $\mathbf{x} \in \Omega$  and  $\mathbf{y} \in \mathbb{R}^m$ , and*

$$\mathbf{B} \stackrel{\text{def}}{=} \mathbf{G}(\mathbf{x})^{-1} \mathbf{G}'(\mathbf{x}) (\boldsymbol{\phi}(\mathbf{x})'')^{-1/2},$$

we have

$$\|\mathbf{B}\mathbf{y}\|_{\mathbf{G}(\mathbf{x})} \leq \frac{2}{1+\alpha} \|\mathbf{y}\|_{\mathbf{G}(\mathbf{x})} \quad \text{and} \quad \|\mathbf{B}\mathbf{y}\|_{\infty} \leq \frac{2}{1+\alpha} \left( \|\mathbf{y}\|_{\infty} + \frac{1+2\alpha}{1+\alpha} \|\mathbf{y}\|_{\mathbf{G}(\mathbf{x})} \right).$$

Therefore

$$\|\mathbf{B}\|_{\mathbf{g}+\infty} \leq \frac{2}{1+\alpha} \left( 1 + \frac{2}{C_{\text{norm}}} \right).$$

*Proof.* Fix an arbitrary  $\mathbf{x} \in \Omega$  and let  $\mathbf{g} \stackrel{\text{def}}{=} \mathbf{g}(\mathbf{x})$ ,  $\boldsymbol{\sigma} \stackrel{\text{def}}{=} \boldsymbol{\sigma}_{\mathbf{A}_x}(\mathbf{g}^{-\alpha}(\mathbf{x}))$ ,  $\boldsymbol{\Sigma} \stackrel{\text{def}}{=} \boldsymbol{\Sigma}_{\mathbf{A}_x}(\mathbf{g}^{-\alpha}(\mathbf{x}))$ ,  $\mathbf{P} \stackrel{\text{def}}{=} \mathbf{P}_{\mathbf{A}_x}(\mathbf{g}^{-\alpha}(\mathbf{x}))$ ,  $\boldsymbol{\Lambda} \stackrel{\text{def}}{=} \boldsymbol{\Lambda}_{\mathbf{A}_x}(\mathbf{g}^{-\alpha}(\mathbf{x}))$ . Also, fix an arbitrary  $\mathbf{y} \in \mathbb{R}^m$  and let  $\mathbf{z} \stackrel{\text{def}}{=} \mathbf{B}\mathbf{y}$ .

By Lemma 6.7.2,  $\mathbf{G}' = -\mathbf{G}(\mathbf{G} + \alpha \boldsymbol{\Lambda})^{-1} \boldsymbol{\Lambda} (\boldsymbol{\Phi}'')^{-1} \boldsymbol{\Phi}'''$  and therefore

$$\begin{aligned} \mathbf{B} &= -\mathbf{G}^{-1} \left( \mathbf{G}(\mathbf{G} + \alpha \boldsymbol{\Lambda})^{-1} \boldsymbol{\Lambda} (\boldsymbol{\Phi}'')^{-1} \boldsymbol{\Phi}''' \right) (\boldsymbol{\Phi}'')^{-1/2} \\ &= (\mathbf{G} + \alpha \boldsymbol{\Lambda})^{-1} (2\boldsymbol{\Lambda}) \text{Diag} \left( \frac{-\boldsymbol{\phi}'''}{2(\boldsymbol{\phi}'')^{3/2}} \right). \end{aligned}$$

Let  $\mathbf{C} \stackrel{\text{def}}{=} (\mathbf{G} + \alpha \boldsymbol{\Lambda})^{-1} (2\boldsymbol{\Lambda})$  and let  $\mathbf{y}' \stackrel{\text{def}}{=} \text{Diag} \left( \frac{-\boldsymbol{\phi}'''}{2(\boldsymbol{\phi}'')^{3/2}} \right) \mathbf{y}$ . By the self concordance of  $\boldsymbol{\phi}$  (6.11), we know that  $\|\mathbf{y}'\| \leq \|\mathbf{y}\|$  for both  $\|\cdot\|_{\mathbf{G}}$  and  $\|\cdot\|_{\infty}$ . Since  $\mathbf{z} = \mathbf{B}\mathbf{y} = \mathbf{C}\mathbf{y}'$ , it suffices to bound  $\|\mathbf{C}\mathbf{y}'\|$  in terms of  $\|\mathbf{y}'\|$  for the necessary norms.

Letting  $\bar{\boldsymbol{\Lambda}} \stackrel{\text{def}}{=} \mathbf{G}^{-1/2} \boldsymbol{\Lambda} \mathbf{G}^{-1/2}$ , we simplify the equation further and note that

$$\|\mathbf{C}\|_{\mathbf{G}} = \|\mathbf{G}^{1/2} (\mathbf{G} + \alpha \boldsymbol{\Lambda})^{-1} (2\boldsymbol{\Lambda}) \mathbf{G}^{-1/2}\|_2 = \|(\mathbf{I} + \alpha \bar{\boldsymbol{\Lambda}})^{-1} (2\bar{\boldsymbol{\Lambda}})\|_2.$$

Now, for any eigenvector,  $\mathbf{v}$ , of  $\bar{\mathbf{A}}$  with eigenvalue  $\lambda$ , we see that  $\mathbf{v}$  is an eigenvector of  $(\mathbf{I} + \alpha\bar{\mathbf{A}})^{-1}(2\bar{\mathbf{A}})$  with eigenvalue  $2\lambda/(1 + \alpha\lambda)$ . Furthermore, since  $\mathbf{0} \preceq \bar{\mathbf{A}} \preceq \mathbf{I}$ , we have that  $\|\mathbf{C}\|_{\mathbf{G}} \leq 2/(1 + \alpha)$  and hence  $\|\mathbf{z}\|_{\mathbf{G}} \leq 2(1 + \alpha)^{-1}\|\mathbf{y}'\|_{\mathbf{G}} \leq 2(1 + \alpha)^{-1}\|\mathbf{y}\|_{\mathbf{G}}$  as desired.

To bound  $\|\mathbf{z}\|_{\infty}$ , we use that  $(\mathbf{G} + \alpha\mathbf{\Lambda})\mathbf{z} = 2\mathbf{\Lambda}\mathbf{y}'$ ,  $\mathbf{\Lambda} = \mathbf{\Sigma} - \mathbf{P}^{(2)}$ , and  $\mathbf{G} = \mathbf{\Sigma} + \beta\mathbf{I}$  to derive

$$(1 + \alpha)\mathbf{\Sigma}\mathbf{z} + \beta\mathbf{z} - \alpha\mathbf{P}^{(2)}\mathbf{z} = 2\mathbf{\Sigma}\mathbf{y}' - 2\mathbf{P}^{(2)}\mathbf{y}'.$$

Looking at the  $i^{\text{th}}$  coordinate of both sides and using that  $\sigma_i \geq 0$ , we have

$$\begin{aligned} & ((1 + \alpha)\sigma_i + \beta) |z_i| \\ & \leq \alpha \left| [\mathbf{P}^{(2)}\mathbf{z}]_i \right| + 2\sigma_i \|\mathbf{y}'\|_{\infty} + 2 \left| [\mathbf{P}^{(2)}\mathbf{y}']_i \right| \\ & \leq \alpha\sigma_i \|\mathbf{z}\|_{\mathbf{\Sigma}} + 2\sigma_i \|\mathbf{y}'\|_{\infty} + 2\sigma_i \|\mathbf{y}'\|_{\mathbf{\Sigma}} \quad (\text{Lemma 6.11.2}) \\ & \leq 2\sigma_i \|\mathbf{y}'\|_{\infty} + \sigma_i \left( \frac{2\alpha}{1 + \alpha} + 2 \right) \|\mathbf{y}'\|_{\mathbf{G}} \quad (\mathbf{\Sigma} \preceq \mathbf{G} \text{ and } \|\mathbf{z}\|_{\mathbf{G}} \leq 2(1 + \alpha)^{-1}\|\mathbf{y}'\|_{\mathbf{G}}) \end{aligned}$$

Hence, we have

$$\begin{aligned} |z_i| & \leq \frac{2}{1 + \alpha} \|\mathbf{y}'\|_{\infty} + \frac{1}{1 + \alpha} \left( \frac{2\alpha}{1 + \alpha} + 2 \right) \|\mathbf{y}'\|_{\mathbf{G}} \\ & \leq \frac{2}{1 + \alpha} [\|\mathbf{y}'\|_{\infty} + 2\|\mathbf{y}'\|_{\mathbf{G}}]. \end{aligned}$$

Therefore,  $\|\mathbf{B}\mathbf{y}\|_{\infty} = \|\mathbf{z}\|_{\infty} \leq 2(1 + \alpha)^{-1}(\|\mathbf{y}'\|_{\infty} + 2\|\mathbf{y}'\|_{\mathbf{G}})$ . Finally, we note that

$$\begin{aligned} \|\mathbf{B}\mathbf{y}\|_{\mathbf{g}_{+\infty}} & = \|\mathbf{B}\mathbf{y}\|_{\infty} + C_{\text{norm}}\|\mathbf{B}\mathbf{y}\|_{\mathbf{G}} \quad (\text{Definition}) \\ & \leq \frac{2}{1 + \alpha} \|\mathbf{y}\|_{\infty} + \frac{2}{1 + \alpha} \cdot 2\|\mathbf{y}\|_{\mathbf{G}} + \frac{2}{1 + \alpha} C_{\text{norm}}\|\mathbf{y}\|_{\mathbf{G}} \\ & \leq \frac{2}{1 + \alpha} \left( 1 + \frac{2}{C_{\text{norm}}} \right) \|\mathbf{y}\|_{\mathbf{g}_{+\infty}}. \end{aligned}$$

□

**Theorem 6.7.4.** *Choosing parameters*

$$\alpha = 1 + \frac{1}{\log_2 \left( \frac{2m}{\text{rank}(\mathbf{A})} \right)}, \quad \beta = \frac{\text{rank}(\mathbf{A})}{2m}, \quad \text{and} \quad C_{\text{norm}} = 18 \log_2 \left( \frac{2m}{\text{rank}(\mathbf{A})} \right)$$

yields

$$c_1(\mathbf{g}) = 2\text{rank}(\mathbf{A}) \quad , \quad c_{\gamma}(\mathbf{g}) = 1 + \frac{1}{9 \log_2 \left( \frac{2m}{\text{rank}(\mathbf{A})} \right)}, \quad \text{and} \quad c_{\delta}(\mathbf{g}) = 1 - \frac{2}{9 \log_2 \left( \frac{2m}{\text{rank}(\mathbf{A})} \right)}.$$

In particular, we have

$$c_{\gamma}(\mathbf{g})c_{\delta}(\mathbf{g}) \leq 1 - \frac{1}{9 \log_2 \left( \frac{2m}{\text{rank}(\mathbf{A})} \right)}.$$

*Proof.* The bounds on  $c_1(\mathbf{g})$  and  $c_{\delta}(\mathbf{g})$  follow immediately from Lemma 6.7.2 and Lemma 6.7.3. Now, we estimate the  $c_{\gamma}(\mathbf{g})$  and let  $\frac{4}{5}\mathbf{g} \leq \mathbf{w} \leq \frac{5}{4}\mathbf{g}$ . Fix an arbitrary  $\mathbf{x} \in \Omega$  and let  $\mathbf{g} \stackrel{\text{def}}{=} \mathbf{g}(\mathbf{x})$ . Recall that by Lemma 6.7.2, we have  $\mathbf{g} \geq \beta$ . Furthermore, since  $\mathbf{g}^{-1} = \mathbf{g}^{\alpha-1}\mathbf{g}^{-\alpha}$  and  $\beta^{\alpha-1} \geq \frac{1}{2}$ , the following holds

$$\frac{4}{10}\mathbf{g}_i^{-\alpha} \leq \frac{4}{5}\beta^{\alpha-1}\mathbf{g}_i^{-\alpha} \leq \frac{4}{10}\mathbf{g}_i^{-1} \leq \mathbf{w}_i^{-1} \quad (6.32)$$

for all  $i$ . Applying this and using the definition of  $\mathbf{P}_{A_x}$  yields

$$\mathbf{A}_x(\mathbf{A}_x^T \mathbf{W}^{-1} \mathbf{A}_x)^{-1} \mathbf{A}_x^T \preceq \frac{10}{4} \mathbf{A}_x(\mathbf{A}_x^T \mathbf{G}^{-\alpha} \mathbf{A}_x)^{-1} \mathbf{A}_x^T = \frac{10}{4} \mathbf{G}^{\alpha/2} \mathbf{P}_{A_x}(\mathbf{g}^{-\alpha}) \mathbf{G}^{\alpha/2}. \quad (6.33)$$

Hence, we have

$$\begin{aligned} \frac{\sigma_i\left(\frac{1}{\mathbf{w}\bar{\phi}''}\right)}{\mathbf{w}_i} &= \frac{\mathbf{1}_i^T \mathbf{A}_x(\mathbf{A}_x^T \mathbf{W}^{-1} \mathbf{A}_x)^{-1} \mathbf{A}_x^T \mathbf{1}_i}{\mathbf{w}_i^2} \\ &\leq \frac{10 \mathbf{1}_i^T \mathbf{G}^{\alpha/2} \mathbf{P}_{A_x}(\mathbf{g}^{-\alpha}) \mathbf{G}^{\alpha/2} \mathbf{1}_i}{4 \mathbf{w}_i^2} \\ &\leq \frac{10}{4} \left(\frac{5}{4}\right)^2 \frac{\sigma_i\left(\frac{1}{\mathbf{g}^\alpha \bar{\phi}''}\right)}{\mathbf{g}_i^{-2\alpha}} < 4. \end{aligned}$$

Since  $\mathbf{P}_{x,w}$  is an orthogonal projection in  $\|\cdot\|_w$ , we have  $\|\mathbf{P}_{x,w}\|_{w \rightarrow w} = 1$ . Let  $\bar{\mathbf{P}}_{x,w} \stackrel{\text{def}}{=} \mathbf{I} - \mathbf{P}_{x,w}$ , we have

$$\begin{aligned} \|\bar{\mathbf{P}}_{x,w}\|_{w \rightarrow \infty} &= \max_{i \in [m]} \max_{\|\mathbf{y}\|_w \leq 1} \mathbf{1}_i^T \bar{\mathbf{P}}_{x,w} \mathbf{y} \\ &\leq \max_{i \in [m]} \|(\mathbf{w})^{-1/2} \bar{\mathbf{P}}_{x,w}^T \mathbf{1}_i\|^2 \\ &= \max_{i \in [m]} \sqrt{\mathbf{1}_i^T \mathbf{W}^{-1} \mathbf{A}_x (\mathbf{A}_x^T \mathbf{W}^{-1} \mathbf{A}_x)^{-1} \mathbf{A}_x^T \mathbf{W}^{-1} \mathbf{1}_i}. \\ &= \max_{i \in [m]} \sqrt{\frac{\sigma_i\left(\frac{1}{\mathbf{w}\bar{\phi}''}\right)}{\mathbf{w}_i}} \leq 2. \end{aligned}$$

For any  $\mathbf{y}$ , we have

$$\begin{aligned} \|\mathbf{P}_{x,w} \mathbf{y}\|_{w+\infty} &\leq \|\mathbf{P}_{x,w} \mathbf{y}\|_\infty + C_{\text{norm}} \|\mathbf{P}_{x,w} \mathbf{y}\|_w \\ &\leq \|\mathbf{y}\|_\infty + \|\bar{\mathbf{P}}_{x,w} \mathbf{y}\|_\infty + C_{\text{norm}} \|\mathbf{y}\|_w \\ &\leq \|\mathbf{y}\|_\infty + (2 + C_{\text{norm}}) \|\mathbf{y}\|_w \\ &\leq \frac{C_{\text{norm}} + 2}{C_{\text{norm}}} \|\mathbf{y}\|_{w+\infty}. \end{aligned}$$

Hence, we have  $c_\gamma \leq \frac{C_{\text{norm}} + 2}{C_{\text{norm}}}$ . Thus, we have picked  $C_{\text{norm}} = \frac{18}{\alpha - 1}$  and have  $c_\gamma \leq 1 + \frac{\alpha - 1}{9}$ .  $\square$

### ■ 6.7.1 Weight computation

Here, we describe how to efficiently compute approximations to the weight function  $\mathbf{g} : \Omega \rightarrow \mathbb{R}_{>0}^m$  as given by (6.30). The two main technical tools we use towards this end are the *gradient descent method*, Theorem 2.3.12, a standard result in convex optimization, and fast numerical methods for estimating leverage scores using the Johnson-Lindenstrauss Lemma, Theorem 2.3.4, a powerful tool in randomized numerical linear algebra.

Since the weight function,  $\mathbf{g}$ , is defined as the minimizer of a convex optimization problem, we could use the gradient descent method directly to minimize  $\hat{f}$  and hence compute  $\mathbf{g}$ . Indeed, in Lemma 6.7.6 we show how applying the gradient descent method in a carefully scaled space allows us to compute  $\mathbf{g}(\mathbf{x})$  to high accuracy in  $\tilde{O}(1)$  iterations. Unfortunately, this result makes two assumptions to compute  $\mathbf{g}(\mathbf{x})$ : (1) we are given a weight  $\mathbf{w} \in \Omega$  that is not too far from  $\mathbf{g}(\mathbf{x})$  and (2) we compute the gradient

of  $\hat{f}$  exactly.

Assumption (1) is not an issue as we always ensure that  $\mathbf{g}$  does not change too much between calls to compute  $\mathbf{g}$  and therefore can always use our previous weights as the approximation to  $\mathbf{g}(\mathbf{x})$ . However, naively computing the gradient of  $\hat{f}$  is computationally expensive and hence assumption (2) is problematic. To deal with this issue we use the fact that by careful application of Johnson-Lindenstrauss one can compute a multiplicative approximation to the gradient efficiently and in Theorem 6.7.7 we show that this suffices to compute an approximation to  $\mathbf{g}$  that suffices to use in our weighted path following scheme.

Recall from Theorem 2.3.12 that if we take repeated projected gradient steps then we can achieve linear convergence up to bounds on how much the hessian of the function changes over the domain of interest. To apply this Theorem 2.3.12 to compute  $\mathbf{g}(\mathbf{s})$  we first need to show that there is a region around the optimal point  $\mathbf{g}(\mathbf{s})$  such that the Hessian of  $\hat{f}$  does not change too much.

**Lemma 6.7.5** (Hessian Approximation). *For  $\|\mathbf{W}_{(0)}^{-1}(\mathbf{g}(\mathbf{x}) - \mathbf{w}^{(0)})\|_\infty \leq \frac{1}{24}$  we have*

$$0.84\mathbf{W}^{-1} \preceq \nabla_{\mathbf{w}\mathbf{w}}^2 f_t(\mathbf{x}, \mathbf{w}) \preceq 3.6\mathbf{W}^{-1}.$$

*Proof.* From Lemma 6.7.1, we know that

$$\nabla_{\mathbf{w}\mathbf{w}}^2 f_t(\mathbf{x}, \mathbf{w}) = \mathbf{W}^{-1} (\boldsymbol{\Sigma} + \beta\mathbf{I} + \alpha\boldsymbol{\Lambda}) \mathbf{W}^{-1}$$

where  $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_{A_x}(\mathbf{w}^{-\alpha})$  and  $\boldsymbol{\Lambda} = \boldsymbol{\Lambda}_{A_x}(\mathbf{w}^{-\alpha})$ . Using  $\mathbf{0} \preceq \boldsymbol{\Lambda} \preceq \boldsymbol{\Sigma}$ , we have

$$\mathbf{W}^{-1} (\boldsymbol{\Sigma} + \beta\mathbf{I}) \mathbf{W}^{-1} \preceq \nabla_{\mathbf{w}\mathbf{w}}^2 f_t(\mathbf{x}, \mathbf{w}) \preceq 3\mathbf{W}^{-1} (\boldsymbol{\Sigma} + \beta\mathbf{I}) \mathbf{W}^{-1}$$

Using that  $\|\mathbf{W}_{(0)}^{-1}(\mathbf{g}(\mathbf{x}) - \mathbf{w}^{(0)})\|_\infty \leq \frac{1}{24}$ , we have

$$\boldsymbol{\Sigma} + \beta\mathbf{I} \preceq \left(1 - \frac{1}{24}\right)^{-4} \boldsymbol{\Sigma}_{A_x}(\mathbf{g}^{-\alpha}) + \beta\mathbf{I} \preceq \left(1 - \frac{1}{24}\right)^{-4} \mathbf{G} \preceq 1.19\mathbf{W}$$

and

$$\boldsymbol{\Sigma} + \beta\mathbf{I} \succeq \left(1 - \frac{1}{24}\right)^4 \boldsymbol{\Sigma}_{A_x}(\mathbf{g}^{-\alpha}) + \beta\mathbf{I} \succeq \left(1 - \frac{1}{24}\right)^4 \mathbf{G} \succeq 0.84\mathbf{W}.$$

□

Combining Theorem 2.3.12 and Lemma 6.7.5, we get the following algorithm to compute the weight function using the exact computation of the gradient of  $\hat{f}$ . Note that this algorithm applies Theorem 2.3.12 multiple times as in each iteration we are taking a gradient step with respect to a different norm.

**Lemma 6.7.6** (Exact Weight Computation). *Given  $\mathbf{w}^{(0)} \in \mathbb{R}_{>0}^m$  such that  $\|\mathbf{W}_{(0)}^{-1}(\mathbf{g}(\mathbf{x}) - \mathbf{w}^{(0)})\|_\infty \leq \frac{1}{48}$ . Let*

$$Q = \left\{ \mathbf{w} \in \mathbb{R}^m \mid \|\mathbf{W}_{(0)}^{-1}(\mathbf{g}(\mathbf{x}) - \mathbf{w}^{(0)})\|_\infty \leq \frac{1}{48} \right\}.$$

*For all  $k \geq 0$  let*

$$\mathbf{w}^{(k+1)} = \arg \min_{\mathbf{w} \in Q} \left\| \mathbf{w} - \frac{1}{2} \left( \mathbf{w}^{(k)} + \boldsymbol{\sigma}_{A_x} \left( \left( \mathbf{w}^{(k)} \right)^{-\alpha} \right) + \beta \right) \right\|_{\mathbf{W}_{(k)}^{-1}}^2$$

*This implies that for all  $k$ ,*

$$\left\| \mathbf{g}(\mathbf{x})^{-1}(\mathbf{g}(\mathbf{x}) - \mathbf{w}^{(k)}) \right\|_\infty^2 \leq 4m^2 \left(1 - \frac{1}{10}\right)^k.$$



*Proof.* Note that iterations of Theorem 2.3.12 can be rewritten as

$$\begin{aligned} \mathbf{w}^{(k+1)} &= \arg \min_{\mathbf{w} \in Q} \left\langle \left( \mathbf{I} - \Sigma_{\mathbf{A}_x} \left( \left( \mathbf{w}^{(k)} \right)^{-\alpha} \right) \mathbf{W}_{(k)}^{-1} - \beta \mathbf{W}_{(k)}^{-1} \right) \mathbf{1}, \mathbf{w} \right\rangle + \frac{4}{2} \left\| \mathbf{w} - \mathbf{w}^{(k)} \right\|_{\mathbf{W}_{(k)}^{-1}}^2 \\ &= \arg \min_{\mathbf{w} \in Q} \left\| \mathbf{w} - \frac{1}{4} \left( \mathbf{w}^{(k)} + \sigma_{\mathbf{A}_x} \left( \left( \mathbf{w}^{(k)} \right)^{-\alpha} \right) + \beta \mathbf{1} \right) \right\|_{\mathbf{W}_{(k)}^{-1}}^2 \end{aligned}$$

which is the same as in the statement of this lemma. To apply Theorem 2.3.12 we note that for any  $\mathbf{w} \in Q$  the definition of  $Q$  and the fact that  $\alpha \in (0, 1)$  implies that  $(1 - \frac{1}{48})\mathbf{W}_{(0)} \preceq \mathbf{W} \preceq (1 + \frac{1}{48})\mathbf{W}_{(0)}$ . Therefore Lemma 6.7.5 shows that for all  $\mathbf{w}^{(k)} \in Q$ ,

$$\frac{4}{5}\mathbf{W}_{(k)}^{-1} \preceq 0.84\mathbf{W}_{(0)}^{-1} \preceq \nabla_{\mathbf{w}\mathbf{w}}^2 f_t(\mathbf{x}, \mathbf{w}) \preceq 3.6\mathbf{W}_{(0)}^{-1} \preceq 4\mathbf{W}_{(k)}^{-1}. \quad (6.34)$$

Hence, Theorem 2.3.12 and inequality (6.34) shows that

$$\left\| \mathbf{w}^{(k+1)} - \mathbf{g} \right\|_{\mathbf{W}_{(k)}^{-1}}^2 \leq \left( 1 - \frac{1}{5} \right) \left\| \mathbf{w}^{(k)} - \mathbf{g} \right\|_{\mathbf{W}_{(k)}^{-1}}^2.$$

Since  $\left\| \mathbf{W}_{(0)}^{-1}(\mathbf{g}(s) - \mathbf{w}^{(0)}) \right\|_{\infty} \leq \frac{1}{48}$  and  $\mathbf{w}^{(k)} \in Q$  we know that  $\mathbf{G} \succeq (1 - \frac{1}{48})^2 \mathbf{W}_{(k)}$ . Hence, we have

$$\begin{aligned} \left\| \mathbf{w}^{(k)} - \mathbf{g} \right\|_{\mathbf{G}^{-1}}^2 &\leq \left( 1 - \frac{1}{48} \right)^{-2} \left( 1 - \frac{1}{5} \right) \left\| \mathbf{w}^{(k-1)} - \mathbf{g} \right\|_{\mathbf{G}^{-1}}^2 \\ &\leq \left( 1 - \frac{1}{10} \right) \left\| \mathbf{w}^{(k-1)} - \mathbf{g} \right\|_{\mathbf{G}^{-1}}^2 \\ &\leq \left( 1 - \frac{1}{10} \right)^k \left\| \mathbf{w}^{(0)} - \mathbf{g} \right\|_{\mathbf{G}^{-1}}^2 \end{aligned}$$

The result follows from the facts that

$$\left\| \mathbf{w}^{(0)} - \mathbf{g} \right\|_{\mathbf{G}^{-1}}^2 \leq m \|\mathbf{g}\|_{\infty} \|\mathbf{g}^{-1}(\mathbf{g} - \mathbf{w}^{(0)})\|_{\infty}^2 \leq \frac{m(1 + \beta)}{\left( 1 - \frac{1}{48} \right)^2} \left\| \mathbf{w}_{(0)}^{-1}(\mathbf{g} - \mathbf{w}^{(0)}) \right\|_{\infty}^2$$

and  $\|\mathbf{g}^{-1}(\mathbf{w}^{(k)} - \mathbf{g})\|_{\infty}^2 \leq \beta^{-1} \left\| \mathbf{w}^{(k)} - \mathbf{g} \right\|_{\mathbf{G}^{-1}}^2$  where  $\beta = \frac{\text{rank}(\mathbf{A})}{2m}$ .  $\square$

Recall that Theorem 2.3.4 showed that we can compute leverage scores,  $\sigma_{\mathbf{A}_x}$ , approximately by solving only  $\tilde{O}(1)$  many regression problems. Now, we show that we can modify Lemma 6.7.6 and use Lemma 2.3.4. Our weight computation and the analysis is as follows.

---

**Algorithm 9:**  $\mathbf{w} = \text{computeWeight}(\mathbf{x}, \mathbf{w}^{(0)}, K)$

---

Let  $\alpha = 1 + \frac{1}{\log_2 \left( \frac{2m}{\text{rank}(\mathbf{A})} \right)}$ ,  $\beta = \frac{\text{rank}(\mathbf{A})}{2m}$ ,  $\varepsilon = \frac{K}{40c_r \log \left( \frac{2m}{K} \right)}$ .

Let  $Q = \left\{ \mathbf{w} \in \mathbb{R}^m \mid \left\| \mathbf{W}_{(0)}^{-1}(\mathbf{w} - \mathbf{w}^{(0)}) \right\|_{\infty} \leq \frac{1}{48} \right\}$ .

**for**  $j = 1$  **to**  $k$  **where**  $k = \lceil 20 \log \left( \frac{4m}{K} \right) \rceil$  **do**

    Let  $\sigma^{(j)}$  be a  $(1 + \varepsilon)$  approximation of  $\sigma_{\mathbf{A}} \left( \left( \mathbf{w}^{(j)} \right)^{-\alpha} \right)$ . (Theorem 2.3.4)

    Set  $\mathbf{w}^{(j)} = \arg \min_{\mathbf{w} \in Q} \left\| \mathbf{w} - \frac{1}{4} \left( \mathbf{w}^{(j-1)} + \sigma^{(j)} + \beta \mathbf{1} \right) \right\|_{\mathbf{w}_{(j-1)}^{-1}}^2$

**end**

**Output**  $\mathbf{w}^{(j)}$ .

---

Note that the convex set  $Q$  is aligned with standard basis and hence the step on  $\mathbf{w}^{(j)}$  can be

computed by explicit formula (6.36).

**Theorem 6.7.7** (Approximate Weight Computation). *Let  $\mathbf{x} \in \Omega$ ,  $\|\mathbf{W}_{(0)}^{-1}(\mathbf{g}(\mathbf{x}) - \mathbf{w}^{(0)})\|_\infty \leq \frac{1}{48}$  and  $K \in (0, 1)$ . The algorithm `computeWeight`( $\mathbf{x}, \mathbf{w}^{(0)}, K$ ) returns  $\mathbf{w}$  such that*

$$\|\mathbf{g}(\mathbf{x})^{-1}(\mathbf{g}(\mathbf{x}) - \mathbf{w})\|_\infty \leq K$$

with probability  $1 - \frac{\tilde{O}(1)}{m}$ .

The running time is dominated by the time needed to solve  $\tilde{O}(\log^3(1/K)/K^2)$  linear systems.

*Proof.* Consider an execution of `computeWeight`( $\mathbf{x}, \mathbf{w}^{(0)}, K$ ) where  $\boldsymbol{\sigma}^{(j)} = \boldsymbol{\sigma}_{\mathbf{A}_s}((\mathbf{w}^{(j)})^{-\alpha})$  exactly, i.e.  $\boldsymbol{\sigma}^{(j)} = \boldsymbol{\sigma}_{\mathbf{A}_s}((\mathbf{w})^{-\alpha})$ , and let  $\mathbf{v}^{(j)}$  denote the  $\mathbf{w}$  computed during this idealized execution of `computeWeight`.

Now suppose that for all  $i \in [m]$  we have

$$(1 - \varepsilon)^M \mathbf{v}_i^{(j)} \leq \mathbf{w}_i^{(j)} \leq (1 + \varepsilon)^M \mathbf{v}_i^{(j)} \quad (6.35)$$

for some  $M \geq 0$  and  $j \in [k-1]$ . Since the objective function and the constraint  $Q$  are axis-aligned we can compute  $\mathbf{w}^{(j)}$  coordinate-wise and we see that

$$\mathbf{w}^{(j+1)} = \text{median} \left( \left(1 - \frac{1}{48}\right) \mathbf{w}^{(0)}, \frac{1}{4} \left( \boldsymbol{\sigma}_{\mathbf{A}_s} \left( (\mathbf{w}^{(j)})^\alpha \right) + \beta \right), \left(1 + \frac{1}{48}\right) \mathbf{w}^{(0)} \right) \quad (6.36)$$

where  $[\text{median}(\mathbf{x}, \mathbf{y}, \mathbf{z})]_i$  is equal to the median of  $x_i$ ,  $y_i$  and  $z_i$  for all  $i \in [m]$ . By (6.35), (6.36), and the fact that  $(1 - \varepsilon) \boldsymbol{\sigma}_{\mathbf{A}_x} \left( (\mathbf{w}^{(j+1)})^\alpha \right)_i \leq \boldsymbol{\sigma}_i^{(j+1)} \leq (1 + \varepsilon) \boldsymbol{\sigma}_{\mathbf{A}_x} \left( (\mathbf{w}^{(j+1)})^\alpha \right)_i$  for all  $i \in [m]$ , we have that

$$(1 - \varepsilon)^{M+1} \mathbf{v}_i^{(j+1)} \leq \mathbf{w}_i^{(j+1)} \leq (1 + \varepsilon)^{M+1} \mathbf{v}_i^{(j+1)}.$$

Since  $\mathbf{v}^{(0)} = \mathbf{w}^{(0)}$  and since  $j \in [k-1]$  was arbitrary we can apply induction and we have that for all  $j \in [k]$

$$(1 - \varepsilon)^j \mathbf{v}_i^{(j)} \leq \mathbf{w}_i^{(j)} \leq (1 + \varepsilon)^j \mathbf{v}_i^{(j)}.$$

Note that  $k\varepsilon \leq \frac{1}{8}$  and therefore by Taylor series expansion we have  $\|\mathbf{V}_{(k)}^{-1}(\mathbf{w}^{(k)} - \mathbf{v}^{(k)})\|_\infty \leq \frac{9}{8}\varepsilon k$ . Furthermore since  $\mathbf{v}^{(k)} \in Q$  we know that  $\mathbf{G}(\mathbf{s}) \succeq \left(1 - \frac{1}{24}\right) \mathbf{V}_{(k)}$ . Putting these together, applying Lemma 6.7.6, and recalling that  $k = \lceil 20 \log \left( \frac{4m}{K} \right) \rceil$  we have

$$\begin{aligned} \|\mathbf{g}(\mathbf{x})^{-1}(\mathbf{g}(\mathbf{x}) - \mathbf{w}^{(k)})\|_\infty &\leq \|\mathbf{g}(\mathbf{x})^{-1}(\mathbf{g}(\mathbf{x}) - \mathbf{v}^{(k)})\|_\infty + \|\mathbf{g}(\mathbf{x})^{-1}(\mathbf{v}^{(k)} - \mathbf{w}^{(k)})\|_\infty \\ &\leq 2m \left(1 - \frac{1}{10}\right)^{\frac{k}{2}} + \left(1 - \frac{1}{24}\right)^{-1} \|\mathbf{V}_{(k)}^{-1}(\mathbf{v}^{(k)} - \mathbf{w}^{(k)})\|_\infty \\ &\leq 2m \cdot \exp\left(-\frac{k}{20}\right) + 1.5k\varepsilon \\ &\leq \frac{K}{2} + 1.5\varepsilon \lceil 20 \log \left( \frac{4m}{K} \right) \rceil \leq K \end{aligned}$$

□

Finally, we show how to compute an initial weight without having an approximate weight to help the computation. The algorithm `computeInitialWeight`( $\mathbf{x}, K$ ) computes an initial weight in  $\tilde{O}(\sqrt{\text{rank} \mathbf{A}})$  iterations of `computeWeight` by computing  $\mathbf{g}$  for a large enough value of  $\beta$  that the

calculation is trivial and then decreasing  $\beta$  gradually.

---

**Algorithm 10:**  $\mathbf{w} = \text{computeInitialWeight}(s, K)$

---

Let  $\alpha = 1 + \frac{1}{\log_2\left(\frac{1}{\frac{2m}{\text{rank}(\mathbf{A})}\right)}$ ,  $\beta = 48$  and  $\mathbf{w} = \beta \mathbf{1}$ .

**while**  $\beta > \frac{\text{rank}(\mathbf{A})}{2m}$  **do**  
     Set  $\mathbf{w} = \text{computeWeight}(\mathbf{x}, \mathbf{w}, \frac{1}{100})$ .  
     Set  $\beta = \max \left\{ \left(1 - \frac{1}{1000\sqrt{\text{rank}(\mathbf{A})}}\right) \beta, \frac{\text{rank}(\mathbf{A})}{2m} \right\}$ .

**end**

**Output**  $\text{computeWeight}(\mathbf{x}, \mathbf{w}, K)$ .

---

**Theorem 6.7.8** (Computating Initial Weights). *For  $\mathbf{x} \in \Omega$  and  $K > 0$ , with constant probability the algorithm  $\text{computeInitialWeight}(\mathbf{x}, K)$  returns  $\mathbf{w} \in \mathbb{R}_{>0}^m$  such that*

$$\|\mathbf{g}(\mathbf{x})^{-1}(\mathbf{g}(\mathbf{x}) - \mathbf{w})\|_\infty \leq K.$$

*The total running time of  $\text{computeInitialWeight}(\mathbf{x}, K)$  is dominated by the time needed to solve  $\tilde{O}(\sqrt{\text{rank}(\mathbf{A})} \log^3(1/K)/K^2)$  linear systems.*

*Proof.* Fix  $\mathbf{x} \in \Omega$ . Let  $\mathbf{g} : \Omega \rightarrow \mathbb{R}^m$  be defined by<sup>7</sup>

$$\mathbf{g}(\beta) \stackrel{\text{def}}{=} \arg \min_{\mathbf{w} \in \mathbb{R}_{>0}^m} \mathbf{1}^T \mathbf{w} + \frac{1}{\alpha} \log \det(\mathbf{A}_x^T \mathbf{W}^{-\alpha} \mathbf{A}_x) - \beta \sum_{i \in [m]} \log w_i.$$

We first prove by induction that at the beginning of each while iteration, we have that

$$\|\mathbf{W}^{-1}(\mathbf{g}(\beta) - \mathbf{w})\|_\infty \leq \frac{1}{48}. \quad (6.37)$$

First, we prove that (6.37) holds at the first iteration. Since  $\mathbf{g}(\beta) = \boldsymbol{\sigma}(\beta) + \beta$  where  $\boldsymbol{\sigma}(\beta) \stackrel{\text{def}}{=} \boldsymbol{\sigma}_{\mathbf{A}_x}(\mathbf{g}^{-\alpha}(\beta))$ , we have that for all  $i \in [m]$

$$\beta \leq \mathbf{g}(\beta)_i \leq 1 + \beta.$$

Therefore, the initial weight,  $\mathbf{w} = \beta \mathbf{1} \in \mathbb{R}_{>0}^m$  satisfies the invariant (6.37).

Now, we assume (6.37) holds at the beginning of a certain iteration. After the step of  $\text{computeWeight}$ , by Theorem 6.7.7 we have

$$\|\mathbf{G}(\beta)^{-1}(\mathbf{g}(\beta) - \mathbf{w})\|_\infty \leq \frac{1}{100}. \quad (6.38)$$

Therefore, it suffices to prove that  $\mathbf{g}(\beta)$  is close to  $\mathbf{g}(\beta - \theta)$  for small  $\theta$ .

To bound how much  $\mathbf{g}(\beta)$  changes for small changes in  $\beta$  we proceed similarly to Lemma 6.7.2. First by the implicit function theorem (Lemma 2.3.9) and direct calculation we know that

$$\frac{d\mathbf{g}}{d\beta} = - \left( \mathbf{J}_w(\nabla_w \hat{f}(\mathbf{x}, \mathbf{w})) \right)^{-1} \left( \mathbf{J}_\beta(\nabla_w \hat{f}(\mathbf{x}, \mathbf{w})) \right) = \mathbf{G}(\beta) (\mathbf{G}(\beta) + \alpha \boldsymbol{\Lambda}_g)^{-1} \mathbf{1} \quad (6.39)$$

where  $\boldsymbol{\Lambda}_g \stackrel{\text{def}}{=} \boldsymbol{\Lambda}_{\mathbf{A}_x}(\mathbf{G}(\beta)^{-\alpha} \mathbf{1})$ . Next to estimate how fast  $\mathbf{g}$  can change as a function of  $\beta$  we estimate (6.39) in a similar manner to Lemma 6.7.3. Note that

$$\mathbf{G}(\beta) + \alpha \boldsymbol{\Lambda}_g \succeq \boldsymbol{\Sigma}(\beta)$$

---

<sup>7</sup>Note that early we assumed that  $\beta < 1$  and here we use much larger values of  $\beta$ . However, this bound on  $\beta$  was primarily to assist in bounding  $c_1$  and does not affect this proof.

where  $\Sigma(\beta) \stackrel{\text{def}}{=} \Sigma_{\mathbf{A}_s}(\mathbf{g}^\alpha(\beta))$ . Consequently,

$$\begin{aligned} \left\| \mathbf{G}(\beta)^{-1} \frac{d\mathbf{g}}{d\beta} \right\|_{\Sigma(\beta)}^2 &\leq \left\| (\mathbf{G}(\beta) + \alpha \mathbf{\Lambda}_g)^{-1} \mathbf{1} \right\|_{\Sigma(\beta)}^2 \\ &\leq \|\mathbf{1}\|_{\Sigma(\beta)}^2 = \text{rank}(\mathbf{A}). \end{aligned} \quad (6.40)$$

Using this estimate of how much  $\mathbf{g}$  changes in the  $\Sigma(\beta)$  norm, we now estimate how much  $\mathbf{g}$  changes in the  $\ell_\infty$  norm. Let  $\mathbf{z} \stackrel{\text{def}}{=} (\mathbf{G}(\beta) + \alpha \mathbf{\Lambda}_g)^{-1} \mathbf{1}$ . Then, we have

$$\begin{aligned} ((1 + \alpha) \sigma_i(\beta) + \beta) |z_i| &\leq \left| \alpha \mathbf{1}_i^T \mathbf{P}^{(2)} \mathbf{z} \right| + 1 \\ &\leq \alpha \sigma_i(\beta) \|\mathbf{z}\|_{\Sigma(\beta)} + 1. \end{aligned}$$

Using (6.40) and  $\alpha < 1$ , we have

$$\left\| \frac{d \ln \mathbf{g}}{d\beta} \right\|_{\infty} = \|\mathbf{z}\|_{\infty} \leq \max \left( \frac{\alpha \|\mathbf{z}\|_{\Sigma(\beta)}}{1 + \alpha}, \frac{1}{\beta} \right) \leq \max \left( \sqrt{\text{rank}(\mathbf{A})}, \frac{1}{\beta} \right).$$

Using (6.38), we have that

$$\|\mathbf{G}(\beta - \theta)^{-1} (\mathbf{g}(\beta - \theta) - \mathbf{w})\|_{\infty} \leq \frac{1}{48}$$

for  $\theta \leq \frac{\beta}{1000 \sqrt{\text{rank}(\mathbf{A})}}$ .

Hence, this proves the while loop preserves the invariant (6.37). Therefore, the assumptions needed for Theorem 6.7.7 throughout and `computeWeight` works as desired. Since each iteration  $\beta$  decreased by  $\tilde{O}\left(1/\sqrt{\text{rank}(\mathbf{A})}\right)$  portion and the initial  $\beta$  is  $\tilde{O}(1)$  we see that the algorithm requires only  $\tilde{O}\left(\sqrt{\text{rank}(\mathbf{A})}\right)$  iterations. Using Theorem 6.7.7 to bound the total number of linear systems solved then yields the result.  $\square$

## ■ 6.8 The Algorithm

Here we show how to use the results of previous sections to solve (6.9) using exact linear system solver. In the next section we will discuss how to relax this assumption. The central goal of this section is to develop an algorithm, `LPSolve`, for which we can prove the following theorem

**Theorem 6.8.1.** *Suppose we have an interior point  $\mathbf{x}_0 \in \Omega$  for the linear program (6.9). Then, the algorithm `LPSolve` outputs  $\mathbf{x}$  such that  $\mathbf{c}^T \mathbf{x} \leq OPT + \varepsilon$  in  $\tilde{O}\left(\sqrt{\text{rank}(\mathbf{A})} (\mathcal{T}_w + \text{nnz}(\mathbf{A})) \log(U/\varepsilon)\right)$  work and  $\tilde{O}\left(\sqrt{\text{rank}(\mathbf{A})} \mathcal{T}_d \log(U/\varepsilon)\right)$  depth where  $U = \max\left(\left\|\frac{\mathbf{u}-\mathbf{l}}{\mathbf{u}-\mathbf{x}_0}\right\|_{\infty}, \left\|\frac{\mathbf{u}-\mathbf{l}}{\mathbf{x}_0-\mathbf{l}}\right\|_{\infty}, \|\mathbf{u}-\mathbf{l}\|_{\infty}, \|\mathbf{c}\|_{\infty}\right)$  and  $\mathcal{T}_w$  and  $\mathcal{T}_d$  is the work and depth needed to compute  $(\mathbf{A}^T \mathbf{D} \mathbf{A})^{-1} \mathbf{q}$  for input positive definite diagonal matrix  $\mathbf{D}$  and vector  $\mathbf{q}$ .*

We break this proof into several parts. First we provide Lemma 6.8.2, and adaptation of a proof from [206, Thm 4.2.7] that allows us to reason about the effects of making progress along the weighted central path. Then we provide Lemma 6.8.3 that we use to bound the distance to the weighted central path in terms of centrality. After that in Lemma 6.8.3, we analyze a subroutine, `pathFollowing`, for following the weighted central path. Using these lemmas we conclude by describing our `LPSolve` algorithm and proving Theorem 6.8.1.

**Lemma 6.8.2** ([206, Theorem 4.2.7]). *Let  $\mathbf{x}^* \in \mathbb{R}^m$  denote an optimal solution to (6.9) and  $\mathbf{x}_t = \arg \min f_t(\mathbf{x}, \mathbf{w})$  for some  $t > 0$  and  $\mathbf{w} \in \mathbb{R}_{>0}^m$ . Then the following holds*

$$\mathbf{c}^T \mathbf{x}_t(\mathbf{w}) - \mathbf{c}^T \mathbf{x}^* \leq \frac{\|\mathbf{w}\|_1}{t}.$$

*Proof.* By the optimality conditions of (6.9) we know that  $\nabla_x f_t(\mathbf{x}_t(\mathbf{w})) = t \cdot \mathbf{c} + \mathbf{w} \vec{\phi}'(\mathbf{x}_t(\mathbf{w}))$  is orthogonal to the kernel of  $\mathbf{A}^T$ . Furthermore since  $\mathbf{x}_t(\mathbf{w}) - \mathbf{x}^* \in \ker(\mathbf{A}^T)$  we have

$$\left( t \cdot \mathbf{c} + \mathbf{w} \vec{\phi}'(\mathbf{x}_t(\mathbf{w})) \right)^T (\mathbf{x}_t(\mathbf{w}) - \mathbf{x}^*) = 0.$$

Using that  $\phi'_i(x_t(\mathbf{w})_i) \cdot (x_i^* - x_t(\mathbf{w})_i) \leq 1$  by Lemma 6.5.2 then yields

$$\mathbf{c}^T (\mathbf{x}_t(\mathbf{w}) - \mathbf{x}^*) = \frac{1}{t} \sum_{i \in [m]} w_i \cdot \phi'_i(x_t(\mathbf{w})_i) \cdot (x_i^* - x_t(\mathbf{w})_i) \leq \frac{\|\mathbf{w}\|_1}{t}.$$

□

**Lemma 6.8.3.** *For  $\delta_t(\mathbf{x}^{(1)}, \mathbf{g}(\mathbf{x}^{(1)})) \leq \frac{1}{960c_k^2 \log(400m)}$  and  $\mathbf{x}_t \stackrel{\text{def}}{=} \arg \min f_t(\mathbf{x}, \mathbf{w})$  we have*

$$\left\| \sqrt{\vec{\phi}''(\mathbf{x}_t)} \left( \mathbf{x}^{(1)} - \mathbf{x}_t \right) \right\|_{\infty} \leq 16c_{\gamma} c_k \delta_t(\mathbf{x}^{(1)}, \mathbf{g}(\mathbf{x}^{(1)})).$$

*Proof.* We use Theorem 6.6.8 with exact weight computation and start with  $\mathbf{x}^{(1)}$  and  $\mathbf{w}^{(1)} = \mathbf{g}(\mathbf{x}^{(1)})$ . In each iteration,  $\delta_t$  is decreased by a factor of  $\left(1 - \frac{1}{4c_k}\right)$ . (6.16) shows that

$$\left\| \sqrt{\vec{\phi}''(\mathbf{x}^{(k)})} \left( \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} \right) \right\|_{\infty} \leq c_{\gamma} \delta_t(\mathbf{x}^{(k)}, \mathbf{w}^{(k)}).$$

The Lemma 6.5.1 shows that

$$\begin{aligned} \left\| \log \left( \phi''(\mathbf{x}^{(k)}) \right) - \log \left( \phi''(\mathbf{x}^{(k+1)}) \right) \right\|_{\infty} &\leq \left( 1 - 2c_{\gamma} \delta_t(\mathbf{x}^{(k)}, \mathbf{w}^{(k)}) \right)^{-1} \\ &\leq e^{4c_{\gamma} \delta_t(\mathbf{x}^{(k)}, \mathbf{w}^{(k)})}. \end{aligned}$$

Therefore, for any  $k$ , we have

$$\left\| \log \left( \phi''(\mathbf{x}^{(1)}) \right) - \log \left( \phi''(\mathbf{x}^{(k)}) \right) \right\|_{\infty} \leq e^{4c_{\gamma} \sum \delta_t(\mathbf{x}^{(k)}, \mathbf{w}^{(k)})} \leq e^{32c_k c_{\gamma} \delta_t(\mathbf{x}^{(1)}, \mathbf{g}(\mathbf{x}^{(1)}))} \leq 2.$$

Hence, for any  $k$ , we have

$$\begin{aligned} \left\| \sqrt{\vec{\phi}''(\mathbf{x}_t)} \left( \mathbf{x}^{(1)} - \mathbf{x}^{(k)} \right) \right\|_{\infty} &\leq \sum 2c_{\gamma} \delta_t(\mathbf{x}^{(k)}, \mathbf{w}^{(k)}) \\ &\leq 16c_{\gamma} c_k \delta_t(\mathbf{x}^{(1)}, \mathbf{g}(\mathbf{x}^{(1)})). \end{aligned}$$

It is clear now  $\mathbf{x}^{(k)}$  forms a Cauchy sequence and converges to  $\mathbf{x}_t$  because  $\delta_t$  continuous and  $\mathbf{x}_t$  is the unique point such that  $\delta_t = 0$ . □

Next, we put together the results of Section 6.6 and analyze the following algorithm for following

the weighted central path.

---

**Algorithm 11:**  $(\mathbf{x}^{(\text{final})}, \mathbf{w}^{(\text{new})}) = \text{pathFollowing}(\mathbf{x}, \mathbf{w}, t_{\text{start}}, t_{\text{end}}, \varepsilon)$

---

Let  $c_k = 9 \log_2 \left( \frac{2m}{\text{rank}(\mathbf{A})} \right)$ ,  $t = t_{\text{start}}$ ,  $K = \frac{1}{20c_k}$ .

**while**  $(t < t_{\text{end}}$  if  $t_{\text{start}} < t_{\text{end}}$ ) or  $(t > t_{\text{end}}$  if  $t_{\text{start}} > t_{\text{end}})$  **do**

    Set  $(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})}) = \text{centeringInexact}(\mathbf{x}, \mathbf{w}, K)$  where it uses the function `computeWeight` to find the approximation of  $\mathbf{g}(\mathbf{x})$ .

    Set  $t \leftarrow t \left( 1 \pm \frac{1}{10^5 c_k^4 \log(400m) \sqrt{\text{rank}(\mathbf{A})}} \right)$  where the sign of  $\pm$  is the sign of  $t_{\text{end}} - t_{\text{start}}$ .

    Set  $\mathbf{x} \leftarrow \mathbf{x}^{(\text{final})}$  and  $\mathbf{w} \leftarrow \mathbf{w}^{(\text{new})}$ .

**end**

**for**  $i = 1, \dots, 4c_k \log(1/\varepsilon)$  **do**

    Set  $(\mathbf{x}, \mathbf{w}) = \text{centeringInexact}(\mathbf{x}, \mathbf{w}, K)$  where it uses the function `computeWeight` to find the approximation of  $\mathbf{g}(\mathbf{x})$ .

**end**

**Output**  $(\mathbf{x}, \mathbf{w})$ .

---

**Theorem 6.8.4.** *Suppose that*

$$\delta_{t_{\text{start}}}(\mathbf{x}, \mathbf{w}) \leq \frac{1}{960c_k^2 \log(400m)} \quad \text{and} \quad \Phi_\mu(\Psi(\mathbf{x}, \mathbf{w})) \leq (400m)^2.$$

where  $\mu = 2 \log(400m) / K$ . Let  $(\mathbf{x}^{(\text{final})}, \mathbf{w}^{(\text{new})}) = \text{pathFollowing}(\mathbf{x}, \mathbf{w}, t_{\text{start}}, t_{\text{end}})$ , then

$$\delta_{t_{\text{end}}}(\mathbf{x}^{(\text{final})}, \mathbf{w}^{(\text{new})}) \leq \varepsilon \quad \text{and} \quad \Phi_\mu(\Psi(\mathbf{x}^{(\text{final})}, \mathbf{w}^{(\text{new})})) \leq (400m)^2.$$

Furthermore,  $\text{pathFollowing}(\mathbf{x}, \mathbf{w}, t_{\text{start}}, t_{\text{end}})$  takes time

$$\tilde{O} \left( \sqrt{\text{rank}(\mathbf{A})} \left( \left| \log \left( \frac{t_{\text{end}}}{t_{\text{start}}} \right) \right| + \log(1/\varepsilon) \right) (\mathcal{T} + m) \right)$$

where  $\mathcal{T}$  is the time needed to solve on linear system.

*Proof.* We use induction to prove the invariant that  $\delta_t(\mathbf{x}, \mathbf{w}) \leq \frac{1}{960c_k^2 \log(400m)}$  and  $\Phi_\alpha(\Psi(\mathbf{x}, \mathbf{w})) \leq (400m)^2$  on the beginning of while iteration. The initial condition holds by assumption.

Now, assume the invariant holds at a certain iteration. Theorem 6.6.8 shows that

$$\|\log(\mathbf{g}(\mathbf{x}^{(\text{new})})) - \log(\mathbf{w})\|_\infty \leq K \leq \frac{1}{20c_k}. \quad (6.41)$$

Thus, the weight satisfies the condition of Theorem 6.7.7 and the algorithm `centeringInexact` can use the function `computeWeight` to find the approximation of  $\mathbf{g}(\mathbf{x}^{(\text{new})})$ . Consequently,

$$\delta_t(\mathbf{x}^{(\text{final})}, \mathbf{w}^{(\text{new})}) \leq \left( 1 - \frac{1}{4c_k} \right) \delta_t \quad \text{and} \quad \Phi_\alpha(\Psi(\mathbf{x}^{(\text{final})}, \mathbf{w}^{(\text{new})})) \leq (400m)^2.$$

Using Lemma 6.6.1, (6.41) and Theorem 6.7.4, we have

$$\delta_t(\mathbf{x}^{(\text{final})}, \mathbf{w}^{(\text{new})}) \leq \frac{1}{960c_k^2 \log(400m)}.$$

Hence, we proved that the invariant holds after one iteration. The  $\delta_t < \varepsilon$  bounds follows from the last loop.  $\square$

---

**Algorithm 12:**  $\mathbf{x}^{(\text{final})} = \text{LPSolve}(\mathbf{x}, \varepsilon)$

---

Let  $\beta = \frac{\text{rank}(\mathbf{A})}{2m}$ ,  $\mathbf{w} = \text{computeInitialWeight}(\mathbf{x}, \frac{1}{10^5 \log^5(400m)})$ ,  $\mathbf{d} = -\mathbf{w}_i \phi'_i(\mathbf{x})$ .

Let  $t_1 = (10^{10} U^2 m^3)^{-1}$ ,  $t_2 = 3m/\varepsilon$ ,  $\varepsilon_1 = \frac{1}{2000c_k^2 \log(400m)}$ ,  $\varepsilon_2 = \frac{\varepsilon}{100^3 m^3 U^2}$ .

Set  $(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})}) = \text{pathFollowing}(\mathbf{x}, \mathbf{w}, 1, t_1, \varepsilon_1)$  with cost vector  $\mathbf{d}$ .

Set  $(\mathbf{x}^{(\text{final})}, \mathbf{w}^{(\text{final})}) = \text{pathFollowing}(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})}, t_1, t_2, \varepsilon_2)$  with cost vector  $\mathbf{c}$ .

**Output**  $\mathbf{x}^{(\text{final})}$ .

---

*Proof of Theorem 6.8.1.* By Theorem 6.7.8, we know the initial weight satisfies

$$\|\mathbf{G}(\mathbf{x})^{-1}(\mathbf{g}(\mathbf{x}) - \mathbf{w})\|_\infty \leq \frac{1}{10^5 \log^5(400m)}.$$

By the definition of  $\mathbf{d}$ , we have  $\mathbf{x}$  is the minimum of  $\min \mathbf{d}^T \mathbf{x} - \sum \mathbf{w}_i \phi_i(\mathbf{x})$  given  $\mathbf{A}^T \mathbf{x} = \mathbf{b}$ . Therefore,  $(\mathbf{x}, \mathbf{w})$  satisfies the assumption of Theorem 6.8.4 because  $\delta_t = 0$  and  $\Phi_\alpha$  is small enough. Hence, we have

$$\delta_{t_1}(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})}) \leq \frac{1}{2000c_k^2 \log(400m)} \quad \text{and} \quad \Phi_\alpha(\Psi(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})})) \leq (400m)^2.$$

Lemma 6.5.2 shows that  $\|\phi'_i(\mathbf{x})\|_\infty \leq U$  and hence  $\|\mathbf{c} - \mathbf{d}\|_\infty \leq 2U$ . Also, Lemma 6.5.1 shows that  $\min_{\mathbf{y}} \sqrt{\vec{\phi}''(\mathbf{y})} \geq \frac{1}{U}$ . Therefore, we have

$$\begin{aligned} \delta_{t_1}^{\mathbf{c}}(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})}) &= \min_{\boldsymbol{\eta} \in \mathbb{R}^n} \left\| \frac{t_1 \mathbf{c} + \mathbf{w} \vec{\phi}'(\mathbf{x}^{(\text{new})}) - \mathbf{A} \boldsymbol{\eta}}{\mathbf{w}^{(\text{new})} \sqrt{\vec{\phi}''(\mathbf{x}^{(\text{new})})}} \right\|_{\mathbf{w}^{(\text{new})} + \infty} \\ &\leq \min_{\boldsymbol{\eta} \in \mathbb{R}^n} \left\| \frac{t_1 \mathbf{d} + \mathbf{w} \vec{\phi}'(\mathbf{x}^{(\text{new})}) - \mathbf{A} \boldsymbol{\eta}}{\mathbf{w}^{(\text{new})} \sqrt{\vec{\phi}''(\mathbf{x}^{(\text{new})})}} \right\|_{\mathbf{w}^{(\text{new})} + \infty} + t_1 \left\| \frac{\mathbf{c} - \mathbf{d}}{\mathbf{w}^{(\text{new})} \sqrt{\vec{\phi}''(\mathbf{x}^{(\text{new})})}} \right\|_{\mathbf{w}^{(\text{new})} + \infty} \\ &\leq \delta_{t_1}^{\mathbf{d}}(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})}) + 4U^2 t_1 \|\mathbf{1}\|_{\mathbf{w} + \infty} \\ &= \delta_{t_1}^{\mathbf{d}}(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})}) + 100mU^2 t_1. \end{aligned}$$

Since we have chosen  $t_1$  small enough, we have  $\delta_{t_1}^{\mathbf{c}}(\mathbf{x}^{(\text{new})}, \mathbf{w}^{(\text{new})})$  is small enough to satisfy the assumption of Theorem 6.8.4. So, we only need to prove how large  $t_2$  should be and how small  $\varepsilon_2$  should be in order to get  $\mathbf{x}$  such that  $\mathbf{c}^T \mathbf{x} \leq \text{OPT} + \varepsilon$ . By Lemma 6.8.2 and  $\|\mathbf{w}^{(\text{final})}\| \leq 3m$ , we have

$$\mathbf{c}^T \mathbf{x}_{t_2} \leq \text{OPT} + \frac{3m}{t_2}.$$

Also, Lemma 6.8.3 shows that we have

$$\left\| \sqrt{\vec{\phi}''(\mathbf{x}_{t_2})} (\mathbf{x}^{(\text{final})} - \mathbf{x}_{t_2}) \right\|_\infty \leq 32\varepsilon_2 c_k.$$

Using  $\min_{\mathbf{y}} \sqrt{\vec{\phi}''(\mathbf{y})} \geq \frac{1}{U}$ , we have  $\|\mathbf{x}^{(\text{final})} - \mathbf{x}_{t_2}\|_\infty \leq 32\varepsilon_2 c_k U$  and hence our choice of  $t_2$  and  $\varepsilon_2$  gives the result

$$\mathbf{c}^T \mathbf{x}^{(\text{final})} \leq \text{OPT} + \frac{3m}{t_2} + 32\varepsilon_2 c_k U^2 \leq \text{OPT} + \varepsilon.$$

□

## ■ 6.9 Generalized Minimum Cost Flow

In this section we show how to use the interior point method in Section 6.8 to solve the maximum flow problem in time  $\tilde{O}(m\sqrt{n}\log^{O(1)}(U))$ , to solve the minimum cost flow problem in time,  $\tilde{O}(m\sqrt{n}\log^{O(1)}(U))$ , and to compute  $\varepsilon$ -approximate solutions to the lossy generalized minimum cost flow problem in time  $\tilde{O}(m\sqrt{n}\log^{O(1)}(U/\varepsilon))$ . Our algorithm for the generalized minimum cost flow problem is essentially the same as our algorithm for the simpler specific case of minimum cost flow and maximum flow and therefore, we present the algorithm for the generalized minimum cost flow problem directly.

The generalized minimum cost flow problem [63] is as follows. Let  $G = (V, E)$  be a connected directed graph where each edge  $e$  has capacity  $c_e > 0$  and multiplier  $1 \geq \gamma_e > 0$ . For each edge  $e$ , there can be only at most  $c_e$  units of flow on that edge and the flow on that edge must be non-negative. Also, for each unit of flow entering edge  $e$ , there are only  $\gamma_e$  units of flow going out. The generalized maximum flow problem is to compute how much flow can be sent into  $t$  given a unlimited source  $s$ . The generalized minimum cost flow is to ask what is the minimum cost of sending the maximum flow given the cost of each edge is  $q_e$ . The maximum flow and the minimum cost flow are the case with  $\gamma_e = 1$  for all edges  $e$ .

Since the generalized minimum cost flow includes all of these cases, we focus on this general formulation. The problem can be written as the following linear program

$$\min_{\mathbf{0} \leq \mathbf{x} \leq \mathbf{c}} \mathbf{q}^T \mathbf{x} \text{ such that } \mathbf{A}\mathbf{x} = F\mathbf{1}_t$$

where  $F$  is the generalized maximum flow value,  $\mathbf{1}_t$  is a indicator vector of size  $(n-1)$  that is non-zero at vertices  $t$  and  $\mathbf{A}$  is a  $|V \setminus \{s\}| \times |E|$  matrix such that for each edge  $e$ , we have

$$\begin{aligned} \mathbf{A}(e_{head}, e) &= \gamma(e), \\ \mathbf{A}(e_{tail}, e) &= -1. \end{aligned}$$

In order words, the constraint  $\mathbf{A}\mathbf{x} = F\mathbf{1}_t$  requires the flow to satisfies the flow conservation at all vertices except  $s$  and  $t$  and requires it flows  $F$  unit of flow to  $t$ . We assume  $c_e$  are integer and  $\gamma_e$  is a rational number. Let  $U$  be the maximum of  $c_e$ ,  $q_e$ , the numerator of  $\gamma_e$  and the denominator of  $\gamma_e$ . For the generalized flow problems, getting an efficient exact algorithm is difficult and we aim for approximation algorithms only.

**Definition 6.9.1.** We call a flow an  $\varepsilon$ -approximate generalized maximum flow if it is a flow satisfies the flow conservation and the flow value is larger than maximum flow value minus  $\varepsilon$ . We call a flow is an  $\varepsilon$ -approximate generalized minimum cost maximum flow if it is an  $\varepsilon$ -approximate maximum flow and has cost not greater than the minimum cost maximum flow value.

Note that  $\text{rank}(\mathbf{A}) = n - 1$  because the graph is connected and hence our algorithm takes only  $\tilde{O}(\sqrt{n}L)$  iterations. Therefore, the problems remaining are to compute  $L$  and bound how much time is required to solve the linear systems involved. However,  $L$  is large in the most general setting and hence we cannot use the standard theory to say how to get the initial point, how to round to the vertex. Furthermore, the condition number of  $\mathbf{A}^T \mathbf{A}$  can be very bad.

In [63], they used dual path following to solve the generalized minimum cost flow problem with the caveats that the dual polytope is not bounded, the problem of getting the initial flow, the problem of rounding it to the a feasible flow. We use there analysis to formulate the problem in a manner more amenable to our algorithms. Since we are doing the primal path following, we will state a reformulation of the LP slightly different.



**Theorem 6.9.2** ([63]). *Given a directed graph  $G$ . We can find a new directed graph  $\tilde{G}$  with  $O(m)$  edges and  $O(n)$  vertices in  $\tilde{O}(m)$  time such that the modified linear program*

$$\min_{0 \leq x_i \leq c_i, 0 \leq y_i \leq 4mU^2, 0 \leq z_i \leq 4mU^2} \mathbf{q}^T \mathbf{x} + \frac{256m^5U^5}{\varepsilon^2} (\mathbf{1}^T \mathbf{y} + \mathbf{1}^T \mathbf{z}) \quad \text{such that } \mathbf{A}\mathbf{x} + \mathbf{y} - \mathbf{z} = F\mathbf{1}_t$$

satisfies the following conditions:

1.  $\mathbf{x} = \frac{\varepsilon}{2}\mathbf{1}$ ,  $\mathbf{y} = 2mU^2\mathbf{1} - (\mathbf{A}\frac{\varepsilon}{2}\mathbf{1})^- + F\vec{\mathbf{1}}_t$ ,  $\mathbf{z} = 2mU^2\mathbf{1} + (\mathbf{A}\frac{\varepsilon}{2}\mathbf{1})^+$  is an interior point of the linear program.
2. Given any  $(\mathbf{x}, \mathbf{y}, \mathbf{z})$  such that  $\|\mathbf{A}\mathbf{x} + \mathbf{y} - \mathbf{z}\|_2 \leq \frac{\varepsilon^2}{128m^2n^2U^3}$  and with cost value within  $\frac{\varepsilon^2}{128m^2n^2U^3}$  of the optimum. Then, one can compute an  $\varepsilon$ -approximate minimum cost maximum flow for graph  $G$  in time  $\tilde{O}(m)$ .
3. The linear system of the linear program is well-conditioned, i.e., the condition number of

$$\begin{bmatrix} \mathbf{A} & \mathbf{I} & -\mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A}^T \\ \mathbf{I} \\ -\mathbf{I} \end{bmatrix}$$

is  $O(mU)$ .

4. The linear system of the linear program can be solve in nearly linear time, i.e. for any diagonal matrix  $\mathbf{S}$  with condition number  $\kappa$  and vector  $b$ , it takes  $\tilde{O}(m \log(\frac{\kappa U}{\delta}))$  time to find  $x$  such that

$$\|x - \mathbf{L}^{-1}b\|_{\mathbf{L}} \leq \delta \|x\|_{\mathbf{L}}$$

$$\text{where } \mathbf{L} = \begin{bmatrix} \mathbf{A} & \mathbf{I} & -\mathbf{I} \end{bmatrix} \mathbf{S} \begin{bmatrix} \mathbf{A}^T \\ \mathbf{I} \\ -\mathbf{I} \end{bmatrix}.$$

The main difference between what stated in [63] and here is that

1. Our linear program solver can support constraint  $l_i \leq x_i \leq u_i$  and hence we do not need to split the flow variable to positive part and negative part.
2. Our linear program solver is primal path following and hence we add the constraint  $y_i \leq 4mU^2$  and  $z_i \leq 4mU^2$ . Since the maximum flow value is at most  $mU^2$ , it does not affect the optimal solution of the linear program.
3. We remove the variable  $\mathbf{x}_3$  in [63] because the purpose of that is to make the dual polytope is bounded and we do not need it here.

Using the reduction mentioned above, one can obtain the promised generalized minimum cost flow algorithm.

**Theorem 6.9.3.** *There is a randomized algorithm to compute an  $\varepsilon$ -approximate generalized minimum cost maximum flow in  $\tilde{O}(\sqrt{n} \log^{O(1)}(U/\varepsilon))$  depth  $\tilde{O}(m\sqrt{n} \log^{O(1)}(U/\varepsilon))$  total work (see Definition 6.9.1). Furthermore, there is an algorithm to compute an exact standard minimum cost maximum flow in  $\tilde{O}(\sqrt{n} \log^{O(1)}(U))$  depth and  $\tilde{O}(m\sqrt{n} \log^{O(1)}(U))$  total work.*

*Proof.* Using the reduction above and Theorem 6.8.1, we get an algorithm of generalized minimum cost flow by solving  $\tilde{O}(\sqrt{n})$  linear systems to  $\tilde{O}(1)$  bit accuracy and the condition number of those systems

are  $\text{poly}(mU/\varepsilon)$ . In [63], they showed that the linear system involved can be reduced to  $\tilde{O}(\log(U/\varepsilon))$  many Laplacian systems and hence we can use a recent nearly linear work polylogarithmic depth Laplacian system solver of Spielman and Peng [223]. In total, it takes  $\tilde{O}(m \log^{O(1)}(\frac{U}{\varepsilon}))$  time to solve each systems.

For the standard minimum cost maximum flow problem, it is known that the solution set is a convex polytope with integer coordinates and we can use Isolation lemma to make sure there is unique minimum. Hence, we only need to take  $\varepsilon = \text{poly}(1/mU)$  and round the solution to the closest integer. See Section 3.5 in [63] for details.  $\square$

## ■ 6.10 Glossary

Here we summarize various linear programming specific notation that we use throughout the chapter. For many quantities we included the typical order of magnitude as they appear during our algorithms.

- Linear program related: constraint matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , cost vector  $\mathbf{c} \in \mathbb{R}^m$ , constraint vector  $\mathbf{b} \in \mathbb{R}^n$ , solution  $\mathbf{x} \in \mathbb{R}^m$ , weights of constraints  $\mathbf{w} \in \mathbb{R}^m$  where  $m$  is the number of variables and  $n$  is the number of constraints.
- Matrix version of variables:  $\mathbf{S}$  is the diagonal matrix corresponds to  $\mathbf{s}$ ,  $\mathbf{W}$  corresponds to  $\mathbf{w}$ ,  $\Phi$  corresponds to  $\phi$ .
- Penalized objective function (6.12):  $f_t(\mathbf{x}, \mathbf{w}) = t \cdot \mathbf{c}^T \mathbf{x} + \sum_{i \in [m]} w_i \phi_i(\mathbf{x}_i)$ .
- Barrier functions (Sec 6.5.1): For  $[l, \infty)$ , we use  $\phi(x) = -\log(x - l)$ . For  $(-\infty, u]$ , we use  $\phi(x) = -\log(u - x)$ . For  $[l, u]$ , we use  $\phi(x) = -\log(ax + b)$  where  $a = \frac{\pi}{u-l}$  and  $b = -\frac{\pi}{2} \frac{u+l}{u-l}$ .
- The projection matrix  $\mathbf{P}_{\mathbf{x}, \mathbf{w}}$  (6.14):  $\mathbf{P}_{\mathbf{x}, \mathbf{w}} = \mathbf{I} - \mathbf{W}^{-1} \mathbf{A}_x (\mathbf{A}_x^T \mathbf{W}^{-1} \mathbf{A}_x)^{-1} \mathbf{A}_x^T$  where  $\mathbf{A}_x \stackrel{\text{def}}{=} \Phi''(\mathbf{x})^{-1/2} \mathbf{A}$ .
- Newton step (6.13):  $\mathbf{h}_t(\mathbf{x}, \mathbf{w}) = -\Phi''(\mathbf{x})^{-1/2} \mathbf{P}_{\mathbf{x}, \mathbf{w}} \mathbf{W}^{-1} \Phi''(\mathbf{x})^{-1/2} \nabla_x f_t(\mathbf{x}, \mathbf{w})$ .
- The mixed norm (6.15):  $\|\mathbf{y}\|_{\mathbf{w}+\infty} = \|\mathbf{y}\|_\infty + C_{\text{norm}} \|\mathbf{y}\|_{\mathbf{W}}$  where  $C_{\text{norm}} \approx \text{polylog}(m)$ .
- Centrality (6.17):  $\delta_t(\mathbf{x}, \mathbf{w}) = \min_{\boldsymbol{\eta} \in \mathbb{R}^n} \left\| \frac{\nabla_x f_t(\mathbf{x}, \mathbf{w}) - \mathbf{A} \boldsymbol{\eta}}{\mathbf{w} \sqrt{\Phi''(\mathbf{x})}} \right\|_{\mathbf{w}+\infty} \approx \frac{1}{\text{polylog}(m)}$ .
- Properties of weight function (Def 6.5.4): **size**  $c_1(\mathbf{g}) = \|\mathbf{g}(\mathbf{x})\|_1 \approx \text{rank}(\mathbf{A})$ , **slack sensitivity**  $c_\gamma(\mathbf{g}) = \|\mathbf{P}_{\mathbf{x}, \mathbf{w}}\|_{\mathbf{w}+\infty} \approx 1 + \frac{1}{\text{polylog}(m)}$ , **step consistency**  $c_\delta(\mathbf{g}) \approx 1 - \frac{1}{\text{polylog}(m)}$ .
- Difference between  $\mathbf{g}$  and  $\mathbf{w}$  (6.29):  $\Psi(\mathbf{x}, \mathbf{w}) = \log(\mathbf{g}(\mathbf{x})) - \log(\mathbf{w})$ .
- Potential function for tracing 0 (Thm 6.6.6):  $\Phi_\mu(\mathbf{x}) = e^{\mu x} + e^{-\mu x} \approx \text{poly}(m)$ .
- The weight function proposed (6.30):

$$\mathbf{g}(\mathbf{x}) = \arg \min_{\mathbf{w} \in \mathbb{R}_{>0}^m} \hat{f}(\mathbf{x}, \mathbf{w}) \quad \text{where} \quad \hat{f}(\mathbf{x}, \mathbf{w}) = \mathbf{1}^T \mathbf{w} + \frac{1}{\alpha} \log \det(\mathbf{A}_x^T \mathbf{W}^{-\alpha} \mathbf{A}_x) - \beta \sum_i \log w_i$$

where  $\mathbf{A}_x = (\Phi''(\mathbf{x}))^{-1/2} \mathbf{A}$ ,  $\alpha \approx 1 + 1/\log_2\left(\frac{m}{\text{rank}(\mathbf{A})}\right)$ ,  $\beta \approx \text{rank}(\mathbf{A})/m$ .

## ■ 6.11 Appendix: Technical Lemmas

**Lemma 6.11.1.** For any norm  $\|\cdot\|$  and  $\|\mathbf{y}\|_Q \stackrel{\text{def}}{=} \min_{\boldsymbol{\eta} \in \mathbb{R}^n} \left\| \mathbf{y} - \frac{A\boldsymbol{\eta}}{w\sqrt{\vec{\phi}''(\mathbf{x})}} \right\|$ , we have

$$\|\mathbf{y}\|_Q \leq \|\mathbf{P}_{x,w}\mathbf{y}\| \leq \|\mathbf{P}_{x,w}\| \cdot \|\mathbf{y}\|_Q.$$

*Proof.* By definition  $\mathbf{P}_{x,w}\mathbf{y} = \mathbf{y} - \frac{A\boldsymbol{\eta}_y}{w\sqrt{\vec{\phi}''(\mathbf{x})}}$  for some  $\boldsymbol{\eta}_y \in \mathbb{R}^n$ . Consequently,

$$\|\mathbf{y}\|_Q = \min_{\boldsymbol{\eta} \in \mathbb{R}^n} \left\| \mathbf{y} - \frac{A\boldsymbol{\eta}}{w\sqrt{\vec{\phi}''(\mathbf{x})}} \right\| \leq \|\mathbf{P}_{x,w}\mathbf{y}\|.$$

On the other hand, let  $\boldsymbol{\eta}_q$  such that  $\|\mathbf{y}\|_Q = \left\| \mathbf{y} - \frac{A\boldsymbol{\eta}_q}{w\sqrt{\vec{\phi}''(\mathbf{x})}} \right\|$ . Since  $\mathbf{P}_{x,w}\mathbf{W}^{-1}(\boldsymbol{\Phi}'')^{-1/2}\mathbf{A} = \mathbf{0}$ , we have

$$\|\mathbf{P}_{x,w}\mathbf{y}\| = \left\| \mathbf{P}_{x,w} \left( \mathbf{y} - \frac{A\boldsymbol{\eta}_q}{w\sqrt{\vec{\phi}''(\mathbf{x})}} \right) \right\| \leq \|\mathbf{P}_{x,w}\| \cdot \left\| \mathbf{y} - \frac{A\boldsymbol{\eta}_q}{w\sqrt{\vec{\phi}''(\mathbf{x})}} \right\| = \|\mathbf{P}_{x,w}\| \cdot \|\mathbf{y}\|_Q.$$

□

**Lemma 6.11.2.** For any projection matrix  $\mathbf{P} \in \mathbb{R}^{m \times m}$ ,  $\boldsymbol{\Sigma} = \text{Diag}(\mathbf{P})$ ,  $i, j \in [m]$ ,  $\mathbf{x} \in \mathbb{R}^m$ , and  $w \in \mathbb{R}_{>0}^m$  we have

1.  $\boldsymbol{\Sigma}_{ii} = \sum_{j \in [m]} \mathbf{P}_{ij}^{(2)}$ ,
2.  $\mathbf{0} \preceq \mathbf{P}^{(2)} \preceq \boldsymbol{\Sigma} \preceq \mathbf{I}$ ,
3.  $\mathbf{P}_{ij}^{(2)} \leq \boldsymbol{\Sigma}_{ii}\boldsymbol{\Sigma}_{jj}$ ,
4.  $|\mathbf{1}_i^T \mathbf{P}^{(2)} \mathbf{x}| \leq \boldsymbol{\Sigma}_{ii} \|\mathbf{x}\|_{\boldsymbol{\Sigma}}$ .
5.  $\nabla_w \log \det(\mathbf{A}^T \mathbf{W} \mathbf{A}) = \boldsymbol{\Sigma}_A(\mathbf{w})\mathbf{w}^{-1}$ .
6.  $\mathbf{J}_w(\boldsymbol{\sigma}_A(\mathbf{w})) = \boldsymbol{\Lambda}_A(\mathbf{w})\mathbf{W}^{-1}$ .

*Proof.* To prove (1), we simply note that by definition of a projection matrix  $\mathbf{P} = \mathbf{P}\mathbf{P}$  and therefore

$$\boldsymbol{\Sigma}_{ii} = \mathbf{P}_{ii} = \mathbf{1}_i^T \mathbf{P} \mathbf{1}_i = \mathbf{1}_i^T \mathbf{P} \mathbf{P} \mathbf{1}_i = \sum_{j \in [n]} \mathbf{P}_{ij}^2 = \sum_{j \in [n]} \mathbf{P}_{ij}^{(2)}$$

To prove (2), we observe that since  $\mathbf{P}$  is a projection matrix, all its eigenvectors are either 0 or 1. Therefore,  $\boldsymbol{\Sigma} \preceq \mathbf{I}$  and by (1)  $\boldsymbol{\Sigma} - \mathbf{P}^{(2)}$  is diagonally dominant. Consequently,  $\boldsymbol{\Sigma} - \mathbf{P}^{(2)} \succeq \mathbf{0}$ . Rearranging terms and using the well known fact that the shur product of two positive semi-definite matrices is positive semi-definite yields (2).

To prove (3), we use  $\mathbf{P} = \mathbf{P}\mathbf{P}$ , Cauchy-Schwarz, and (1) to derive

$$\mathbf{P}_{ij} = \sum_{k \in [n]} \mathbf{P}_{ik} \mathbf{P}_{kj} \leq \sqrt{\left( \sum_{k \in [n]} \mathbf{P}_{ik}^2 \right) \left( \sum_{k \in [n]} \mathbf{P}_{kj}^2 \right)} = \sqrt{\boldsymbol{\Sigma}_{ii} \boldsymbol{\Sigma}_{jj}}.$$

Squaring then yields (3).

To prove (4), we note that by the definition of  $\mathbf{P}^{(2)}$  and Cauchy-Schwarz, we have

$$\left| \mathbf{1}_i^T \mathbf{P}^{(2)} \mathbf{x} \right| = \left| \sum_{j \in [n]} \mathbf{P}_{ij}^{(2)} \bar{x}_j \right| \leq \sqrt{\left( \sum_{j \in [n]} \Sigma_{jj} \bar{x}_j^2 \right)} \cdot \sum_{j \in [n]} \frac{\mathbf{P}_{ij}^{(4)}}{\Sigma_{jj}} \quad (6.42)$$

Now, by (1) and (3), we know that

$$\sum_{j \in [n]} \frac{\mathbf{P}_{ij}^4}{\Sigma_{jj}} \leq \sum_{j \in [n]} \frac{\mathbf{P}_{ij}^2 \Sigma_{ii} \Sigma_{jj}}{\Sigma_{jj}} = \Sigma_{ii} \sum_{j \in [n]} \mathbf{P}_{ij}^2 = \Sigma_{ii}^2 \quad (6.43)$$

Since  $\|\mathbf{x}\|_{\Sigma} \stackrel{\text{def}}{=} \sqrt{\sum_{j \in [n]} \Sigma_{jj} \mathbf{x}_j^2}$ , combining (6.42) and (6.43) yields  $\left| \mathbf{1}_i^T \mathbf{P}^{(2)} \mathbf{x} \right| \leq \Sigma_{ii} \|\mathbf{x}\|_{\Sigma}$  as desired.

To prove (5), we let  $f(\mathbf{w}) = \log \det(\mathbf{A}^T \mathbf{W} \mathbf{A})$ . Lemma 2.3.8 shows that

$$\begin{aligned} \frac{\partial}{\partial w_i} f(\mathbf{w}) &= \text{Tr} \left( (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{a}_i \mathbf{a}_i^T \right) \\ &= \Sigma_{\mathbf{A}}(\mathbf{w})_{ii} w_i^{-1}. \end{aligned}$$

To prove (6), we let  $\mathbf{R}_{\mathbf{A}}(\mathbf{w})_{ij} \stackrel{\text{def}}{=} \mathbf{a}_i^T (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{a}_j$ . Lemma 2.3.8 shows that

$$\frac{\partial}{\partial w_i} (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} = -(\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{a}_i \mathbf{a}_i^T (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1}$$

and hence  $\frac{\partial}{\partial w_k} [\mathbf{R}_{\mathbf{A}}(\mathbf{w})]_{ij} = -\mathbf{R}_{\mathbf{A}}(\mathbf{w})_{ik} \mathbf{R}_{\mathbf{A}}(\mathbf{w})_{kj}$ . Since by definition  $\sigma_{\mathbf{A}}(\mathbf{w})_i = w_i \mathbf{R}_{\mathbf{A}}(\mathbf{w})_{ii}$  by the previous lemma, we have that

$$\frac{\partial}{\partial w_j} \sigma_{\mathbf{A}}(\mathbf{w})_i = \mathbf{1}_{i=j} \mathbf{R}_{\mathbf{A}}(\mathbf{w})_{ii} - w_i \mathbf{R}_{\mathbf{A}}(\mathbf{w})_{ij}^{(2)}.$$

Writing this in matrix form and recalling the definition of the Jacobian then yields

$$\mathbf{J}_{\mathbf{w}}(\sigma_{\mathbf{A}}(\mathbf{w})) = \text{Diag}(\mathbf{R}_{\mathbf{A}}(\mathbf{w})) - \mathbf{W} \mathbf{R}_{\mathbf{A}}(\mathbf{w})^{(2)}.$$

Right multiplying by  $\mathbf{I} = \mathbf{W} \mathbf{W}^{-1}$  and recalling the definition of  $\mathbf{\Lambda}_{\mathbf{A}}$  then yields the result.  $\square$

**Lemma 6.11.3.** *For any  $x, \varepsilon$  and  $\lambda > 0$ , we have*

$$\frac{\pi}{2} |x| - \pi \varepsilon - \frac{1}{\lambda} \leq x \tan^{-1}(\lambda(x + \varepsilon)) \leq \frac{\pi}{2} |x|.$$

*Proof.* We first consider the case  $\varepsilon = 0$ . Note that

$$x \tan^{-1}(\lambda x) \leq \frac{\pi}{2} |x|.$$

Also, we note that

$$x \tan^{-1}(\lambda x) \geq |x| \left( \frac{\pi}{2} - \frac{1}{\lambda |x|} \right)$$

because

$$\left| \tan\left(\frac{\pi}{2} - \frac{1}{\lambda |x|}\right) \right| = \left| \frac{\cos\left(\frac{1}{\lambda |x|}\right)}{\sin\left(\frac{1}{\lambda |x|}\right)} \right| \leq \lambda |x|.$$

Hence, we have

$$\frac{\pi}{2} |x| - \frac{1}{\lambda} \leq x \tan^{-1}(\lambda x) \leq \frac{\pi}{2} |x|.$$

For  $\varepsilon \neq 0$ , we have

$$\frac{\pi}{2} |x + \varepsilon| - \frac{1}{\lambda} \leq (x + \varepsilon) \tan^{-1}(\lambda(x + \varepsilon)) \leq \frac{\pi}{2} |x + \varepsilon|.$$

Thus, we have

$$\frac{\pi}{2} |x| - \pi\varepsilon - \frac{1}{\lambda} \leq x \tan^{-1}(\lambda(x + \varepsilon)) \leq \frac{\pi}{2} |x|.$$

□

## ■ 6.12 Appendix: Projection on Mixed Norm Ball

In this section, we give an algorithm to solve the problem

$$\max_{\|\mathbf{x}\|_2 + \|\mathbf{l}^{-1}\mathbf{x}\|_\infty \leq 1} \langle \mathbf{a}, \mathbf{x} \rangle \quad (6.44)$$

for some given vector  $\mathbf{l}$  and  $\mathbf{a}$ . This is used in Section 6.7.1 to compute the weights. We note that

$$\begin{aligned} \max_{\|\mathbf{x}\|_2 + \|\mathbf{l}^{-1}\mathbf{x}\|_\infty \leq 1} \langle \mathbf{a}, \mathbf{x} \rangle &= \max_{0 \leq t \leq 1} \max_{\|\mathbf{x}\|_2 \leq 1-t \text{ and } -tl_i \leq x_i \leq tl_i} \langle \mathbf{a}, \mathbf{x} \rangle \\ &= \max_{0 \leq t \leq 1} (1-t) \max_{\|\mathbf{x}\|_2 \leq 1 \text{ and } -\frac{t}{1-t}l_i \leq x_i \leq \frac{t}{1-t}l_i} \langle \mathbf{a}, \mathbf{x} \rangle \\ &= \max_{0 \leq t \leq 1} (1-t) f\left(\frac{t}{1-t}\right) \end{aligned} \quad (6.45)$$

where

$$f\left(\frac{t}{1-t}\right) = \max_{\|\mathbf{x}\|_2 \leq 1, -tl_i \leq x_i \leq tl_i} \langle \mathbf{a}, \mathbf{x} \rangle.$$

In [165], we showed that the maximizer  $\mathbf{x}$  in the problem  $f$  must be of the form

$$\mathbf{x}^{i_t} = \begin{cases} \frac{\frac{t}{1-t} \text{sign}(a_j) l_j}{\sqrt{\frac{1 - \left(\frac{t}{1-t}\right)^2 \sum_{k=0}^{i_t} l_k^2}{1 - \sum_{k=0}^{i_t} a_k^2}}} & \text{if } j \in \{1, 2, \dots, i_t\} \\ \mathbf{a}_j & \text{otherwise} \end{cases}. \quad (6.46)$$

where  $i_t$  be the first coordinate such that

$$\frac{1 - \left(\frac{t}{1-t}\right)^2 \sum_{k=0}^{i_t} l_k^2}{1 - \sum_{k=0}^{i_t} a_k^2} \leq \frac{\left(\frac{t}{1-t}\right)^2 l_i^2}{a_i^2}.$$

Note that  $i_t \geq i_s$  if  $t \leq s$ . Therefore, we have that the set of  $t$  such that  $i_t = j$  is simply an interval given by<sup>8</sup>

$$\frac{|a_j|}{\sqrt{l_j^2 \left(1 - \sum_{k=0}^j a_k^2\right) + a_j^2 \sum_{k=0}^j l_k^2}} \leq \frac{t}{1-t} < \frac{|a_{j-1}|}{\sqrt{l_{j-1}^2 \left(1 - \sum_{k=0}^{j-1} a_k^2\right) + a_{j-1}^2 \sum_{k=0}^{j-1} l_k^2}}. \quad (6.47)$$

<sup>8</sup>There are some boundary cases we ignored for simplicity.

Therefore, we know that

$$\begin{aligned} f\left(\frac{t}{1-t}\right) &= \langle \mathbf{a}, \mathbf{x}^{(i_t)} \rangle \\ &= \frac{t}{1-t} \sum_{j=1}^{i_t} |a_j| |l_j| + \sqrt{1 - \left(\frac{t}{1-t}\right)^2 \sum_{k=0}^{i_t} l_k^2} \sqrt{1 - \sum_{k=0}^{i_t} a_k^2}. \end{aligned}$$

Putting this into 6.45, we have that

$$\max_{\|\mathbf{x}\|_2 + \|\mathbf{l}^{-1}\mathbf{x}\|_\infty \leq 1} \langle \mathbf{a}, \mathbf{x} \rangle = \max_{0 \leq t \leq 1} t \sum_{j=1}^{i_t} |a_j| |l_j| + \sqrt{(1-t)^2 - t^2 \sum_{k=0}^{i_t} l_k^2} \sqrt{1 - \sum_{k=0}^{i_t} a_k^2}.$$

Note that the function  $t \sum_{j=1}^i |a_j| |l_j| + \sqrt{(1-t)^2 - t^2 \sum_{k=0}^i l_k^2} \sqrt{1 - \sum_{k=0}^i a_k^2}$  is concave and the solution has a close form. Therefore, one can compute the maximum value for each interval of  $t$  (6.47) and find which is the best. Hence, we get the following algorithm.

$\mathbf{x} = \text{projectOntoMixedNormBallParallel}(\mathbf{a}, \mathbf{l})$
1. Set $\mathbf{a} = \mathbf{a} / \ \mathbf{a}\ _2$ .
2. Sort the coordinate such that $ a_i  / l_i$ is in descending order.
3. Precompute $\sum_{k=0}^i l_k^2$ , $\sum_{k=0}^i a_k^2$ and $\sum_{j=1}^i  a_j   l_j $ for all $i$ .
4. Let $g_i(t) = t \sum_{j=1}^i  a_j   l_j  + \sqrt{(1-t)^2 - t^2 \sum_{k=0}^i l_k^2} \sqrt{1 - \sum_{k=0}^i a_k^2}$ .
5. For each $j \in \{1, \dots, n\}$ , Find $t_j = \arg \max_{i=j} g_j(t)$ using (6.47)
6. Find $i = \arg \max_i g_i(t_i)$ .
7. Output $(1 - t_i)\mathbf{x}^{(i)}$ defined by 6.46.

The discussion above leads to the following theorem.

**Theorem 6.12.1.** *The algorithm `projectOntoMixedNormBallParallel` outputs a solution to*

$$\max_{\|\mathbf{x}\|_2 + \|\mathbf{l}^{-1}\mathbf{x}\|_\infty \leq 1} \langle \mathbf{a}, \mathbf{x} \rangle$$

in total work  $\tilde{O}(m)$  and depth  $\tilde{O}(1)$ .

## ■ 6.13 Appendix: A $\tilde{O}(n)$ -Self-Concordant Barrier Function

In this section, we proved the barrier function introduced in Subsection 6.3.5 is an  $\tilde{O}(n)$ -self-concordant barrier. We separate the proof of Theorem 6.3.5 into two lemma, Lemma 6.13.3 and Lemma 6.13.8. The proof is pretty straightforward, but it requires some familiarity with matrix calculus.

### ■ 6.13.1 Notation

Here, we define some notations that is used throughout in this section. They different from previous sections because this section our domain is  $\Omega \stackrel{\text{def}}{=} \{\mathbf{A}\mathbf{x} > \mathbf{b}\}$  instead of  $\mathbf{A}^\top \mathbf{x} = \mathbf{b}$ .

1. Let  $f(x, w) = \ln \det(\mathbf{A}^\top \mathbf{S}^{-1} \mathbf{W} \mathbf{S}^{-1} \mathbf{A}) - \frac{n}{m} \sum_{i=1}^m w_i \ln w_i - \frac{n}{m} \sum_{i=1}^m \ln s_i$ .
2. Let  $\mathbf{S} = \text{Diag} \mathbf{A}\mathbf{x} - \mathbf{b}$  and  $\mathbf{A}_x = \mathbf{S}^{-1} \mathbf{A}$ .

3. Let  $\mathbf{P}(x, w) = \sqrt{\mathbf{W}} \mathbf{A}_x (\mathbf{A}_x^T \mathbf{W} \mathbf{A}_x)^{-1} \mathbf{A}_x^T \sqrt{\mathbf{W}}$ ,  $\sigma(x, w) = \mathbf{diag}(\mathbf{P}(x, w))$ .
4. Let  $\mathbf{P}^{(2)}(x, w)$  be a  $m \times m$  matrix such that  $(\mathbf{P}^{(2)}(x, w))_{ij} = (\mathbf{P}(x, w))_{ij}^2$ .
5. Let  $\mathbf{\Lambda}(x, w) = \mathbf{\Sigma}(x, w) - \mathbf{P}^{(2)}(x, w)$ .
6. Let  $w_x$  be the unique maximizer of  $f(x, w)$ , i.e.  $w_x \stackrel{\text{def}}{=} \arg \max_{w_i \geq 0} f(x, w)$ .
7. Let  $\mathbf{P}_x = \mathbf{P}(x, w)$ ,  $\sigma_x = \sigma(x, w_x)$ ,  $\mathbf{\Sigma}_x = \mathbf{\Sigma}(x, w_x)$ ,  $\mathbf{P}_x^{(2)} = \mathbf{P}^{(2)}(x, w_x)$  and  $\mathbf{\Lambda}_x = \mathbf{\Lambda}(x, w_x)$ .

Similar to lemma 6.11.2, we have  $\mathbf{0} \preceq \mathbf{P}^{(2)}(x, w) \preceq \mathbf{\Sigma}(x, w) \preceq \mathbf{I}$ ,  $\mathbf{0} \preceq \mathbf{\Lambda}(x, w)$  and  $\sum_i \sigma_i(x, w) = n$ .

### ■ 6.13.2 Basic properties of $f$ and $p$

First, we compute the gradient and Hessian of  $f$  and  $p$ .

**Lemma 6.13.1.** *For any interior point  $x \in \Omega$  and any positive weight  $w > 0$ , we have*

$$\begin{aligned}
\nabla_x f(x, w) &= -2\mathbf{A}_x^T \sigma(x, w) - \frac{n}{m} \mathbf{A}_x^T \mathbf{1}, \\
\nabla_w f(x, w) &= \mathbf{W}^{-1} \sigma(x, w) - \frac{n}{m} \ln(ew), \\
\nabla_{xx}^2 f(x, w) &= \mathbf{A}_x^T \left( 6\mathbf{\Sigma}(x, w) - 4\mathbf{P}^{(2)}(x, w) + \frac{n}{m} \mathbf{I} \right) \mathbf{A}_x, \\
\nabla_{ww}^2 f(x, w) &= -\mathbf{W}^{-1} \mathbf{P}^{(2)}(x, w) \mathbf{W}^{-1} - \frac{n}{m} \mathbf{W}^{-1}, \\
\nabla_{xw}^2 f(x, w) &= -2\mathbf{A}_x^T \left( \mathbf{\Sigma}(x, w) - \mathbf{P}^{(2)}(x, w) \right) \mathbf{W}^{-1}.
\end{aligned} \tag{6.48}$$

For  $x \in \Omega$ ,  $w_x$  uniquely exist in the region  $\{\frac{1}{e} \leq w_i \leq \frac{m}{n}\}$  and  $p$  satisfies

$$\begin{aligned}
\nabla p(x) &= -2\mathbf{A}_x^T \sigma_x - \frac{n}{m} \mathbf{A}_x^T \mathbf{1}. \\
\nabla^2 p(x) &= \mathbf{A}_x^T \left( 6\mathbf{\Sigma}_x - 4\mathbf{P}_x^{(2)} + \frac{n}{m} \mathbf{I} \right) \mathbf{A}_x + 4\mathbf{A}_x^T \mathbf{\Lambda}_x \left( \mathbf{P}_x^{(2)} + \frac{n}{m} \mathbf{W}_x \right)^{-1} \mathbf{\Lambda}_x \mathbf{A}_x.
\end{aligned} \tag{6.49}$$

Also, we have

$$\mathbf{J}_x w = -2\mathbf{W} \left( \mathbf{P}^{(2)}(x, w) + \frac{n}{m} \mathbf{W} \right)^{-1} \left( \mathbf{\Sigma}(x, w) - \mathbf{P}^{(2)}(x, w) \right) \mathbf{A}_x. \tag{6.50}$$

*Proof.* For  $\nabla_x f$ , we have

$$\begin{aligned}
\frac{\partial f}{\partial s_i} &= \text{Tr} \left( (\mathbf{A}^T \mathbf{S}^{-1} \mathbf{W} \mathbf{S}^{-1} \mathbf{A})^{-1} \left( \frac{\partial}{\partial s_i} \right) (\mathbf{A}^T \mathbf{S}^{-1} \mathbf{W} \mathbf{S}^{-1} \mathbf{A}) \right) - \frac{n}{m} \frac{1}{s_i} \\
&= -2\text{Tr} \left( \mathbf{A} (\mathbf{A}^T \mathbf{S}^{-1} \mathbf{W} \mathbf{S}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \frac{w_i}{s_i^3} \mathbf{1}_i \mathbf{1}_i^T \right) - \frac{n}{m} \frac{1}{s_i} \\
&= -2 \frac{\sigma_i(x, w)}{s_i} - \frac{n}{m} \frac{1}{s_i}
\end{aligned}$$

By Chain rule, we have

$$\nabla_x f(x, w) = -2\mathbf{A}_x^T \sigma(x, w) - \frac{n}{m} \mathbf{A}_x^T \mathbf{1}.$$

For  $\nabla_w f$ , we have

$$\begin{aligned}\frac{\partial f}{\partial w_i} &= \text{Tr} \left( (\mathbf{A}^T \mathbf{S}^{-1} \mathbf{W} \mathbf{S}^{-1} \mathbf{A})^{-1} \left( \frac{\partial}{\partial w_i} \right) (\mathbf{A}^T \mathbf{S}^{-1} \mathbf{W} \mathbf{S}^{-1} \mathbf{A}) \right) - \frac{n}{m} \ln w_i - \frac{n}{m} \\ &= \frac{\sigma_i(x, w)}{w_i} - \frac{n}{m} \ln(ew_i).\end{aligned}$$

Hence, we have

$$\nabla_w f(x, w) = \mathbf{W}^{-1} \sigma - \frac{n}{m} \ln(ew).$$

For  $\nabla_{xx} f$ , taking partial derivative on  $\frac{\partial f}{\partial s_i}$ , we have

$$\begin{aligned}\frac{\partial^2 f}{\partial s_i \partial s_j} &= \frac{\partial}{\partial s_j} \left( -2 \text{Tr} \left( \mathbf{A} (\mathbf{A}^T \mathbf{S}^{-1} \mathbf{W} \mathbf{S}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \frac{w_i}{s_i^3} \mathbf{1}_i \mathbf{1}_i^T \right) - \frac{n}{m} \frac{1}{s_i} \right) \\ &= -4 \text{Tr} \left( \mathbf{A} (\mathbf{A}_x^T \mathbf{W} \mathbf{A}_x)^{-1} \mathbf{A}^T \frac{w_j}{s_j^3} \mathbf{1}_j \mathbf{1}_j^T \mathbf{A} (\mathbf{A}_x^T \mathbf{W}_x \mathbf{A})^{-1} \mathbf{A}^T \frac{w_i}{s_i^3} \mathbf{1}_i \mathbf{1}_i^T \right) \\ &\quad 6 \text{Tr} \left( \mathbf{A} (\mathbf{A}_x^T \mathbf{W} \mathbf{A}_x)^{-1} \mathbf{A}^T \frac{w_i}{s_i^4} \mathbf{1}_i \mathbf{1}_i^T \mathbf{1}_{i=j} \right) + \frac{n}{m} \frac{1}{s_i^2} \mathbf{1}_{i=j} \\ &= 6 \frac{\sigma_i(x, w)}{s_i^2} \mathbf{1}_{i=j} - 4 \frac{(\mathbf{P}(x, w))_{ij}^2}{s_i s_j} + \frac{n}{m} \frac{1}{s_i^2} \mathbf{1}_{i=j}.\end{aligned}$$

By Chain rule, we have

$$\nabla_{xx} f(x, w) = \mathbf{A}_x^T \left( 6 \Sigma(x, w) - 4 \mathbf{P}^{(2)}(x, w) + \frac{n}{m} \mathbf{I} \right) \mathbf{A}_x.$$

For  $\nabla_{ww} f$ , taking partial derivative on  $\frac{\partial f}{\partial w_i}$ , we have

$$\begin{aligned}\frac{\partial^2 f}{\partial w_i \partial w_j} &= \frac{\partial}{\partial w_j} \left( \text{Tr} \left( \mathbf{A}_x (\mathbf{A}_x^T \mathbf{W} \mathbf{A}_x)^{-1} \mathbf{A}_x^T e_i e_i^T \right) - \frac{n}{m} \ln w_i - \frac{n}{m} \right) \\ &= -\text{Tr} \left( \mathbf{A}_x (\mathbf{A}_x^T \mathbf{W} \mathbf{A}_x)^{-1} \mathbf{A}_x^T e_j e_j^T \mathbf{A}_x (\mathbf{A}_x^T \mathbf{W} \mathbf{A}_x)^{-1} \mathbf{A}_x^T e_i e_i^T \right) - \frac{n}{m} \frac{1}{w_j} \mathbf{1}_{i=j} \\ &= -\frac{(\mathbf{P}(x, w))_{ij}^2}{w_i w_j} - \frac{n}{m} \frac{1}{w_j} \mathbf{1}_{i=j}.\end{aligned}$$

Hence, we have

$$\nabla_{ww} f(x, w) = -\mathbf{W}^{-1} \mathbf{P}^{(2)}(x, w) \mathbf{W}^{-1} - \frac{n}{m} \mathbf{W}^{-1}.$$

For  $\nabla_{xw} f$ , taking partial derivative on  $\frac{\partial f}{\partial w_i}$ , we have

$$\begin{aligned}\frac{\partial^2 f}{\partial s_j \partial w_i} &= \frac{\partial}{\partial s_j} \left( \text{Tr} \left( \mathbf{A} (\mathbf{A}_x^T \mathbf{W} \mathbf{A}_x)^{-1} \mathbf{A}^T \frac{e_i e_i^T}{s_i^2} \right) - \frac{n}{m} \ln w_i - \frac{n}{m} \right) \\ &= 2 \text{Tr} \left( \mathbf{A} (\mathbf{A}_x^T \mathbf{W} \mathbf{A}_x)^{-1} \mathbf{A}^T \frac{w_j e_j e_j^T}{s_j^3} \mathbf{A} (\mathbf{A}_x^T \mathbf{W} \mathbf{A}_x)^{-1} \mathbf{A}^T \frac{e_i e_i^T}{s_i^2} \right) \\ &\quad - 2 \text{Tr} \left( \mathbf{A} (\mathbf{A}_x^T \mathbf{W} \mathbf{A}_x)^{-1} \mathbf{A}^T \frac{e_i e_i^T}{s_i^3} \mathbf{1}_{i=j} \right) \\ &= 2 \frac{(\mathbf{P}(x, w))_{ij}}{s_j w_i} - 2 \frac{\sigma_i(x, w)}{s_j w_i} \mathbf{1}_{i=j}.\end{aligned}$$



By Chain rule, we have

$$\nabla_{xw}f(x, w) = -2\mathbf{A}_x^T \left( \boldsymbol{\Sigma}(x, w) - \mathbf{P}^{(2)}(x, w) \right) \mathbf{W}^{-1}.$$

Now, we prove that  $w_x$  uniquely exists in the region  $\{\frac{1}{e} \leq w_i \leq \frac{m}{n}\}$ . For any  $w$  such that  $w_i > \frac{m}{n}$  for some  $i$ , we have

$$\begin{aligned} \frac{\partial f}{\partial w_i} &= \frac{\sigma_i(x, w)}{w_i} - \frac{n}{m} \ln(ew_i) \\ &\leq \frac{1}{w_i} - \frac{n}{m} \ln(ew_i) \\ &\leq \frac{n}{m} - \frac{n}{m} \ln\left(e\frac{m}{n}\right) < 0. \end{aligned}$$

Hence, we have  $\sup_{w_i \geq 0} f(x, w) = \sup_{\frac{m}{n} \geq w_i \geq 0} f(x, w)$ . Since the region  $0 \leq w_i \leq \frac{m}{n}$  is a compact set and  $f$  is a continuous function, there is  $w_x$  such that attains  $\sup_{w_i \geq 0} f(x, w)$ . For the uniqueness, we note that  $\nabla_{ww}f(x, w) \prec \mathbf{0}$  for all  $w$  and hence  $f$  is strictly concave and the maximizer is unique. Now, we note that for any  $w$  such that  $w_i < \frac{1}{e}$  for some  $i$ , we have

$$\frac{\partial f}{\partial w_i} = \frac{\sigma_i(x, w)}{w_i} - \frac{n}{m} \ln(ew_i) > 0.$$

Therefore, the minimizer is in the region  $\{\frac{1}{e} \leq w_i \leq \frac{m}{n}\}$ .

To compute  $p$ , we note that  $f(x, w)$  is concave in  $w$  and the maximizer is in the interior of the set  $\{w_i \geq 0\}$  and therefore the optimality conditions imply that

$$\nabla_w f(x, w_x) = 0. \quad (6.51)$$

Since  $\nabla_{ww}^2 f(x, w_x)$  is invertible, implicit function theorem (Lemma 2.3.9) shows that  $w_x$  is differentiable and satisfies

$$\mathbf{J}_x w = -(\nabla_{ww}^2 f(x, w_x))^{-1} (\nabla_{wx}^2 f(x, w_x)).$$

This gives (6.50).

Using that  $p(x) = f(x, w_x)$  and applying the chain rule, we have

$$\begin{aligned} \nabla p(x) &= \nabla_x f(x, w_x) + \nabla_w f(x, w_x) \mathbf{J}_x w_x \\ &= \nabla_x f(x, w_x) \end{aligned}$$

because the optimality condition. Taking derivative again, we have

$$\begin{aligned} \nabla^2 p(x) &= \nabla_{xx}^2 f(x, w_x) + \nabla_{xw}^2 f(x, w_x) \mathbf{J}_x w_x \\ &= \nabla_{xx}^2 f(x, w_x) - \nabla_{xw}^2 f(x, w_x) (\nabla_{ww}^2 f(x, w_x))^{-1} (\nabla_{wx}^2 f(x, w_x)) \end{aligned}$$

Substituting in the computed values for  $\nabla_{xx}^2 f(x, w_x)$ ,  $\nabla_{xw}^2 f(x, w_x)$ , and  $\nabla_{ww}^2 f(x, w_x)$  then yields the result.  $\square$

Now, we discuss properties of the optimal weight  $w_x$ .

**Lemma 6.13.2.** *For  $x \in \Omega$ , we have  $\sigma_x = \frac{n}{m} w_x \ln(ew_x)$ . Also, we have  $\max_i \frac{\sigma_{x,i}}{w_{x,i}} \leq \frac{n}{m} \ln\left(\frac{me}{n}\right)$ .*

*Proof.* The optimality conditions (6.51) implies

$$\sigma_x = \frac{n}{m} w_x \ln(ew_x).$$

Lemma 6.13.1 shows that  $\frac{1}{e} \leq w_x \leq \frac{m}{n}$  and hence

$$\frac{\sigma_{x,i}}{w_{x,i}} \leq \frac{n}{m} \ln\left(\frac{me}{n}\right).$$

□

Using these, we can bound  $(\nabla p)^T (\nabla^2 p)^{-1} (\nabla p)$ .

**Lemma 6.13.3.** *For  $x \in \Omega$ , we have*

$$(\nabla p(x))^T (\nabla^2 p(x))^{-1} (\nabla p(x)) \leq 6n.$$

*Proof.* Using  $\nabla p(x) = -2\mathbf{A}_x^T \sigma_x - \frac{n}{m} \mathbf{A}_x^T \mathbf{1}$ , we have

$$\begin{aligned} & (\nabla p(x))^T (\nabla^2 p(x))^{-1} (\nabla p(x)) \\ & \leq 8\sigma_x^T \mathbf{A}_x (\nabla^2 p(x))^{-1} \mathbf{A}_x^T \sigma_x + 2\left(\frac{n}{m}\right)^2 \mathbf{1}^T \mathbf{A}_x (\nabla^2 p(x))^{-1} \mathbf{A}_x^T \mathbf{1}. \end{aligned}$$

For the first term, we use  $\mathbf{\Lambda}_x \succeq 0$  and  $\mathbf{P}_x^{(2)} \succeq 0$  and get

$$\begin{aligned} & \sigma_x^T \mathbf{A}_x (\nabla^2 p(x))^{-1} \mathbf{A}_x^T \sigma_x \\ & \leq \sigma_x^T \mathbf{A}_x \left( \mathbf{A}_x^T \left( 6\mathbf{\Sigma}_x - 4\mathbf{P}_x^{(2)} + \frac{n}{m} \mathbf{I}_m \right) \mathbf{A}_x \right)^{-1} \mathbf{A}_x^T \sigma_x \\ & \leq \frac{1}{2} \sigma_x^T \mathbf{A}_x \left( \mathbf{A}_x^T \mathbf{\Sigma}_x \mathbf{A}_x \right)^{-1} \mathbf{A}_x^T \sigma_x \\ & \leq \frac{1}{2} \sum_{i=1}^m \sigma_{x,i} \\ & = \frac{n}{2}. \end{aligned}$$

For the second term, we have

$$\begin{aligned} & \mathbf{1}_m^T \mathbf{A}_x (\nabla^2 p(x))^{-1} \mathbf{A}_x^T \mathbf{1}_m \\ & \leq \mathbf{1}_m^T \mathbf{A}_x \left( \mathbf{A}_x^T \left( 6\mathbf{\Sigma}_x - 4\mathbf{P}_x^{(2)} + \frac{n}{m} \mathbf{I}_m \right) \mathbf{A}_x \right)^{-1} \mathbf{A}_x^T \mathbf{1}_m \\ & \leq \frac{m}{n} \mathbf{1}_m^T \mathbf{A}_x \left( \mathbf{A}_x^T \mathbf{A}_x \right)^{-1} \mathbf{A}_x^T \mathbf{1}_m \\ & \leq \frac{m^2}{n}. \end{aligned}$$

Combining both terms, we have the result. □

### ■ 6.13.3 $w_x$ , $\mathbf{\Sigma}_x$ and $\mathbf{P}_x^{(2)}$ are stable

In this subsection, we write  $x_t = x + th$ ,  $\mathbf{A}_t = \mathbf{A}_{x_t}$ ,  $s_t = \mathbf{A}_t x_t - b$ ,  $\mathbf{S}_t = \mathbf{Diag}(s_t)$ ,  $w_t = w_{x_t}$ ,  $\mathbf{W}_t = \mathbf{Diag}(w_t)$ ,  $\mathbf{P}_t = \mathbf{P}(x_t, w_t)$ ,  $\sigma_t = \sigma(x_t, w_t)$ ,  $\mathbf{\Sigma}_t = \mathbf{\Sigma}(x_t, w_t)$ ,  $\mathbf{P}_t^{(2)} = \mathbf{P}^{(2)}(x_t, w_t)$  and  $\mathbf{\Lambda}_t = \mathbf{\Lambda}(x_t, w_t)$ . First of all, we show that  $\mathbf{A}_x h$  is small.

**Lemma 6.13.4.** *For any  $x \in \Omega$  and  $h \in \mathbb{R}^n$ , we have*

$$\|\mathbf{A}_x h\|_\infty \leq \frac{1}{\sqrt{2}} \ln\left(\frac{me}{n}\right) \|h\|_{\nabla^2 p(x)}.$$

Therefore, we have

$$\|\mathbf{S}_t^{-1} \frac{d}{dt} s_t\|_\infty \leq \frac{1}{\sqrt{2}} \ln\left(\frac{me}{n}\right) \|h\|_{\nabla^2 p(x_t)}.$$

*Proof.* Note that

$$\begin{aligned} \|\mathbf{A}_x h\|_\infty &= \max_i \langle e_i, \mathbf{A}_x h \rangle \\ &\leq \max_i \left\langle (\mathbf{A}_x^T \mathbf{W} \mathbf{A}_x)^{-1/2} \mathbf{A}_x^T \mathbf{1}_i, (\mathbf{A}_x^T \mathbf{W} \mathbf{A}_x)^{1/2} h \right\rangle \\ &\leq \left( \max_i \frac{\sigma_{x,i}}{w_{x,i}} \right) \|h\|_{\mathbf{A}_x^T \mathbf{W} \mathbf{A}_x}. \end{aligned}$$

Lemma 6.13.2 shows that  $\max_i \frac{\sigma_{x,i}}{w_{x,i}} \leq \frac{n}{m} \ln\left(\frac{me}{n}\right)$  and hence

$$\|\mathbf{A}_x h\|_\infty \leq \frac{n}{m} \ln\left(\frac{me}{n}\right) \|h\|_{\mathbf{A}_x^T \mathbf{W} \mathbf{A}_x}.$$

On the other hand, we have

$$\begin{aligned} \nabla^2 p &= \mathbf{A}_x^T \left( 6\boldsymbol{\Sigma}_x - 4\mathbf{P}_x^{(2)} + \frac{n}{m} \mathbf{I} \right) \mathbf{A}_x + 4\mathbf{A}_x^T \boldsymbol{\Lambda}_x \left( \frac{n}{m} \mathbf{W}_x + \mathbf{P}_x^{(2)} \right)^{-1} \boldsymbol{\Lambda}_x \mathbf{A}_x \\ &\succeq \mathbf{A}_x^T \left( 2\boldsymbol{\Sigma}_x + \frac{n}{m} \mathbf{I} \right) \mathbf{A}_x \end{aligned}$$

Lemma 6.13.2 shows that  $\sigma_x = \frac{n}{m} w \ln(ew)$  and hence

$$2\sigma + \frac{n}{m} = \frac{n}{m} (2w \ln(ew) + 1).$$

Using  $\frac{2w \ln(ew) + 1}{w} \geq 4 - \ln(4) \geq 2$  for all  $w > 0$ , we have

$$\nabla^2 p(x) \succeq 2\mathbf{A}_x^T \mathbf{W}_x \mathbf{A}_x$$

Hence, we have

$$\|\mathbf{A}_x h\|_\infty \leq \frac{1}{\sqrt{2}} \ln\left(\frac{me}{n}\right) \|h\|_{\nabla^2 p(x)}.$$

The last claim follows from the equality

$$\mathbf{S}_t^{-1} \frac{d}{dt} s_t = \mathbf{A}_t h.$$

□

In the next lemma, we show that  $w_t$  is stable both  $\mathbf{P}^{(2)}$  norm,  $\boldsymbol{\Sigma}$  norm and  $\ell^\infty$  norm.

**Lemma 6.13.5.** *We have*

$$\begin{aligned} \|\mathbf{W}_t^{-1} \frac{d}{dt} w_t\|_{\mathbf{P}_t^{(2)}} &\leq \frac{1}{2} \|h\|_{\nabla^2 p(x_t)}, \\ \|\mathbf{W}_t^{-1} \frac{d}{dt} w_t\|_{\boldsymbol{\Sigma}_t} &\leq \sqrt{\ln\left(\frac{m}{n}e\right)} \|h\|_{\nabla^2 p(x_t)}, \\ \|\mathbf{W}_t^{-1} \frac{d}{dt} w_t\|_\infty &\leq \ln\left(\frac{me}{n}\right) \left( 1 + \sqrt{2} \ln\left(\frac{me^2}{n}\right) \right) \|h\|_{\nabla^2 p(x_t)}. \end{aligned}$$

*Proof.* Using (6.50), we have

$$\frac{d}{dt}w_t = -2\mathbf{W}_t \left( \frac{n}{m}\mathbf{W}_t + \mathbf{P}_t^{(2)} \right)^{-1} \mathbf{\Lambda}_t \mathbf{A}_t h. \quad (6.52)$$

Let  $z = \left( \frac{n}{m}\mathbf{W}_t + \mathbf{P}_t^{(2)} \right)^{-1} \mathbf{\Lambda}_t \mathbf{A}_t h$ . First, we note that

$$\|z\|_{\mathbf{P}_t^{(2)}}^2 = h^T \mathbf{A}_t^T \mathbf{\Lambda}_t \left( \frac{n}{m}\mathbf{W}_t + \mathbf{P}_t^{(2)} \right)^{-1} \mathbf{P}_t^{(2)} \left( \frac{n}{m}\mathbf{W}_t + \mathbf{P}_t^{(2)} \right)^{-1} \mathbf{\Lambda}_t \mathbf{A}_t h.$$

Since  $\mathbf{P}_t^{(2)} \preceq \frac{n}{m}\mathbf{W}_t + \mathbf{P}_t^{(2)}$  and  $\nabla^2 p \succeq 4\mathbf{A}_t^T \mathbf{\Lambda}_t \left( \frac{n}{m}\mathbf{W}_t + \mathbf{P}_t^{(2)} \right)^{-1} \mathbf{\Lambda}_t \mathbf{A}_t$  (6.49), we have

$$\|z\|_{\mathbf{P}_t^{(2)}}^2 \leq h^T \mathbf{A}_t^T \mathbf{\Lambda}_t \left( \frac{n}{m}\mathbf{W}_t + \mathbf{P}_t^{(2)} \right)^{-1} \mathbf{\Lambda}_t \mathbf{A}_t h \leq \frac{1}{4} \|h\|_{\nabla^2 p(x_t)}^2.$$

For the  $\mathbf{\Sigma}_t$  bound, note that

$$\|z\|_{\mathbf{\Sigma}_t}^2 = h^T \mathbf{A}_t^T \mathbf{\Lambda}_t \left( \frac{n}{m}\mathbf{W}_t + \mathbf{P}_t^{(2)} \right)^{-1} \mathbf{\Sigma}_t \left( \frac{n}{m}\mathbf{W}_t + \mathbf{P}_t^{(2)} \right)^{-1} \mathbf{\Lambda}_t \mathbf{A}_t h.$$

Using Lemma 6.13.2, we have  $\max_i \frac{\sigma_{t,i}}{w_{t,i}} \leq \frac{n}{m} \ln \left( \frac{me}{n} \right)$  and hence

$$\begin{aligned} \|z\|_{\mathbf{\Sigma}_t}^2 &\leq \ln \left( \frac{me}{n} \right) h^T \mathbf{A}_t^T \mathbf{\Lambda}_t \left( \frac{n}{m}\mathbf{W}_t + \mathbf{P}_t^{(2)} \right)^{-1} \left( \frac{n}{m}\mathbf{W}_t \right) \left( \frac{n}{m}\mathbf{W}_t + \mathbf{P}_t^{(2)} \right)^{-1} \mathbf{\Lambda}_t \mathbf{A}_t h \\ &\leq \ln \left( \frac{me}{n} \right) h^T \mathbf{A}_t^T \mathbf{\Lambda}_t \left( \frac{n}{m}\mathbf{W}_t + \mathbf{P}_t^{(2)} \right)^{-1} \mathbf{\Lambda}_t \mathbf{A}_t h \leq \frac{\ln \left( \frac{me}{n} \right)}{4} \|h\|_{\nabla^2 p(x_t)}^2. \end{aligned}$$

For the  $\ell^\infty$  bound, note that

$$\left( \frac{n}{m}\mathbf{W}_t + \mathbf{P}_t^{(2)} \right) z = \mathbf{\Sigma}_t \mathbf{A}_t h - \mathbf{P}_t^{(2)} \mathbf{A}_t h.$$

Hence, we have

$$\frac{n}{m} w_{t,i} |z_i| \leq \left| \mathbf{P}_t^{(2)} z \right|_i + \sigma_{t,i} |\mathbf{A}_t h|_i + \left| \mathbf{P}_t^{(2)} \mathbf{A}_t h \right|_i. \quad (6.53)$$

For the term  $\left| \mathbf{P}_t^{(2)} z \right|_i$ , note that

$$\left| \mathbf{P}_t^{(2)} z \right|_i = \sqrt{1_i^T \mathbf{P}_t^{(2)} 1_i} \sqrt{\|z\|_{\mathbf{P}_t^{(2)}}^2} = \sigma_{t,i} \|z\|_{\mathbf{P}_t^{(2)}} \leq \frac{\sigma_{t,i}}{2} \|h\|_{\nabla^2 p(x_t)}.$$

Therefore,  $\left| \mathbf{P}_t^{(2)} z \right|_i \leq \frac{\sigma_{t,i}}{2} \|h\|_{\nabla^2 p(x_t)}$ .

For the term  $\left| \mathbf{P}_t^{(2)} \mathbf{A}_t h \right|_i$ , using  $\nabla^2 p(x_t) \succeq 2\mathbf{A}_t^T \mathbf{\Sigma}_t \mathbf{A}_t$  and  $\mathbf{\Sigma}_t \succeq \mathbf{P}_t^{(2)}$ , we have

$$\begin{aligned} \left| \mathbf{P}_t^{(2)} \mathbf{A}_t h \right|_i &= \sqrt{1_i^T \mathbf{P}_t^{(2)} 1_i} \sqrt{\|\mathbf{A}_t h\|_{\mathbf{P}_t^{(2)}}^2} \\ &\leq \sigma_{t,i} \|\mathbf{A}_t h\|_{\mathbf{\Sigma}_t} \\ &\leq \frac{\sigma_{t,i}}{\sqrt{2}} \|h\|_{\nabla^2 p(x_t)}. \end{aligned}$$

For the term  $|\mathbf{A}_t h|_i$ , Lemma 6.13.4 shows that  $|\mathbf{A}_t h|_i \leq \frac{1}{\sqrt{2}} \ln \left( \frac{me}{n} \right) \|h\|_{\nabla^2 p(x_t)}$ .

Combining these terms into (6.53), we have

$$\frac{n}{m} w_{t,i} |z_i| \leq \sigma_{t,i} \left( \frac{1}{2} + \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{2}} \ln \left( \frac{me}{n} \right) \right) \|h\|_{\nabla^2 p(x_t)}.$$

Using Lemma 6.13.2, we have  $\max_i \frac{\sigma_{t,i}}{w_{t,i}} \leq \frac{n}{m} \ln\left(\frac{me}{n}\right)$  and hence

$$|z_i| \leq \ln\left(\frac{me}{n}\right) \left( \frac{1}{2} + \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{2}} \ln\left(\frac{me}{n}\right) \right) \|h\|_{\nabla^2 p(x_t)}.$$

Using the formula (6.52) for  $w_t$ , we have the result.  $\square$

In the next lemma, we show that  $\Sigma_t$  is stable.

**Lemma 6.13.6.** *We have*

$$\left| \frac{d}{dt} \Sigma_t \right| \preceq 2 \ln\left(\frac{me}{n}\right) \left( 1 + \sqrt{2} \ln\left(\frac{me^2}{n}\right) \right) \|h\|_{\nabla^2 p(x_t)} \Sigma_t.$$

*Proof.* Fix any  $y \in \mathbb{R}^n$ . Let  $\alpha(t) = y^T \mathbf{A}_t^T \Sigma_t \mathbf{A}_t y$ . The optimality condition (6.48) gives

$$\sigma_t = \frac{n}{m} w_t \ln(ew_t).$$

Taking derivative on both side, we have

$$\begin{aligned} \frac{d}{dt} \sigma_t &= \frac{n}{m} \left( \frac{d}{dt} w_t \right) \ln(ew_t) + \frac{n}{m} \left( \frac{d}{dt} w_t \right) \\ &= \frac{n}{m} \left( \frac{d}{dt} w_t \right) \ln(e^2 w_t) \end{aligned}$$

Using Lemma 6.13.5, we have

$$\|\mathbf{W}_t^{-1} \frac{d}{dt} w_t\|_\infty \leq \ln\left(\frac{me}{n}\right) \left( 1 + \sqrt{2} \ln\left(\frac{me^2}{n}\right) \right) \|h\|_{\nabla^2 p(x_t)}.$$

Therefore, for all  $i$ , we have

$$\begin{aligned} \left| \frac{d}{dt} \sigma_{t,i} \right| &\leq \ln\left(\frac{me}{n}\right) \left( 1 + \sqrt{2} \ln\left(\frac{me^2}{n}\right) \right) \frac{n}{m} w_{t,i} \ln(e^2 w_t) \\ &\leq 2 \ln\left(\frac{me}{n}\right) \left( 1 + \sqrt{2} \ln\left(\frac{me^2}{n}\right) \right) \sigma_{t,i} \|h\|_{\nabla^2 p(x_t)}. \end{aligned}$$

This gives the result.  $\square$

Now, we want to prove  $\mathbf{P}_t^{(2)}$  is stable.

**Lemma 6.13.7.** *We have*

$$\left| \frac{d}{dt} y^T \mathbf{P}_t^{(2)} y \right| \preceq \left( \ln\left(\frac{me}{n}\right) \left( 1 + \sqrt{2} \ln\left(\frac{me^3}{n}\right) \right) + 2\sqrt{\ln\left(\frac{m}{n}e\right) + 4} \right) \|y\|_{\Sigma_t}^2 \|h\|_{\nabla^2 p(x_t)}.$$

*Proof.* Fix any  $y \in \mathbb{R}^n$ . Let  $\alpha(t) = y^T \mathbf{P}_t^{(2)} y$ . Then, we have

$$\frac{d\alpha}{dt} = \frac{d}{dt} \sum_{i,j} (\mathbf{P}_t)_{i,j}^2 y_i y_j = 2 \sum_{i,j} (\mathbf{P}_t)_{i,j} \frac{d}{dt} (\mathbf{P}_t)_{i,j} y_i y_j.$$

Let  $z_t = \frac{d}{dt} \ln(w_t/s_t^2)$ . Note that

$$\begin{aligned} \frac{d}{dt} (\mathbf{P}_t)_{i,j} &= \frac{d}{dt} \mathbf{1}_i^T \sqrt{\mathbf{W}_t} \mathbf{A}_t (\mathbf{A}_t^T \mathbf{W}_t \mathbf{A}_t)^{-1} \mathbf{A}_t^T \sqrt{\mathbf{W}_t} \mathbf{1}_j \\ &= \frac{1}{2} \frac{dz_{t,i}}{dt} (\mathbf{P}_t)_{i,j} + \frac{1}{2} \frac{dz_{t,j}}{dt} (\mathbf{P}_t)_{i,j} - \mathbf{1}_i^T \cdot \mathbf{P}_t \cdot \mathbf{Diag} \frac{dz_{t,k}}{dt} \cdot \mathbf{P}_t \cdot \mathbf{1}_j. \end{aligned}$$

Hence, we have

$$\begin{aligned} \frac{d\alpha}{dt} &= \sum_{i,j} (\mathbf{P}_t)_{i,j} \frac{dz_{t,i}}{dt} (\mathbf{P}_t)_{i,j} y_i y_j + \sum_{i,j} (\mathbf{P}_t)_{i,j} \frac{dz_{t,j}}{dt} (\mathbf{P}_t)_{i,j} y_i y_j \\ &\quad - 2 \sum_{i,j} (\mathbf{P}_t)_{i,j} \mathbf{1}_i^T \cdot \mathbf{P}_t \cdot \mathbf{Diag} \left( \frac{dz_{t,k}}{dt} \right) \cdot \mathbf{P}_t \cdot \mathbf{1}_j y_i y_j \\ &= 2y^T \mathbf{P}_t^{(2)} \mathbf{Diag} \left( \frac{dz_{t,j}}{dt} \right) y - 2 \sum_{i,j,k} (\mathbf{P}_t)_{i,j} (\mathbf{P}_t)_{i,k} (\mathbf{P}_t)_{j,k} y_i y_j \frac{dz_{t,k}}{dt}. \end{aligned} \quad (6.54)$$

For the first term  $y^T \mathbf{P}_t^{(2)} \mathbf{diag} \left( \frac{dz_{t,j}}{dt} \right) y$ , using  $\mathbf{P}_t^{(2)} \preceq \Sigma_t$ , we have

$$\left| y^T \mathbf{P}_t^{(2)} \mathbf{Diag} \left( \frac{dz_{t,j}}{dt} \right) y \right| \leq \|y\|_{\mathbf{P}_t^{(2)}} \|\mathbf{Diag} \left( \frac{dz_{t,j}}{dt} \right) y\|_{\mathbf{P}_t^{(2)}} \leq \|y\|_{\Sigma_t}^2 \left\| \frac{dz_t}{dt} \right\|_{\infty}.$$

Using Lemma 6.13.4 and 6.13.5, we have

$$\begin{aligned} \left\| \frac{dz_t}{dt} \right\|_{\infty} &\leq \left\| \frac{d}{dt} \ln(w_t) \right\|_{\infty} + 2 \left\| \frac{d}{dt} \ln(s_t) \right\|_{\infty} \\ &\leq \ln\left(\frac{me}{n}\right) \left( 1 + \sqrt{2} \ln\left(\frac{me^3}{n}\right) \right) \|h\|_{\nabla^2 p(x_t)}. \end{aligned}$$

Hence, we have

$$\left| y^T \mathbf{P}_t^{(2)} \mathbf{Diag} \left( \frac{dz_{t,j}}{dt} \right) y \right| \leq \left( \ln\left(\frac{me}{n}\right) \left( 1 + \sqrt{2} \ln\left(\frac{me^3}{n}\right) \right) \right) \|y\|_{\Sigma_t}^2 \|h\|_{\nabla^2 p(x_t)}.$$

For the second term  $\sum_{i,j,k} (\mathbf{P}_t)_{i,j} (\mathbf{P}_t)_{i,k} (\mathbf{P}_t)_{j,k} y_i y_j \frac{d}{dt} z_{t,k}$ , we let  $a_{i,j} = (\mathbf{P}_t)_{i,j} \sqrt{|y_i y_j|}$ ,  $b_{i,k} = (\mathbf{P}_t)_{i,k} \sqrt{|y_i \frac{d}{dt} z_{t,k}|}$  and  $c_{j,k} = (\mathbf{P}_t)_{j,k} \sqrt{|y_j \frac{d}{dt} z_{t,k}|}$  and get

$$\left| \sum_{i,j,k} (\mathbf{P}_t)_{i,j} (\mathbf{P}_t)_{i,k} (\mathbf{P}_t)_{j,k} y_i y_j \frac{d}{dt} z_{t,k} \right| \leq \sum_{i,j,k} |a_{i,j} b_{i,k} c_{j,k}|.$$

Using Cauchy Schwarz Inequality twice and  $\mathbf{P}_t^{(2)} \preceq \Sigma_t$ , we have

$$\begin{aligned} \left| \sum_{i,j,k} (\mathbf{P}_t)_{i,j} (\mathbf{P}_t)_{i,k} (\mathbf{P}_t)_{j,k} y_i y_j \frac{d}{dt} z_{t,k} \right| &\leq \sqrt{\sum_{i,j} a_{i,j}^2} \sqrt{\sum_{i,k} b_{i,k}^2} \sqrt{\sum_{j,k} c_{j,k}^2} \\ &= \| |y| \|_{\mathbf{P}_t^{(2)}} \left\langle |y|, \left| \frac{d}{dt} z_t \right| \right\rangle_{\mathbf{P}_t^{(2)}}^2 \\ &\leq \|y\|_{\Sigma_t}^2 \left\| \frac{d}{dt} z_t \right\|_{\Sigma_t}. \end{aligned}$$

Using Lemma 6.13.5, we have

$$\left| \sum_{i,j,k} (\mathbf{P}_t)_{i,j} (\mathbf{P}_t)_{i,k} (\mathbf{P}_t)_{j,k} y_i y_j \frac{d}{dt} z_{t,k} \right| \leq \left( \sqrt{\ln \left( \frac{m}{n} e \right)} + 2 \right) \|y\|_{\Sigma_t}^2 \|h\|_{\nabla^2 p(x_t)}.$$

Combining two cases into (6.54), we have

$$\left| \frac{d\alpha}{dt} \right| \leq \left( \ln \left( \frac{me}{n} \right) \left( 1 + \sqrt{2} \ln \left( \frac{me^3}{n} \right) \right) + 2 \sqrt{\ln \left( \frac{m}{n} e \right)} + 4 \right) \|y\|_{\Sigma_t}^2 \|h\|_{\nabla^2 p(x_t)}.$$

□

### ■ 6.13.4 $\nabla^2 p$ is stable

Now, we can combine all lemmas before to prove the wanted statement.

**Lemma 6.13.8.** *We have*

$$\left| \frac{d}{dt} \nabla^2 p \right| \preceq O \left( \ln^3 \left( \frac{me}{n} \right) \right) \|h\|_{\nabla^2 p(x_t)} \nabla^2 p.$$

*Proof.* Note that

$$\nabla^2 p(x_t) = \mathbf{A}_t^T \left( 6\Sigma_t - 4\mathbf{P}_t^{(2)} + \frac{n}{m} I \right) \mathbf{A}_t + 4\mathbf{A}_t^T \Lambda_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \Lambda_t \mathbf{A}_t.$$

□

Fix any  $y \in \mathbb{R}^n$ . Let  $\alpha(t) = y^T \nabla^2 p(x_t) y$ . Then, we have

$$\begin{aligned} \left| \frac{d\alpha}{dt} \right| &\leq 6 \left| \frac{d}{dt} y^T \mathbf{A}_t^T \Sigma_t \mathbf{A}_t y \right| + 4 \left| \frac{d}{dt} y^T \mathbf{A}_t^T \mathbf{P}_t^{(2)} \mathbf{A}_t y \right| \\ &\quad + 4 \left| \frac{d}{dt} \mathbf{A}_t^T \Lambda_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \Lambda_t \mathbf{A}_t y \right| \\ &\leq 12 \left| y^T \mathbf{A}_t^T \Sigma_t \left( \frac{d}{dt} \mathbf{A}_t \right) y \right| + 6 \left| y^T \mathbf{A}_t^T \left( \frac{d}{dt} \Sigma_t \right) \mathbf{A}_t y \right| \\ &\quad + 8 \left| y^T \mathbf{A}_t^T \mathbf{P}_t^{(2)} \left( \frac{d}{dt} \mathbf{A}_t \right) y \right| + 4 \left| y^T \mathbf{A}_t^T \left( \frac{d}{dt} \mathbf{P}_t^{(2)} \right) \mathbf{A}_t y \right| \\ &\quad + 8 \left| y^T \mathbf{A}_t^T \Lambda_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \Lambda_t \left( \frac{d}{dt} \mathbf{A}_t \right) y \right| \\ &\quad + 8 \left| y^T \mathbf{A}_t^T \Lambda_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \left( \frac{d}{dt} \mathbf{P}_t^{(2)} \right) \mathbf{A}_t y \right| \\ &\quad + 8 \left| y^T \mathbf{A}_t^T \Lambda_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \left( \frac{d}{dt} \Sigma_t \right) \mathbf{A}_t y \right| \\ &\quad + 4 \left| y^T \mathbf{A}_t^T \Lambda_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \left( \frac{d}{dt} \mathbf{P}_t^{(2)} \right) \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \Lambda_t \mathbf{A}_t y \right| \\ &\quad + 4 \left| y^T \mathbf{A}_t^T \Lambda_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \left( \frac{d}{dt} \frac{n}{m} \mathbf{W}_t \right) \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \Lambda_t \mathbf{A}_t y \right|. \end{aligned}$$

Now, we need to bound the terms above one by one. Using Lemma 6.13.4, we have

$$\begin{aligned} \left| y^T \mathbf{A}_t^T \Sigma_t \left( \frac{d}{dt} \mathbf{A}_t \right) y \right| &\leq O\left(\ln\left(\frac{me}{n}\right)\right) \|h\|_{\nabla^2 p(x_t)} |y^T \mathbf{A}_t^T \Sigma_t \mathbf{A}_t y| \\ &= O\left(\ln\left(\frac{me}{n}\right)\right) \|h\|_{\nabla^2 p(x_t)} \|y\|_{\nabla^2 p(x_t)}^2. \end{aligned}$$

Using Lemma 6.13.6, we have

$$\begin{aligned} \left| y^T \mathbf{A}_t^T \left( \frac{d}{dt} \Sigma_t \right) \mathbf{A}_t y \right| &\leq O\left(\ln^2\left(\frac{me}{n}\right)\right) \|h\|_{\nabla^2 p(x_t)} |y^T \mathbf{A}_t^T \Sigma_t \mathbf{A}_t y| \\ &\leq O\left(\ln^2\left(\frac{me}{n}\right)\right) \|h\|_{\nabla^2 p(x_t)} \|y\|_{\nabla^2 p(x_t)}^2. \end{aligned}$$

Using Lemma 6.13.4, we have

$$\begin{aligned} \left| y^T \mathbf{A}_t^T \mathbf{P}_t^{(2)} \left( \frac{d}{dt} \mathbf{A}_t \right) y \right| &\leq \sqrt{y^T \mathbf{A}_t^T \mathbf{P}_t^{(2)} \mathbf{A}_t y} \sqrt{y^T \left( \frac{d}{dt} \mathbf{A}_t \right) \mathbf{P}_t^{(2)} \left( \frac{d}{dt} \mathbf{A}_t \right) y} \\ &\leq O(1) \sqrt{\|y\|_{\nabla^2 p(x_t)}^2} \sqrt{\left(\ln\left(\frac{me}{n}\right)\right)^2 \|h\|_{\nabla^2 p(x_t)}^2 \|y\|_{\nabla^2 p(x_t)}^2} \\ &= O\left(\ln\left(\frac{me}{n}\right)\right) \|h\|_{\nabla^2 p(x_t)} \|y\|_{\nabla^2 p(x_t)}^2. \end{aligned}$$

Using Lemma 6.13.7, we have

$$\begin{aligned} \left| y^T \mathbf{A}_t^T \left( \frac{d}{dt} \mathbf{P}_t^{(2)} \right) \mathbf{A}_t y \right| &\leq O\left(\ln^2\left(\frac{me}{n}\right)\right) \|\mathbf{A}_t y\|_{\Sigma_t}^2 \|h\|_{\nabla^2 p(x_t)} \\ &= O\left(\ln^2\left(\frac{me}{n}\right)\right) \|h\|_{\nabla^2 p(x_t)} \|y\|_{\nabla^2 p(x_t)}^2. \end{aligned}$$

Using Lemma 6.13.4, we have

$$\begin{aligned} &\left| y^T \mathbf{A}_t^T \Lambda_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \Lambda_t \left( \frac{d}{dt} \mathbf{A}_t \right) y \right| \\ &\leq \sqrt{y^T \mathbf{A}_t^T \Lambda_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \Lambda_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \mathbf{A}_t^T y} \\ &\quad \times \sqrt{y^T \left( \frac{d}{dt} \mathbf{A}_t \right)^T \Lambda_t \left( \frac{d}{dt} \mathbf{A}_t \right) y} \\ &= O(1) \sqrt{\ln\left(\frac{me}{n}\right) \|y\|_{\nabla^2 p(x_t)}^2} \sqrt{y^T \left( \frac{d}{dt} \mathbf{A}_t \right)^T \Sigma_t \left( \frac{d}{dt} \mathbf{A}_t \right) y} \\ &= O\left(\ln^{3/2}\left(\frac{me}{n}\right)\right) \|h\|_{\nabla^2 p(x_t)} \|y\|_{\nabla^2 p(x_t)}^2 \end{aligned}$$



Using Lemma 6.13.6, we have

$$\begin{aligned}
& \left| y^T \mathbf{A}_t^T \boldsymbol{\Lambda}_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \left( \frac{d}{dt} \boldsymbol{\Sigma}_t \right) \mathbf{A}_t y \right| \\
& \leq \sqrt{y^T \mathbf{A}_t^T \boldsymbol{\Lambda}_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \left( \frac{d}{dt} \boldsymbol{\Sigma}_t \right) \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \mathbf{A}_t^T y} \\
& \quad \times \sqrt{y^T \mathbf{A}_t^T \left( \frac{d}{dt} \boldsymbol{\Sigma}_t \right) \mathbf{A}_t y} \\
& \leq O \left( \ln^{3/2} \left( \frac{me}{n} \right) \right) \|h\|_{\nabla^2 p(x_t)} \|y\|_{\nabla^2 p(x_t)}^2
\end{aligned}$$

Let  $z = \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \boldsymbol{\Lambda}_t \mathbf{A}_t y$ , using Lemma 6.13.7, we have

$$\begin{aligned}
& \left| y^T \mathbf{A}_t^T \boldsymbol{\Lambda}_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \left( \frac{d}{dt} \mathbf{P}_t^{(2)} \right) \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \boldsymbol{\Lambda}_t \mathbf{A}_t y \right| \\
& \leq O \left( \ln^2 \left( \frac{me}{n} \right) \right) \|z\|_{\boldsymbol{\Sigma}_t}^2 \|h\|_{\nabla^2 p(x_t)}.
\end{aligned}$$

Note that

$$\begin{aligned}
\|z\|_{\boldsymbol{\Sigma}_t}^2 &= y^T \mathbf{A}_t^T \boldsymbol{\Lambda}_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \boldsymbol{\Sigma}_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \boldsymbol{\Lambda}_t \mathbf{A}_t y \\
&\leq \ln \left( \frac{me}{n} \right) y^T \mathbf{A}_t^T \boldsymbol{\Lambda}_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \boldsymbol{\Lambda}_t \mathbf{A}_t y \leq O \left( \ln \left( \frac{me}{n} \right) \right) \|y\|_{\nabla^2 p(x_t)}^2.
\end{aligned}$$

Hence, we have

$$\begin{aligned}
& \left| y^T \mathbf{A}_t^T \boldsymbol{\Lambda}_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \left( \frac{d}{dt} \mathbf{P}_t^{(2)} \right) \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \boldsymbol{\Lambda}_t \mathbf{A}_t y \right| \\
& \leq O \left( \ln^3 \left( \frac{me}{n} \right) \right) \|y\|_{\nabla^2 p(x_t)}^2 \|h\|_{\nabla^2 p(x_t)}.
\end{aligned}$$

Using Lemma 6.13.5, we have

$$\begin{aligned}
& \left| y^T \mathbf{A}_t^T \boldsymbol{\Lambda}_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \left( \frac{d}{dt} \frac{n}{m} \mathbf{W}_t \right) \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \boldsymbol{\Lambda}_t \mathbf{A}_t y \right| \\
& \leq O \left( \ln^2 \left( \frac{me}{n} \right) \right) \|y\|_{\nabla^2 p(x_t)}^2 \|h\|_{\nabla^2 p(x_t)}
\end{aligned}$$

Using Lemma 6.13.7, we have

$$\begin{aligned}
& \left| y^T \mathbf{A}_t^T \boldsymbol{\Lambda}_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \left( \frac{d}{dt} \mathbf{P}_t^{(2)} \right) \mathbf{A}_t y \right| \\
& \leq \sqrt{y^T \mathbf{A}_t^T \boldsymbol{\Lambda}_t \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \left( \frac{d}{dt} \mathbf{P}_t^{(2)} \right) \left( \mathbf{P}_t^{(2)} + \frac{n}{m} \mathbf{W}_t \right)^{-1} \boldsymbol{\Lambda}_t \mathbf{A}_t y} \sqrt{y^T \mathbf{A}_t^T \left( \frac{d}{dt} \mathbf{P}_t^{(2)} \right) \mathbf{A}_t y} \\
& \leq \sqrt{O \left( \ln^3 \left( \frac{me}{n} \right) \right) \|y\|_{\nabla^2 p(x_t)}^2 \|h\|_{\nabla^2 p(x_t)}} \sqrt{O \left( \ln^2 \left( \frac{me}{n} \right) \|y\|_{\nabla^2 p(x_t)}^2 \|h\|_{\nabla^2 p(x_t)} \right)} \\
& = O \left( \ln^{5/2} \left( \frac{me}{n} \right) \right) \|y\|_{\nabla^2 p(x_t)}^2 \|h\|_{\nabla^2 p(x_t)}.
\end{aligned}$$

Combining all terms together, we have  $\left| \frac{d\alpha}{dt} \right| = O \left( \ln^3 \left( \frac{me}{n} \right) \right) \|y\|_{\nabla^2 p(x_t)}^2 \|h\|_{\nabla^2 p(x_t)}$ .



## Geometric Median In Nearly-Linear Time

### ■ 7.1 Introduction

In this chapter, we study Fermat-Weber problem, one of the oldest easily-stated nontrivial problems in computational geometry. Given a set of  $n$  points in  $d$  dimensions,  $a^{(1)}, \dots, a^{(n)} \in \mathbb{R}^d$ , find a point  $x_* \in \mathbb{R}^d$  that minimizes the sum of Euclidean distances to them:

$$x_* \in \arg \min_{x \in \mathbb{R}^d} f(x) \quad \text{where} \quad f(x) \stackrel{\text{def}}{=} \sum_{i \in [n]} \|x - a^{(i)}\|_2$$

This problem, also known as the *geometric median problem*, is well studied and has numerous applications. It is often considered over low dimensional spaces in the context of the facility location problem [267] and over higher dimensional spaces it has applications to clustering in machine learning and data analysis. For example, computing the geometric median is a subroutine in popular expectation maximization heuristics for  $k$ -medians clustering.

The problem is also important to robust estimation, where we like to find a point representative of given set of points that is resistant to outliers. The geometric median is a rotation and translation invariant estimator that achieves the optimal *breakdown point* of 0.5, i.e. it is a good estimator even when up to half of the input data is arbitrarily corrupted [172]. Moreover, if a large constant fraction of the points lie in a ball of diameter  $\varepsilon$  then the geometric median lies in that ball with diameter  $O(\varepsilon)$  (see Lemma 7.7.1). Consequently, the geometric median can be used to turn expected results into high probability results: e.g. if the  $a^{(i)}$  are drawn independently such that  $\mathbf{E}\|x - x_*\|_2 \leq \varepsilon$  for some  $\varepsilon > 0$  and  $x \in \mathbb{R}^d$  then this fact, Markov bound, and Chernoff Bound, imply  $\|x_* - a^{(i)}\|_2 = O(\varepsilon)$  with high probability in  $n$ .

Despite the ancient nature of the Fermat-Weber problem and its many uses there are relatively few theoretical guarantees for solving it (see Table 7.1). To compute a  $(1 + \varepsilon)$ -approximate solution, i.e.  $x \in \mathbb{R}^d$  with  $f(x) \leq (1 + \varepsilon)f(x_*)$ , the previous fastest running times were either  $O(d \cdot n^{4/3} \varepsilon^{-8/3})$  by [50],  $\tilde{O}(d \exp \varepsilon^{-4} \log \varepsilon^{-1})$  by [25],  $\tilde{O}(nd + \text{poly}(d, \varepsilon^{-1}))$  by [86], or  $O((nd)^{O(1)} \log \frac{1}{\varepsilon})$  time by [222, 271]. In this chapter, we improve upon these running times by providing an  $O(nd \log^3 \frac{n}{\varepsilon})$  time algorithm<sup>1</sup> as well as an  $O(d/\varepsilon^2)$  time algorithm, provided we have an oracle for sampling a random  $a^{(i)}$ . Picking the faster algorithm for the particular value of  $\varepsilon$  improves the running time to  $O(nd \log^3 \frac{1}{\varepsilon})$ . We also extend these results to compute a  $(1 + \varepsilon)$ -approximate solution to the more general Weber's problem,  $\min_{x \in \mathbb{R}^d} \sum_{i \in [n]} w_i \|x - a^{(i)}\|_2$  for non-negative  $w_i$ , in time  $O(nd \log^3 \frac{1}{\varepsilon})$  (see Section 7.10).

Our  $O(nd \log^3 \frac{n}{\varepsilon})$  time algorithm is a careful modification of standard interior point methods for solving the geometric median problem. We provide a long step interior point method tailored to the geometric median problem for which we can implement every iteration in nearly linear time. While our analysis starts with a simple  $O((nd)^{O(1)} \log \frac{1}{\varepsilon})$  time interior point method and shows how to improve it,

<sup>1</sup>If  $z$  is the total number of nonzero entries in the coordinates of the  $a^{(i)}$  then a careful analysis of our algorithm improves our running time to  $O(z \log^3 \frac{n}{\varepsilon})$ .

our final algorithm is quite non-standard from the perspective of interior point literature. Our result is one of very few cases we are aware of outperforming traditional interior point theory (Chapter 6 and [183]) and the only we are aware of using interior point methods to obtain a nearly linear time algorithm for a canonical optimization problem that traditionally requires superlinear time. We hope our work leads to further improvements in this line of research.

Our  $O(d\varepsilon^{-2})$  algorithm is a relatively straightforward application of sampling techniques and stochastic subgradient descent. Some additional insight is required simply to provide a rigorous analysis of the robustness of the geometric median and use this to streamline our application of stochastic subgradient descent. We include it for completeness however, we defer its proof to Section 7.7. The bulk of the work in this chapter is focused on developing our  $O(nd \log^3 \frac{n}{\varepsilon})$  time algorithm which we believe uses a set of techniques of independent interest.

### ■ 7.1.1 Previous Work

The geometric median problem was first formulated for the case of three points in the early 1600s by Pierre de Fermat [148, 69]. A simple elegant ruler and compass construction was given in the same century by Evangelista Torricelli. Such a construction does not generalize when a larger number of points is considered: Bajaj has shown the even for five points, the geometric median is not expressible by radicals over the rationals [27]. Hence, the  $(1 + \varepsilon)$ -approximate problem has been studied for larger values of  $n$ .

Many authors have proposed algorithms with runtime polynomial in  $n$ ,  $d$  and  $1/\varepsilon$ . The most cited and used algorithm is Weiszfeld's 1937 algorithm [268]. Unfortunately Weiszfeld's algorithm may not converge and if it does it may do so very slowly. There have been many proposed modifications to Weiszfeld's algorithm [58, 224, 220, 28, 260, 155] that generally give non-asymptotic runtime guarantees. In light of more modern multiplicative weights methods his algorithm can be viewed as a re-weighted least squares iteration. Chin et al. [50] considered the more general  $L_2$  embedding problem: placing the vertices of a graph into  $\mathbb{R}^d$ , where some of the vertices have fixed positions while the remaining vertices are allowed to float, with the objective of minimizing the sum of the Euclidean edge lengths. Using the multiplicative weights method, they obtained a run time of  $O(d \cdot n^{4/3} \varepsilon^{-8/3})$  for a broad class of problems, including the geometric median problem.<sup>2</sup>

Many authors consider problems that generalize the Fermat-Weber problem, and obtain algorithms for finding the geometric median as a specialization. For example, Badoiu et al. give an approximate  $k$ -median algorithm by sub-sampling with the runtime for  $k = 1$  of  $\tilde{O}(d \cdot \exp(O(\varepsilon^{-4})))$  [25]. Parrilo and Sturmfels demonstrated that the problem can be reduced to semidefinite programming, thus obtaining a runtime of  $\tilde{O}(\text{poly}(n, d) \log \varepsilon^{-1})$  [222]. Furthermore, Bose et al. gave a linear time algorithm for fixed  $d$  and  $\varepsilon^{-1}$ , based on low-dimensional data structures [37] and it has been show how to obtain running times of  $\tilde{O}(nd + \text{poly}(d, \varepsilon^{-1}))$  for this problem and a more general class of problems.[113, 86].

An approach very related to ours was studied by Xue and Ye [271]. They give an interior point method with barrier analysis that runs in time  $\tilde{O}((d^3 + d^2n)\sqrt{n} \log \varepsilon^{-1})$ .

### ■ 7.1.2 Overview of $O(nd \log^3 \frac{n}{\varepsilon})$ Time Algorithm

#### Interior Point Primer

Our algorithm is broadly inspired by interior point methods, a broad class of methods for efficiently solving convex optimization problems [273, 211]. Given an instance of the geometric median problem

---

<sup>2</sup>The result of [50] was stated in more general terms than given here. However, it easy to formulate the geometric median problem in their model.

Year	Authors	Runtime	Comments
1659	Torricelli [263]	-	Assuming $n = 3$
1937	Weiszfeld [268]	-	Does not always converge
1990	Chandrasekaran and Tamir [47]	$\tilde{O}(n \cdot \text{poly}(d) \log \varepsilon^{-1})$	Ellipsoid method
1997	Xue and Ye [271]	$\tilde{O}((d^3 + d^2n) \sqrt{n} \log \varepsilon^{-1})$	Interior point with barrier method
2000	Indyk [116]	$\tilde{O}(dn \cdot \varepsilon^{-2})$	Optimizes only over $x$ in the input
2001	Parrilo and Sturmfels [222]	$\tilde{O}(\text{poly}(n, d) \log \varepsilon^{-1})$	Reduction to SDP
2002	Badoiu et al. [25]	$\tilde{O}(d \cdot \exp(O(\varepsilon^{-4})))$	Sampling
2003	Bose et al. [37]	$\tilde{O}(n)$	Assuming $d, \varepsilon^{-1} = O(1)$
2005	Har-Peled and Kushal [113]	$\tilde{O}(n + \text{poly}(\varepsilon^{-1}))$	Assuming $d = O(1)$
2011	Feldman and Langberg [86]	$\tilde{O}(nd + \text{poly}(d, \varepsilon^{-1}))$	Coreset
2013	Chin et al. [50]	$\tilde{O}(dn^{4/3} \cdot \varepsilon^{-8/3})$	Multiplicative weights
-	<b>This Chapter</b>	$O(nd \log^3(n/\varepsilon))$	Interior point with custom analysis
-	<b>This Chapter</b>	$O(d\varepsilon^{-2})$	Stochastic gradient descent

Table 7.1: Selected Previous Results.

we first put the problem in a more natural form for applying interior point methods. Rather than writing the problem as minimizing a convex function over  $\mathbb{R}^d$

$$\min_{x \in \mathbb{R}^d} f(x) \quad \text{where} \quad f(x) \stackrel{\text{def}}{=} \sum_{i \in [n]} \|x - a^{(i)}\|_2 \quad (7.1)$$

we instead write the problem as minimizing a linear function over a larger convex space:

$$\min_{\{\alpha, x\} \in S} 1^\top \alpha \quad \text{where} \quad S = \left\{ \alpha \in \mathbb{R}^n, x \in \mathbb{R}^d \mid \|x^{(i)} - a^{(i)}\|_2 \leq \alpha_i \text{ for all } i \in [n] \right\}. \quad (7.2)$$

Clearly, these problems are the same as at optimality  $\alpha_i = \|x^{(i)} - a^{(i)}\|_2$ .

To solve problems of the form (7.2) interior point methods relax the constraint  $\{\alpha, x\} \in S$  through the introduction of a *barrier function*. In particular they assume that there is a real valued function  $p$  such that as  $\{\alpha, x\}$  moves towards the boundary of  $S$  the value of  $p$  goes to infinity. A popular class of interior point methods known as *path following methods* [227, 107], they consider relaxations of (7.2) of the form  $\min_{\{\alpha, x\} \in \mathbb{R}^n \times \mathbb{R}^d} t \cdot 1^\top \alpha + p(\alpha, x)$ . The minimizers of this function form a path, known as the central path, parameterized by  $t$ . The methods then use variants of Newton's method to follow the path until  $t$  is large enough that a high quality approximate solution is obtained. The number of iterations of these methods are then typically governed by a property of  $p$  known as its self concordance  $\nu$ . Given a  $\nu$ -self concordant barrier, typically interior point methods require  $O(\sqrt{\nu} \log \frac{1}{\varepsilon})$  iterations to compute a  $(1 + \varepsilon)$ -approximate solution.

For our particular convex set, the construction of our barrier function is particularly simple, we consider each constraint  $\|x - a^{(i)}\|_2 \leq \alpha_i$  individually. In particular, it is known that the function  $p^{(i)}(\alpha, x) = -\ln \alpha_i - \ln(\alpha_i^2 - \|x - a^{(i)}\|_2)$  is a  $O(1)$  self-concordant barrier function for the set  $S^{(i)} = \{x \in \mathbb{R}^d, \alpha \in \mathbb{R}^n \mid \|x - a^{(i)}\|_2 \leq \alpha_i\}$  [206]. Since  $\cap_{i \in [n]} S^{(i)} = S$  we can use the barrier  $\sum_{i \in [n]} p^{(i)}(\alpha, x)$  for  $p(\alpha, x)$  and standard self-concordance theory shows that this is an  $O(n)$  self concordant barrier for  $S$ . Consequently, this easily yields an interior point method for solving the geometric median problem in  $O((nd)^{O(1)} \log \frac{1}{\varepsilon})$  time.

## Difficulties

Unfortunately obtaining a nearly linear time algorithm for geometric median using interior point methods as presented poses numerous difficulties. Particularly troubling is the number of iterations required by standard interior point algorithms. The approach outlined in the previous section produced an  $O(n)$ -self concordant barrier and even if we use more advanced self concordance machinery, i.e. the universal barrier [211], the best known self concordance of barrier for the convex set  $\sum_{i \in [n]} \|x - a^{(i)}\|_2 \leq c$  is  $O(d)$ . An interesting open question still left open by our work is to determine what is the minimal self concordance of a barrier for this set.

Consequently, even if we could implement every iteration of an interior point scheme in nearly linear time it is unclear whether one should hope for a nearly linear time interior point algorithm for the geometric median. While there are a instances of outperforming standard self-concordance analysis (Chapter 6), these instances are few, complex, and to varying degrees specialized to the problems they solve. Moreover, we are unaware of any interior point scheme providing a provable nearly linear time for a general nontrivial convex optimization problem.

## Beyond Standard Interior Point

Despite these difficulties we do obtain a nearly linear time interior point based algorithm that only requires  $O(\log \frac{n}{\varepsilon})$  iterations, i.e. increases to the path parameter. After choosing the natural penalty functions  $p^{(i)}$  described above, we optimize in closed form over the  $\alpha_i$  to obtain the following penalized objective function:<sup>3</sup>

$$\min_x f_t(x) \quad \text{where} \quad f_t(x) = \sum_{i \in [n]} \sqrt{1 + t^2 \|x - a^{(i)}\|_2^2} - \ln \left[ 1 + \sqrt{1 + t^2 \|x - a^{(i)}\|_2^2} \right]$$

We then approximately minimize  $f_t(x)$  for increasing  $t$ . We let  $x_t \stackrel{\text{def}}{=} \arg \min_{x \in \mathbb{R}^d} f_t(x)$  for  $x \geq 0$ , and thinking of  $\{x_t : t \geq 0\}$  as a continuous curve known as the central path, we show how to approximately follow this path. As  $\lim_{t \rightarrow \infty} x_t = x_*$  this approach yields a  $(1 + \varepsilon)$ -approximation.

So far our analysis is standard and interior point theory yields an  $\Omega(\sqrt{n})$  iteration interior point scheme. To overcome this we take a more detailed look at  $x_t$ . We note that for any  $t$  if there is any rapid change in  $x_t$  it must occur in the direction of the smallest eigenvector of  $\nabla^2 f_t(x)$ , denoted  $v_t$ , what we henceforth may refer to as the *bad direction* at  $x_t$ . Furthermore, we show that this bad direction is *stable* in the sense that for all directions  $d \perp v_t$  it is the case that  $d^\top(x_t - x_{t'})$  is small for  $t' \leq ct$  for a small constant  $c$ .

In fact, we show that this movement over such a *long step*, i.e. constant increase in  $t$ , in the directions orthogonal to the bad direction is small enough that for any movement around a ball of this size the Hessian of  $f_t$  only changes by a small multiplicative constant. In short, starting at  $x_t$  there exists a point  $y$  obtained just by moving from  $x_t$  in the bad direction, such that  $y$  is close enough to  $x_{t'}$  that standard first order method will converge quickly to  $x_{t'}$ ! Thus, we might hope to find such a  $y$ , quickly converge to  $x_{t'}$  and repeat. If we increase  $t$  by a multiplicative constant in every such iterations, standard interior point theory suggests that  $O(\log \frac{n}{\varepsilon})$  iterations suffices.

## Building an Algorithm

To turn the structural result in the previous section into a fast algorithm there are several further issues we need to address. We need to

<sup>3</sup>It is unclear how to extend our proof for the simpler function:  $\sum_{i \in [n]} \sqrt{1 + t^2 \|x - a^{(i)}\|_2^2}$ .

- (1) Show how to find the point along the bad direction that is close to  $x_{t'}$
- (2) Show how to solve linear systems in the Hessian to actually converge quickly to  $x_{t'}$
- (3) Show how to find the bad direction
- (4) Bound the accuracy required by these computations

Deferring (1) for the moment, our solution to the rest are relatively straightforward. Careful inspection of the Hessian of  $f_t$  reveals that it is well approximated by a multiple of the identity matrix minus a rank 1 matrix. Consequently using explicit formulas for the inverse of of matrix under rank 1 updates, i.e. the Sherman-Morrison formula, we can solve such systems in nearly linear time thereby addressing (2). For (3), we show that the well known power method carefully applied to the Hessian yields the bad direction if it exists. Finally, for (4) we show that a constant approximate geometric median is near enough to the central path for  $t = \Theta(\frac{1}{f(x_*)})$  and that it suffices to compute a central path point at  $t = O(\frac{n}{f(x_*)\varepsilon})$  to compute a  $1 + \varepsilon$ -geometric median. Moreover, for these values of  $t$ , the precision needed in other operations is clear.

The more difficult operation is (1). Given  $x_t$  and the bad direction exactly, it is still not clear how to find the point along the bad direction line from  $x_t$  that is close to  $x_{t'}$ . Just performing binary search on the objective function a priori might not yield such a point due to discrepancies between a ball in Euclidean norm and a ball in hessian norm and the size of the distance from the optimal point in euclidean norm. To overcome this issue we still line search on the bad direction, however rather than simply using  $f(x_t + \alpha \cdot v_t)$  as the objective function to line search on, we use the function  $g(\alpha) = \min_{\|x - x_t - \alpha \cdot v_t\|_2 \leq c} f(x)$  for some constant  $c$ , that is given an  $\alpha$  we move  $\alpha$  in the bad direction and take the best objective function value in a ball around that point. For appropriate choice of  $c$  the minimizers of  $\alpha$  will include the optimal point we are looking for. Moreover, we can show that  $g$  is convex and that it suffices to perform the minimization approximately.

Putting these pieces together yields our result. We perform  $O(\log \frac{n}{\varepsilon})$  iterations of interior point (i.e. increasing  $t$ ), where in each iteration we spend  $O(nd \log \frac{n}{\varepsilon})$  time to compute a high quality approximation to the bad direction, and then we perform  $O(\log \frac{n}{\varepsilon})$  approximate evaluations on  $g(\alpha)$  to binary search on the bad direction line, and then to approximately evaluate  $g$  we perform gradient descent in approximate Hessian norm to high precision which again takes  $O(nd \log \frac{n}{\varepsilon})$  time. Altogether this yields a  $O(nd \log^3 \frac{n}{\varepsilon})$  time algorithm to compute a  $1 + \varepsilon$  geometric median. Here we made minimal effort to improve the log factors and plan to investigate this further in future work.

### ■ 7.1.3 Overview of $O(d\varepsilon^{-2})$ Time Algorithm

In addition to providing a nearly linear time algorithm we provide a stand alone result on quickly computing a crude  $(1 + \varepsilon)$ -approximate geometric median in Section 7.7. In particular, given an oracle for sampling a random  $a^{(i)}$  we provide an  $O(d\varepsilon^{-2})$ , i.e. sublinear, time algorithm that computes such an approximate median. Our algorithm for this result is fairly straightforward. First, we show that random sampling can be used to obtain some constant approximate information about the optimal point in constant time. In particular we show how this can be used to deduce an Euclidean ball which contains the optimal point. Second, we perform stochastic subgradient descent within this ball to achieve our desired result.

### ■ 7.1.4 Organization

The rest of the chapter is structured as follows. After covering preliminaries in Section 2.3, in Section 7.3 we provide various results about the central path that we use to derive our nearly linear time

algorithm. In Section 7.4 we then provide our nearly linear time algorithm. All the proofs and supporting lemmas for these sections are deferred to Section 7.5 and Section 7.6. In Section 7.7 we provide our  $O(d/\varepsilon^2)$  algorithm, in Section 7.8 we provide the derivation of our penalized objective function, in Section 7.9 we provide general technical machinery we use throughout, and in Section 7.10 we show how to extend our results to Weber's problem, i.e. weighted geometric median.

## ■ 7.2 Preliminaries

### ■ 7.2.1 General Notation

For a symmetric positive semidefinite matrix (PSD),  $\mathbf{A}$ , we let  $\lambda_1(\mathbf{A}) \geq \dots \geq \lambda_n(\mathbf{A}) \geq 0$  denote the eigenvalues of  $\mathbf{A}$  and let  $v_1(\mathbf{A}), \dots, v_n(\mathbf{A})$  denote corresponding eigenvectors.

### ■ 7.2.2 Problem Notation

The central problem of this chapter is as follows: we are given points  $a^{(1)}, \dots, a^{(n)} \in \mathbb{R}^d$  and we wish to compute a geometric median, i.e.  $x_* \in \arg \min_{x \in \mathbb{R}^d} f(x)$  where  $f(x) = \sum_{i \in [n]} \|a^{(i)} - x\|_2$ . We call a point  $x \in \mathbb{R}^d$  an  $(1 + \varepsilon)$ -approximate geometric median if  $f(x) \leq (1 + \varepsilon)f(x_*)$ .

### ■ 7.2.3 Penalized Objective Notation

To solve this problem we relax the objective function  $f$  and instead consider the following family of *penalized objective functions* parameterized by  $t > 0$

$$\min_{x \in \mathbb{R}^d} f_t(x) \quad \text{where} \quad f_t(x) = \sum_{i \in [n]} \sqrt{1 + t^2 \|x - a^{(i)}\|_2^2} - \ln \left[ 1 + \sqrt{1 + t^2 \|x - a^{(i)}\|_2^2} \right]$$

This penalized objective function is derived from a natural interior point formulation of the geometric median problem (See Section 7.8). For all *path parameters*  $t > 0$ , we let  $x_t \stackrel{\text{def}}{=} \arg \min_x f_t(x)$ . Our primary goal is to obtain good approximations to the *central path*  $\{x_t : t > 0\}$  for increasing values of  $t$ .

We let  $g_t^{(i)}(x) \stackrel{\text{def}}{=} \sqrt{1 + t^2 \|x - a^{(i)}\|_2^2}$  and  $f_t^{(i)}(x) \stackrel{\text{def}}{=} g_t^{(i)}(x) - \ln(1 + g_t^{(i)}(x))$  so  $f_t(x) = \sum_{i \in [n]} f_t^{(i)}(x)$ . We refer to the quantity  $w_t(x) \stackrel{\text{def}}{=} \sum_{i \in [n]} \frac{1}{1 + g_t^{(i)}(x)}$  as *weight* as it is a natural measure of total contribution of the  $a^{(i)}$  to  $\nabla^2 f_t(x)$ . We let

$$\bar{g}_t(x) \stackrel{\text{def}}{=} w_t(x) \left[ \sum_{i \in [n]} \frac{1}{(1 + g_t^{(i)}(x_t))g_t^{(i)}(x_t)} \right]^{-1}$$

denote a natural term that helps upper bound the rate of change of the central path. Furthermore, we let  $u^{(i)}(x)$  denote the unit vector corresponding to  $x - a^{(i)}$ , i.e.  $u^{(i)}(x) \stackrel{\text{def}}{=} (x - a^{(i)}) / \|x - a^{(i)}\|_2$  when  $\|x - a^{(i)}\|_2 \neq 0$  and  $u^{(i)}(x) \stackrel{\text{def}}{=} 0$  otherwise. Finally we let  $\mu_t(x) \stackrel{\text{def}}{=} \lambda_d(\nabla^2 f_t(x))$ , denote the minimum eigenvalue of  $\nabla^2 f_t(x)$ , and let  $v_t(x)$  denote a corresponding eigenvector. To simplify notation we often drop the  $(x)$  in these definitions when  $x = x_t$  and  $t$  is clear from context.



## ■ 7.3 Properties of the Central Path

Here provide various facts regarding the penalized objective function and the central path. While we use the lemmas in this section throughout the chapter, the main contribution of this section is Lemma 7.3.5 in Section 7.3.3. There we prove that with the exception of a single direction, the change in the central path is small over a constant multiplicative change in the path parameter. In addition, we show that our penalized objective function is stable under changes in a  $O(\frac{1}{t})$  Euclidean ball (Section 7.3.1), we bound the change in the Hessian over the central path (Section 7.3.2), and we relate  $f(x_t)$  to  $f(x_*)$  (Section 7.3.4).

### ■ 7.3.1 How Much Does the Hessian Change in General?

Here, we show that the Hessian of the penalized objective function is stable under changes in a  $O(\frac{1}{t})$  sized Euclidean ball. This shows that if we have a point which is close to a central path point in Euclidean norm, then we can use Newton method to find it.

*Lemma 7.3.1.* Suppose that  $\|x - y\|_2 \leq \frac{\varepsilon}{t}$  with  $\varepsilon \leq \frac{1}{20}$ . Then, we have

$$(1 - 6\varepsilon^{2/3})\nabla^2 f_t(x) \preceq \nabla^2 f_t(y) \preceq (1 + 6\varepsilon^{2/3})\nabla^2 f_t(x).$$

### ■ 7.3.2 How Much Does the Hessian Change Along the Path?

Here we bound how much the Hessian of the penalized objective function can change along the central path. First we provide the following lemma bound several aspects of the penalized objective function and proving that the weight,  $w_t$ , only changes by a small amount multiplicatively given small multiplicative changes in the path parameter,  $t$ .

*Lemma 7.3.2.* For all  $t \geq 0$  and  $i \in [n]$  the following hold

$$\left\| \frac{d}{dt} x_t \right\|_2 \leq \frac{1}{t^2} \bar{g}_t(x_t) \quad , \quad \left| \frac{d}{dt} g_t^{(i)}(x_t) \right| \leq \frac{1}{t} \left( g_t^{(i)}(x_t) + \bar{g}_t \right) \quad , \quad \text{and} \quad \left| \frac{d}{dt} w_t \right| \leq \frac{2}{t} w_t$$

Consequently, for all  $t' \geq t$  we have that  $(\frac{t}{t'})^2 w_t \leq w_{t'} \leq (\frac{t'}{t})^2 w_t$ .

Next we use this lemma to bound the change in the Hessian with respect to  $t$ .

*Lemma 7.3.3.* For all  $t \geq 0$  we have

$$-12 \cdot t \cdot w_t \mathbf{I} \preceq \frac{d}{dt} [\nabla^2 f_t(x_t)] \preceq 12 \cdot t \cdot w_t \mathbf{I} \tag{7.3}$$

and therefore for all  $\beta \in [0, \frac{1}{8}]$

$$\nabla^2 f(x_t) - 15\beta t^2 w_t \mathbf{I} \preceq \nabla^2 f(x_{t(1+\beta)}) \preceq \nabla^2 f(x_t) + 15\beta t^2 w_t \mathbf{I}. \tag{7.4}$$

### ■ 7.3.3 Where is the Next Optimal Point?

Here we prove our main result of this section. We prove that over a long step the central path moves very little in directions orthogonal to the smallest eigenvector of the Hessian. We begin by noting the Hessian is approximately a scaled identity minus a rank 1 matrix.

*Lemma 7.3.4.* For all  $t$  we have  $\frac{1}{2} [t^2 \cdot w_t \mathbf{I} - (t^2 \cdot w_t - \mu_t) v_t v_t^\top] \preceq \nabla^2 f_t(x_t) \preceq t^2 \cdot w_t \mathbf{I} - (t^2 \cdot w_t - \mu_t) v_t v_t^\top$

Using this and the lemmas of the previous section we bound the amount  $x_t$  can move in every direction far from  $v_t$ .

*Lemma 7.3.5* (The Central Path is Almost Straight). For all  $t \geq 0$ ,  $\beta \in [0, \frac{1}{600}]$ , and any unit vector  $y$  with  $|\langle y, v_t \rangle| \leq \frac{1}{t^2 \kappa}$  where  $\kappa = \max_{\delta \in [t, (1+\beta)t]} \frac{w_\delta}{\mu_\delta}$ , we have  $y^\top (x_{(1+\beta)t} - x_t) \leq \frac{6\beta}{t}$ .

### ■ 7.3.4 Where is the End?

In this section, we bound the quality of the central path with respect to the geometric median objective. In particular, we show that if we can solve the problem for some  $t = \frac{2n}{\varepsilon f(x_*)}$  then we obtain an  $(1 + \varepsilon)$ -approximate solution. As our algorithm ultimately starts from an initial  $t = 1/O(f(x_*))$  and increases  $t$  by a multiplicative constant in every iteration, this yields an  $O(\log \frac{n}{\varepsilon})$  iteration algorithm.

*Lemma 7.3.6.*  $f(x_t) - f(x_*) \leq \frac{2n}{t}$  for all  $t > 0$ .

## ■ 7.4 Overview of the Algorithm

---

### Algorithm 13: AccurateMedian( $\varepsilon$ )

---

**Input:** desired accuracy  $\varepsilon \in (0, 1)$

```

// Compute a 2-approximate geometric median and use it to center
Compute  $x^{(0)} := \frac{1}{n} \sum_{i \in [n]} a^{(i)}$  and  $\tilde{f}_* := f(x^{(0)})$  // Note  $\tilde{f}_* \leq 2f(x_*)$  by Lemma 7.5.6
Let  $t_i = \frac{1}{400\tilde{f}_*} (1 + \frac{1}{600})^{i-1}$ ,  $\tilde{\varepsilon}_* = \frac{1}{3}\varepsilon$ , and  $\tilde{t}_* = \frac{2n}{\tilde{\varepsilon}_* \tilde{f}_*}$ 
Let  $\varepsilon_v = \frac{1}{8} (\frac{\tilde{\varepsilon}_*}{7n})^2$  and let  $\varepsilon_c = (\frac{\varepsilon_v}{36})^{\frac{3}{2}}$ 
 $x^{(1)} = \text{LineSearch}(x^{(0)}, t_1, t_1, 0, \varepsilon_c)$ 

// Iteratively improve quality of approximation
Let  $k = \max_{i \in \mathbb{Z}} t_i \leq \tilde{t}_*$ 
for  $i \in [1, k]$  do
    // Compute  $\varepsilon_v$ -approximate minimum eigenvalue and eigenvector of  $\nabla^2 f_{t_i}(x^{(i)})$ 
     $(\lambda^{(i)}, u^{(i)}) = \text{ApproxMinEig}(x^{(i)}, t_i, \varepsilon_v)$ 

    // Line search to find  $x^{(i+1)}$  such that  $\|x^{(i+1)} - x_{t_{i+1}}\|_2 \leq \frac{\varepsilon_c}{t_{i+1}}$ 
     $x^{(i+1)} = \text{LineSearch}(x^{(i)}, t_i, t_{i+1}, u^{(i)}, \varepsilon_c)$  w.
end
Output:  $\varepsilon$ -approximate geometric median  $x^{(k)}$ 

```

---

Here we show how to use the structural results from the previous section to obtain a nearly linear time algorithm for computing the geometric median. Our algorithm follows a simple structure (See Algorithm 13). First we use simply average the  $a^{(i)}$  to compute a  $(1 + \frac{1}{n})$ -approximate median, denoted  $x^{(0)}$ . Then for a number of iterations we repeatedly move closer to  $x_t$  for some path parameter  $t$ , compute the minimum eigenvector of the Hessian, and line search in that direction to find an approximation to a point further along the central path. Ultimately, this yields a point  $x^{(k)}$  that is precise enough approximation to a point along the central path with large enough  $t$  that we can simply out  $x^{(k)}$  as our  $(1 + \varepsilon)$ -approximate geometric median.

We split the remainder of the algorithm specification and its analysis into several parts. First in Section 7.4.1 we show how to compute an approximate minimum eigenvector and eigenvalue of the Hessian of the penalized objective function. Then in Section 7.4.2 we show how to use this eigenvector to line search for the next central path point. Finally, in Section 7.4.3 we put these results together to obtain our nearly linear time algorithm. Throughout this section we will want an upper bound to

$f(x_*)$  and a slight lower bound on  $\varepsilon$ , the geometric median accuracy we are aiming for. We use an easily computed  $\tilde{f}_* \leq 2f(x_*)$  for the former and  $\tilde{\varepsilon}_* = \frac{1}{3}\varepsilon$  throughout the section.

### ■ 7.4.1 Eigenvector Computation and Hessian Approximation

Here we show how to compute the minimum eigenvector of  $\nabla^2 f_t(x)$  and thereby obtain a concise approximation to  $\nabla^2 f_t(x)$ . Our main algorithmic tool is the well known power method and the fact that it converges quickly on a matrix with a large eigenvalue gap. To improve our logarithmic terms we need a slightly non-standard analysis of the method and therefore we provide and analyze this method for completeness in Section 7.6.1. Using this tool we estimate the top eigenvector as follows.

---

#### Algorithm 14: ApproxMinEig( $x, t, \varepsilon$ )

---

**Input:** Point  $x \in \mathbb{R}^d$ , path parameter  $t$ , and target accuracy  $\varepsilon$ .

Let  $\mathbf{A} = \sum_{i \in [n]} \frac{t^4(x-a^{(i)})(x-a^{(i)})^\top}{(1+g_t^{(i)}(x))^2 g_t^{(i)}(x)}$

Let  $u := \text{PowerMethod}(\mathbf{A}, \Theta(\log(\frac{n}{\varepsilon})))$

Let  $\lambda = u^\top \nabla^2 f_t(x) u$

**Output:**  $(\lambda, u)$

---

*Lemma 7.4.1* (Computing Hessian Approximation). Let  $x \in \mathbb{R}^d$ ,  $t > 0$ , and  $\varepsilon \in (0, \frac{1}{4})$ . The algorithm `ApproxMinEig`( $x, t, \varepsilon$ ) outputs  $(\lambda, u)$  in  $O(nd \log \frac{n}{\varepsilon})$  time such that if  $\mu_t(x) \leq \frac{1}{4}t^2 w_t(x)$  then  $\langle v_t(x), u \rangle^2 \geq 1 - \varepsilon$  with high probability in  $n$ . Furthermore, if  $\varepsilon \leq \frac{\mu_t(x)}{8t^2 \cdot w_t(x)}$  then  $\frac{1}{4}\mathbf{Q} \preceq \nabla^2 f_t(x) \preceq 4\mathbf{Q}$  with high probability in  $n$  where  $\mathbf{Q} \stackrel{\text{def}}{=} t^2 \cdot w_t(x) - (t^2 \cdot w_t(x) - \lambda) uu^\top$ .

Furthermore, we show that the  $v^{(i)}$  computed by this algorithm is sufficiently close to the bad direction. Combining 7.4.1 with the structural results from the previous section and Lemma 7.9.4, a minor technical lemma regarding the transitivity of large inner products, we provide the following lemma.

*Lemma 7.4.2.* Let  $u = \text{ApproxMinEig}(x, t, \varepsilon_v)$  for  $\varepsilon_v < \frac{1}{8}$  and  $x$  such that  $\|x - x_t\|_2 \leq \frac{\varepsilon_c}{t}$  for  $\varepsilon_c \leq (\frac{\varepsilon_v}{36})^{\frac{3}{2}}$ .  $\mu_t \leq \frac{1}{4}t^2 \cdot w_t$ . For all unit vectors  $y \perp u$ , we have  $\langle y, v_t \rangle^2 \leq 8\varepsilon_v$ .

Note that this lemma assumes  $\mu_t$  is small. When  $\mu_t$  is large, we instead show that the next central path point is close to the current point and hence we do not need to compute the bad direction to center quickly.

*Lemma 7.4.3.* Suppose  $\mu_t \geq \frac{1}{4}t^2 \cdot w_t$  and let  $t' \in [t, (1 + \frac{1}{600})t]$  then  $\|x_{t'} - x_t\|_2 \leq \frac{1}{100t}$ .

### ■ 7.4.2 Line Searching

Here we show how to line search along the bad direction to find the next point on the central path. Unfortunately, it is not clear if you can binary search on the objective function directly. If we search over  $\alpha$  to minimize  $f_{t+1}(y^{(i)} + \alpha v^{(i)})$  directly it is unclear if we actually obtain a point close to  $x_{t+1}$ . It might be the case that even after minimizing  $\alpha$  we would be unable to move towards  $x_{t+1}$  efficiently.

To overcome this difficulty, we use the fact that over the region  $\|x - y\|_2 = O(\frac{1}{t})$  the Hessian changes by at most a constant and therefore we can minimize  $f_t(x)$  over this region extremely quickly. Therefore, we instead line search on the following function

$$g_{t,y,v}(\alpha) \stackrel{\text{def}}{=} \min_{\|x - (y + \alpha v)\|_2 \leq \frac{1}{49t}} f_t(x) \quad (7.5)$$

and use that we can evaluate  $g_{t,y,v}(\alpha)$  approximately by using an appropriate centering procedure. We can show (See Lemma 7.9.6) that  $g_{t,y,v}(\alpha)$  is convex and therefore we can minimize it efficiently just by doing an appropriate binary search. By finding the approximately minimizing  $\alpha$  and outputting the corresponding approximately minimizing  $x$ , we can obtain  $x^{(i+1)}$  that is close enough to  $x_{t_{i+1}}$ . For notational convenience, we simply write  $g(\alpha)$  if  $t, y, v$  is clear from the context.

First, we show how we can locally center and provide error analysis for that algorithm.

---

**Algorithm 15: LocalCenter** $(y, t, \varepsilon)$ 


---

**Input:** Point  $y \in \mathbb{R}^d$ , path parameter  $t > 0$ , target accuracy  $\varepsilon > 0$ .

Let  $(\lambda, v) := \text{ApproxMinEig}(x, t, \varepsilon)$ .

Let  $Q = t^2 \cdot w_t(y)I - (t^2 \cdot w_t(y) - \lambda)vv^\top$

Let  $x^{(0)} = y$

**for**  $i = 1, \dots, k = 64 \log \frac{1}{\varepsilon}$  **do**

    | Let  $x^{(i)} = \min_{\|x-y\|_2 \leq \frac{1}{49t}} f(x^{(i-1)}) + \langle \nabla f_t(x^{(i-1)}), x - x^{(i-1)} \rangle + 4\|x - x^{(i-1)}\|_2^2$ .

**end**

**Output:**  $x^{(k)}$

---

*Lemma 7.4.4.* Given some  $y \in \mathbb{R}^d$ ,  $t > 0$  and  $\varepsilon \in (0, \min\{\frac{1}{4}, \frac{\mu_t(x)}{8t^2 \cdot w_t(x)}\})$ . In  $O(nd \log(\frac{n}{\varepsilon}))$  time  $\text{LocalCenter}(y, t, \varepsilon)$  computes  $x^{(k)}$  such that with high probability in  $n/\varepsilon$ .

$$f_t(x^{(k)}) - \min_{\|x-y\|_2 \leq \frac{1}{49t}} f_t(x) \leq \varepsilon \left( f_t(y) - \min_{\|x-y\|_2 \leq \frac{1}{49t}} f_t(x) \right).$$

Using this local centering algorithm as well as a general result for minimizing one dimensional convex functions using a noisy oracle (See Section 7.9.3) we obtain our line search algorithm.

---

**Algorithm 16: LineSearch** $(y, t, t', u, \varepsilon)$ 


---

**Input:** Point  $y \in \mathbb{R}^d$ , current path parameter  $t$ , next path parameter  $t'$ , bad direction  $u$ , target accuracy  $\varepsilon$

Let  $\varepsilon_O = \frac{\varepsilon^2 \cdot \tilde{\varepsilon}_*^2}{300n}$ ,  $\ell = -6\tilde{f}_*$ ,  $u = 6\tilde{f}_*$ .

Define the oracle  $q: \mathbb{R} \rightarrow \mathbb{R}$  by  $q(\alpha) = f_{t'}(\text{LocalCenter}(y + \alpha u, t', \varepsilon_O))$

Let  $\alpha' = \text{OneDimMinimizer}(\ell, u, \varepsilon_O, q, tn)$

**Output:**  $x' = \text{LocalCenter}(y + \alpha u, t', \varepsilon_O)$

---

*Lemma 7.4.5.* Let  $\frac{1}{400f(x_*)} \leq t \leq t' \leq (1 + \frac{1}{600})t \leq \frac{2n}{\varepsilon \cdot \tilde{f}_*}$  and let  $u = \text{ApproxMinEig}(y, t, \varepsilon_v)$  for  $\varepsilon_v \leq \frac{1}{8}(\frac{\tilde{\varepsilon}_*}{7n})^2$  and  $y \in \mathbb{R}^d$  such that  $\|y - x_t\|_2 \leq \frac{1}{t}(\frac{\varepsilon_v}{36})^{\frac{3}{2}}$ . In  $O(nd \log^2(\frac{n}{\varepsilon \cdot \tilde{\varepsilon}_*}))$  time,  $\text{LineSearch}(y, t, t', u, \varepsilon)$  outputs  $x'$  such that  $\|x' - x_{t'}\|_2 \leq \frac{\varepsilon}{t'}$  with high probability in  $n$ .

We also provide the following lemma useful for finding the first center.

*Lemma 7.4.6.* Let  $\frac{1}{400f(x_*)} \leq t \leq t' \leq (1 + \frac{1}{600})t \leq \frac{2n}{\varepsilon_* \cdot \tilde{f}_*}$  and let  $x \in \mathbb{R}^d$  satisfy  $\|x - x_t\|_2 \leq \frac{1}{100t}$ . Then, in  $O(nd \log^2(\frac{n}{\varepsilon \cdot \tilde{\varepsilon}_*}))$  time,  $\text{LineSearch}(x, t, t', u, \varepsilon)$  output  $y$  such that  $\|y - x_t\|_2 \leq \frac{\varepsilon}{t}$  for any vector  $u \in \mathbb{R}^d$ .

### ■ 7.4.3 Putting It All Together

Combining the results of the previous sections to prove our main theorem.

**Theorem 7.4.1.** *In  $O(nd \log^3(\frac{n}{\varepsilon}))$  time, Algorithm 13 outputs an  $(1 + \varepsilon)$ -approximate geometric median with constant probability.*

## ■ 7.5 Analysis of the Central Path

Here we provide proofs of the claims in Section 7.3 as well as additional technical lemmas we use throughout the chapter.

### ■ 7.5.1 Basic Facts

Here we provide basic facts regarding the central path that we will use throughout our analysis. First we compute various derivatives of the penalized objective function.

*Lemma 7.5.1 (Path Derivatives).* We have

$$\nabla f_t(x) = \sum_{i \in [n]} \frac{t^2(x - a^{(i)})}{1 + g_t^{(i)}(x)}, \quad \nabla^2 f_t(x) = \sum_{i \in [n]} \frac{t^2}{1 + g_t^{(i)}(x)} \left( \mathbf{I} - \frac{t^2(x - a^{(i)})(x - a^{(i)})^\top}{g_t^{(i)}(x)(1 + g_t^{(i)}(x))} \right), \text{ and}$$

$$\frac{d}{dt} x_t = -(\nabla^2 f_t(x_t))^{-1} \sum_{i \in [n]} \frac{t(x_t - a^{(i)})}{(1 + g_t^{(i)}(x_t))g_t^{(i)}(x_t)}$$

*Proof of Lemma 7.5.1.* Direct calculation shows that

$$\begin{aligned} \nabla f_t^{(i)}(x) &= \frac{t^2(x - a^{(i)})}{\sqrt{1 + t^2\|x - a^{(i)}\|_2^2}} - \frac{1}{1 + \sqrt{1 + t^2\|x - a^{(i)}\|_2^2}} \left( \frac{t^2(x - a^{(i)})}{\sqrt{1 + t^2\|x - a^{(i)}\|_2^2}} \right) \\ &= \frac{t^2(x - a^{(i)})}{1 + \sqrt{1 + t^2\|x - a^{(i)}\|_2^2}} = \frac{t^2(x - a^{(i)})}{1 + g_t^{(i)}(x)} \end{aligned}$$

and

$$\begin{aligned} \nabla^2 f_t^{(i)}(x) &= \frac{t^2}{1 + \sqrt{1 + t^2\|x - a^{(i)}\|_2^2}} \mathbf{I} - \left( \frac{1}{1 + \sqrt{1 + t^2\|x - a^{(i)}\|_2^2}} \right)^2 \frac{t^4(x - a^{(i)})(x - a^{(i)})^\top}{\sqrt{1 + t^2\|x - a^{(i)}\|_2^2}} \\ &= \frac{t^2}{1 + g_t^{(i)}(x)} \left( \mathbf{I} - \frac{t^2(x - a^{(i)})(x - a^{(i)})^\top}{g_t^{(i)}(x)(1 + g_t^{(i)}(x))} \right) \end{aligned}$$

and

$$\begin{aligned} \left( \frac{d}{dt} \nabla f_t^{(i)} \right) (x) &= \frac{2t(x - a^{(i)})}{1 + \sqrt{1 + t^2\|x - a^{(i)}\|_2^2}} - \frac{t^2 \cdot (x - a^{(i)}) \cdot t\|x - a^{(i)}\|_2^2}{\left(1 + \sqrt{1 + t^2\|x - a^{(i)}\|_2^2}\right)^2 \sqrt{1 + t^2\|x - a^{(i)}\|_2^2}} \\ &= \frac{t \cdot (x - a^{(i)})}{1 + g_t^{(i)}(x)} \left( 2 - \frac{g_t^{(i)}(x)^2 - 1}{(1 + g_t^{(i)}(x))g_t^{(i)}(x)} \right) \\ &= \frac{t \cdot (x - a^{(i)})}{1 + g_t^{(i)}(x)} \left( \frac{2g_t^{(i)}(x) - (g_t^{(i)}(x) - 1)}{g_t^{(i)}(x)} \right) = \frac{t \cdot (x - a^{(i)})}{g_t^{(i)}(x)} \end{aligned}$$

Finally, by the optimality of  $x_t$  we have that  $\nabla f_t(x_t) = 0$ . Consequently,

$$\nabla^2 f_t(x_t) \frac{d}{dt} x_t + \left( \frac{d}{dt} \nabla f_t \right) (x_t) = 0.$$

and solving for  $\frac{d}{dt} x_t$  then yields

$$\begin{aligned} \frac{d}{dt} x_t &= - (\nabla^2 f_t(x_t))^{-1} \left( \left( \frac{d}{dt} \nabla f_t \right) (x_t) \right) \\ &= - (\nabla^2 f_t(x_t))^{-1} \left( \left( \frac{d}{dt} \nabla f_t \right) (x_t) - \frac{1}{t} \nabla f_t(x_t) \right) \\ &= - (\nabla^2 f_t(x_t))^{-1} \left( \sum_{i \in [n]} \left[ \frac{t}{g_t^{(i)}} - \frac{t}{1 + g_t^{(i)}} \right] (x_t - a^{(i)}) \right). \end{aligned}$$

□

Next, in we provide simple facts regarding the Hessian of the penalized objective function.

*Lemma 7.5.2.* For all  $t > 0$  and  $x \in \mathbb{R}^d$

$$\nabla^2 f_t(x) = \sum_{i \in [n]} \frac{t^2}{1 + g_t^{(i)}(x)} \left( \mathbf{I} - \left( 1 - \frac{1}{g_t^{(i)}(x)} \right) u^{(i)}(x) u^{(i)}(x)^\top \right)$$

and therefore

$$\sum_{i \in [n]} \frac{t^2}{(1 + g_t^{(i)}(x)) g_t^{(i)}(x)} \mathbf{I} \preceq \nabla^2 f_t(x) \preceq \sum_{i \in [n]} \frac{t^2}{1 + g_t^{(i)}(x)} \mathbf{I}$$

*Proof of Lemma 7.5.2.* We have that

$$\begin{aligned} \nabla^2 f_t(x) &= \sum_{i \in [n]} \frac{t^2}{1 + g_t^{(i)}(x)} \left( \mathbf{I} - \frac{t^2 (x - a^{(i)})(x - a^{(i)})^\top}{g_t^{(i)}(x)(1 + g_t^{(i)}(x))} \right) \\ &= \sum_{i \in [n]} \frac{t^2}{1 + g_t^{(i)}(x)} \left( \mathbf{I} - \frac{t^2 \|x - a^{(i)}\|_2^2}{(1 + g_t^{(i)}(x)) g_t^{(i)}(x)} u^{(i)}(x) u^{(i)}(x)^\top \right) \\ &= \sum_{i \in [n]} \frac{t^2}{1 + g_t^{(i)}(x)} \left( \mathbf{I} - u^{(i)}(x) u^{(i)}(x)^\top \right) \end{aligned}$$

Since

$$\frac{g_t^{(i)}(x)^2 - 1}{g_t^{(i)}(x)(1 + g_t^{(i)}(x))} = \frac{(g_t^{(i)}(x) + 1)(g_t^{(i)}(x) - 1)}{(1 + g_t^{(i)}(x)) g_t^{(i)}(x)} = 1 - \frac{1}{g_t^{(i)}(x)}$$

the result follows. □

## ■ 7.5.2 Stability of Hessian

Here we show that moving a point  $x \in \mathbb{R}^d$  in  $\ell_2$ , does not change the Hessian,  $\nabla^2 f_t(x)$ , too much spectrally. First we show that such changes do not change  $g_t^{(i)}(x)$  by too much (Lemma 7.5.3) and then we use this to prove the claim, i.e. we prove Lemma 7.3.1.

*Lemma 7.5.3* (Stability of  $g$ ). For all  $x, y \in \mathbb{R}^d$  and  $t > 0$ , we have

$$g_t^{(i)}(x) - t\|x - y\|_2 \leq g_t^{(i)}(y) \leq g_t^{(i)}(x) + t\|x - y\|_2$$

*Proof of Lemma 7.5.3.* Direct calculation reveals that

$$\begin{aligned} g_t^{(i)}(y)^2 &= 1 + t^2\|x - a^{(i)} + y - x\|_2^2 \\ &= 1 + t^2\|x - a^{(i)}\|_2^2 + 2t^2(x - a^{(i)})^\top(y - x) + t^2\|y - x\|_2^2 \\ &= g_t^{(i)}(x)^2 + 2t^2(x - a^{(i)})^\top(y - x) + t^2\|y - x\|_2^2. \end{aligned}$$

Consequently by Cauchy Schwarz

$$g_t^{(i)}(y)^2 \leq g_t^{(i)}(x)^2 + 2t^2\|x - a^{(i)}\|_2 \cdot \|y - x\|_2 + t^2\|y - x\|_2^2 \leq \left(g_t^{(i)}(x) + t\|y - x\|_2\right)^2$$

and

$$g_t^{(i)}(y)^2 \geq g_t^{(i)}(x)^2 - 2t^2\|x - a^{(i)}\|_2 \cdot \|y - x\|_2 + t^2\|y - x\|_2^2 \geq \left(g_t^{(i)}(x) - t\|y - x\|_2\right)^2.$$

□

*Lemma 7.3.1.* Suppose that  $\|x - y\|_2 \leq \frac{\varepsilon}{t}$  with  $\varepsilon \leq \frac{1}{20}$ . Then, we have

$$(1 - 6\varepsilon^{2/3})\nabla^2 f_t(x) \preceq \nabla^2 f_t(y) \preceq (1 + 6\varepsilon^{2/3})\nabla^2 f_t(x).$$

*Proof of Lemma 7.3.1.* Here we prove the following stronger statement, for all  $i \in [n]$

$$(1 - 6\varepsilon^{2/3})\nabla^2 f_t^{(i)}(x) \preceq \nabla^2 f_t^{(i)}(y) \preceq (1 + 6\varepsilon^{2/3})\nabla^2 f_t^{(i)}(x).$$

Without loss of generality let  $y - x = \alpha v + \beta u^{(i)}(x)$  for some  $v \perp u^{(i)}(x)$  with  $\|v\|_2 = 1$ . Since  $\|x - y\|_2^2 \leq \frac{\varepsilon^2}{t^2}$ , we know that  $\alpha^2, \beta^2 \leq \frac{\varepsilon^2}{t^2}$ . Also, let  $\bar{x} = x + \beta u^{(i)}(x)$ , so that clearly,  $u^{(i)}(x) = u^{(i)}(\bar{x})$ . Now some manipulation reveals that for all unit vectors  $z \in \mathbb{R}^d$  the following holds (so long as  $u^{(i)}(x) \neq 0$  and  $u^{(i)}(y) \neq 0$ )

$$\begin{aligned} & \left| \left[ u^{(i)}(x)^\top z \right]^2 - \left[ u^{(i)}(y)^\top z \right]^2 \right| \\ &= \left| \left[ u^{(i)}(\bar{x})^\top z \right]^2 - \left[ u^{(i)}(y)^\top z \right]^2 \right| \\ &= \left| \left[ \frac{(\bar{x} - a^{(i)})^\top z}{\|\bar{x} - a^{(i)}\|_2} \right]^2 - \left[ \frac{(y - a^{(i)})^\top z}{\|y - a^{(i)}\|_2} \right]^2 \right| \\ &\leq \left| \left[ \frac{(\bar{x} - a^{(i)})^\top z}{\|\bar{x} - a^{(i)}\|_2} \right]^2 - \left[ \frac{(\bar{x} - a^{(i)})^\top z}{\|y - a^{(i)}\|_2} \right]^2 \right| + \left| \left[ \frac{(\bar{x} - a^{(i)})^\top z}{\|y - a^{(i)}\|_2} \right]^2 - \left[ \frac{(y - a^{(i)})^\top z}{\|y - a^{(i)}\|_2} \right]^2 \right| \\ &\leq \left| 1 - \frac{\|\bar{x} - a^{(i)}\|_2^2}{\|y - a^{(i)}\|_2^2} \right| + \frac{\left| [(\bar{x} - a^{(i)} + \alpha v)^\top z]^2 - [(\bar{x} - a^{(i)})^\top z]^2 \right|}{\|y - a^{(i)}\|_2^2} \\ &= \frac{\alpha^2 + \left| 2 [(\bar{x} - a^{(i)})^\top z] \cdot [\alpha v^\top z] + [\alpha v^\top z]^2 \right|}{\|\bar{x} - a^{(i)}\|_2^2 + \alpha_i^2}. \end{aligned}$$

Where we used that  $y = \bar{x} + \alpha v$  and  $\|y - a^{(i)}\|_2^2 = \alpha^2 + \|\bar{x} - a^{(i)}\|_2^2$  (since  $v \perp (\bar{x} - a^{(i)})$ ). Now we know that  $\alpha^2 \leq \frac{\varepsilon^2}{t^2}$  and therefore, by Young's inequality and Cauchy Schwarz we have that for all  $\gamma > 0$

$$\begin{aligned}
\left| \left[ u^{(i)}(x)^\top z \right]^2 - \left[ u^{(i)}(y)z \right]^2 \right| &\leq \frac{2\alpha^2 + 2 \left| [(\bar{x} - a^{(i)})^\top z] \cdot [\alpha v^\top z] \right|}{\|\bar{x} - a^{(i)}\|_2^2 + \alpha^2} \\
&\leq \frac{2\alpha^2 + \gamma [(\bar{x} - a^{(i)})^\top z]^2 + \gamma^{-1}\alpha^2 [v^\top z]^2}{\|\bar{x} - a^{(i)}\|_2^2 + \alpha^2} \\
&\leq \frac{\alpha^2 \left( 2 + \gamma^{-1} (v^\top z)^2 \right)}{\|\bar{x} - a^{(i)}\|_2^2 + \alpha^2} + \gamma \left[ (u^{(i)}(x))^\top z \right]^2 \\
&\leq \frac{\varepsilon^2}{t^2 \|\bar{x} - a^{(i)}\|_2^2 + \varepsilon^2} \left( 2 + \frac{1}{\gamma} (v^\top z)^2 \right) + \gamma \left[ (u^{(i)}(x))^\top z \right]^2. \tag{7.6}
\end{aligned}$$

Note that

$$\begin{aligned}
t^2 \|\bar{x} - a^{(i)}\|_2^2 &= t^2 \left( \|x - a^{(i)}\|_2^2 + 2\beta(x - a^{(i)})^\top u^{(i)}(x) + \beta^2 \right) = \left( t\|x - a^{(i)}\|_2 + t\beta \right)^2 \\
&\geq \left( \max \left\{ t\|x - a^{(i)}\|_2 - \varepsilon, 0 \right\} \right)^2.
\end{aligned}$$

Consequently, if  $t\|x - a^{(i)}\|_2 \geq 2\varepsilon^{1/3} \sqrt{g_t^{(i)}(x)}$  then since  $\varepsilon \leq \frac{1}{20}$  we have that  $\|x - a^{(i)}\|_2 \geq 2\varepsilon$  and  $\|y - a^{(i)}\| \geq \varepsilon$ , justifying our assumption that  $u^{(i)}(x) \neq 0$  and  $u^{(i)}(y) \neq 0$ . Furthermore, this implies that  $t^2 \|\bar{x} - a^{(i)}\|_2^2 \geq 2\varepsilon^{2/3} g_t^{(i)}(x)$  and therefore letting  $\gamma = \frac{\varepsilon^{2/3}}{g_t^{(i)}(x)}$  yields

$$\begin{aligned}
\left| \left[ u_t^{(i)}(x)^\top z \right]^2 - \left[ u_t^{(i)}(y)z \right]^2 \right| &\leq \frac{\varepsilon^{4/3}}{2g_t^{(i)}(x)} \left( 2 + \frac{g_t^{(i)}(x)}{\varepsilon^{2/3}} [v^\top z]^2 \right) + \frac{\varepsilon^{2/3}}{g_t^{(i)}(x)} \left[ (u^{(i)}(x))^\top z \right]^2 \\
&\leq 2\varepsilon^{2/3} [v^\top z]^2 + \frac{\varepsilon^{2/3}}{g_t^{(i)}(x)} \left[ (u^{(i)}(x))^\top z \right]^2 \\
&\leq 2\varepsilon^{2/3} \left( \frac{1 + g_t^{(i)}(x)}{t^2} \right) \|z\|_{\nabla^2 f_t^{(i)}(x)}^2
\end{aligned}$$

and therefore if we let

$$\mathbf{H} \stackrel{\text{def}}{=} \frac{t^2}{1 + g_t^{(i)}(x)} \left( \mathbf{I} - \left( 1 - \frac{1}{g_t^{(i)}(x)} \right) u^{(i)}(y)(u^{(i)}(y))^\top \right),$$

we see that for unit vectors  $z$ ,

$$\left| z^\top \left( \mathbf{H} - \nabla^2 f_t^{(i)}(x) \right) z \right| \leq 2\varepsilon^{2/3} \|z\|_{\nabla^2 f_t^{(i)}(x)}^2$$

Otherwise,  $t\|x - a^{(i)}\|_2 < 2\varepsilon^{1/3} \sqrt{g_t^{(i)}(x)}$  and therefore

$$g_t^{(i)}(x)^2 = 1 + t^2 \|x - a^{(i)}\|_2^2 \leq 1 + 4\varepsilon^{2/3} g_t^{(i)}(x)$$

Therefore, we have

$$g_t^{(i)}(x) \leq \frac{4\varepsilon^{2/3} + \sqrt{(4\varepsilon^{2/3})^2 + 4}}{2} \leq 1 + 4\varepsilon^{2/3}.$$



Therefore independent of (7.6) and the assumption that  $u^{(i)}(x) \neq 0$  and  $u^{(i)}(y) \neq 0$  we have

$$\frac{1}{1 + 4\varepsilon^{2/3}} \mathbf{H} \preceq \frac{t^2}{(1 + g_t^{(i)}(x))g_t^{(i)}(x)} \mathbf{I} \preceq \nabla^2 f_t^{(i)}(x) \preceq \frac{t^2}{(1 + g_t^{(i)}(x))} \mathbf{I} \preceq (1 + 4\varepsilon^{2/3}) \mathbf{H}.$$

In either case, we have that

$$\left| z^\top \left( \mathbf{H} - \nabla^2 f_t^{(i)}(x) \right) z \right| \leq 4\varepsilon^{2/3} \|z\|_{\nabla^2 f_t^{(i)}(x)}^2.$$

Now, we note that  $\|x - y\|_2 \leq \frac{\varepsilon}{t} \leq \varepsilon \cdot \frac{g_t^{(i)}(x)}{t}$ . Therefore, by Lemma 7.5.3 we have that

$$(1 - \varepsilon)g_t^{(i)}(x) \leq g_t^{(i)}(y) \leq (1 + \varepsilon)g_t^{(i)}(x)$$

Therefore, we have

$$\frac{1 - 4\varepsilon^{2/3}}{(1 + \varepsilon)^2} \nabla^2 f_t^{(i)}(x) \preceq \frac{1}{(1 + \varepsilon)^2} \mathbf{H} \preceq \nabla^2 f_t^{(i)}(y) \preceq \frac{1}{(1 - \varepsilon)^2} \mathbf{H} \preceq \frac{1 + 4\varepsilon^{2/3}}{(1 - \varepsilon)^2} \nabla^2 f_t^{(i)}(x)$$

Since  $\varepsilon < \frac{1}{20}$ , the result follows.  $\square$

Consequently, so long as we have a point within a  $O(\frac{1}{t})$  sized Euclidean ball of some  $x_t$  Newton's method (or an appropriately transformed first order method) within the ball will converge quickly.

### ■ 7.5.3 How Much Does the Hessian Change Along the Path?

*Lemma 7.3.2.* For all  $t \geq 0$  and  $i \in [n]$  the following hold

$$\left\| \frac{d}{dt} x_t \right\|_2 \leq \frac{1}{t^2} \bar{g}_t(x_t) \quad , \quad \left| \frac{d}{dt} g_t^{(i)}(x_t) \right| \leq \frac{1}{t} \left( g_t^{(i)}(x_t) + \bar{g}_t \right) \quad , \quad \text{and} \quad \left| \frac{d}{dt} w_t \right| \leq \frac{2}{t} w_t$$

Consequently, for all  $t' \geq t$  we have that  $(\frac{t}{t'})^2 w_t \leq w_{t'} \leq (\frac{t'}{t})^2 w_t$ .

*Proof of Lemma 7.3.2.* From Lemma 7.5.1 we know that

$$\frac{d}{dt} x_t = -(\nabla^2 f_t(x_t))^{-1} \sum_{i \in [n]} \frac{t(x_t - a^{(i)})}{(1 + g_t^{(i)}(x_t))g_t^{(i)}(x_t)}$$

and by Lemma 7.5.2 we know that

$$\nabla^2 f_t(x_t) \succeq \sum_{i \in [n]} \frac{t^2}{(1 + g_t^{(i)}(x_t))g_t^{(i)}(x_t)} \mathbf{I} = \frac{t^2}{\bar{g}_t(x_t)} \sum_{i \in [n]} \frac{1}{1 + g_t^{(i)}(x_t)} \mathbf{I}.$$

Using this fact and the fact that  $t\|x_t - a^{(i)}\|_2 \leq g_t^{(i)}$  we have

$$\begin{aligned} \left\| \frac{d}{dt} x_t \right\|_2 &= \left\| -(\nabla^2 f_t(x_t))^{-1} \frac{d}{dt} \nabla f_t(x_t) \right\|_2 \\ &\leq \left( \frac{t^2}{\bar{g}_t(x_t)} \sum_{i \in [n]} \frac{1}{1 + g_t^{(i)}(x_t)} \right)^{-1} \sum_{i \in [n]} \left\| \frac{t(x_t - a^{(i)})}{g_t^{(i)}(x_t)(1 + g_t^{(i)}(x_t))} \right\|_2 \leq \frac{\bar{g}_t(x_t)}{t^2}. \end{aligned}$$

Next, we have

$$\begin{aligned} \frac{d}{dt}g_t^{(i)}(x_t) &= \frac{d}{dt} \left(1 + t^2\|x_t - a^{(i)}\|_2^2\right)^{\frac{1}{2}} \\ &= \frac{1}{2} \cdot g_t^{(i)}(x_t)^{-1} \left(2t\|x_t - a^{(i)}\|_2^2 + 2t^2(x_t - a^{(i)})^\top \frac{d}{dt}x_t\right) \end{aligned}$$

which by Cauchy Schwarz and that  $t\|x_t - a^{(i)}\|_2 \leq g_t^{(i)}(x_t)$  yields the second equation. Furthermore,

$$\begin{aligned} \left|\frac{d}{dt}w_t\right| &= \left|\frac{d}{dt} \sum_{i \in [n]} \frac{1}{1 + g_t^{(i)}(x_t)}\right| \leq \sum_{i \in [n]} \left|\frac{d}{dt} \frac{1}{1 + g_t^{(i)}(x_t)}\right| = \sum_{i \in [n]} \left|\frac{1}{(1 + g_t^{(i)}(x_t))^2} \frac{d}{dt}g_t^{(i)}(x_t)\right| \\ &\leq \frac{1}{t} \sum_{i \in [n]} \frac{g_t^{(i)}(x_t) + \bar{g}_t}{(1 + g_t^{(i)}(x_t))g_t^{(i)}(x_t)} \leq 2\frac{w_t}{t} \end{aligned}$$

which yields the third equation.

Finally, using our earlier results and Jensen's inequality yields that

$$|\ln w_{t'} - \ln w_t| = \left|\int_t^{t'} \frac{\frac{d}{d\alpha}w_\alpha}{w_\alpha} d\alpha\right| \leq \int_t^{t'} \frac{(2\frac{w_\alpha}{\alpha})}{w_\alpha} d\alpha = 2 \int_t^{t'} \frac{1}{\alpha} d\alpha = \ln \left(\frac{t'}{t}\right)^2.$$

Exponentiating the above inequality yields the final inequality.  $\square$

*Lemma 7.3.3.* For all  $t \geq 0$  we have

$$-12 \cdot t \cdot w_t \mathbf{I} \preceq \frac{d}{dt} [\nabla^2 f_t(x_t)] \preceq 12 \cdot t \cdot w_t \mathbf{I} \quad (7.3)$$

and therefore for all  $\beta \in [0, \frac{1}{8}]$

$$\nabla^2 f(x_t) - 15\beta t^2 w_t \mathbf{I} \preceq \nabla^2 f(x_{t(1+\beta)}) \preceq \nabla^2 f(x_t) + 15\beta t^2 w_t \mathbf{I}. \quad (7.4)$$

*Proof of Lemma 7.3.3.* Let

$$\mathbf{A}_t^{(i)} \stackrel{\text{def}}{=} \frac{t^2(x_t - a^{(i)})(x_t - a^{(i)})^\top}{(1 + g_t^{(i)})g_t^{(i)}}$$

and recall that  $\nabla^2 f_t(x_t) = \sum_{i \in [n]} \frac{t^2}{1 + g_t^{(i)}} (\mathbf{I} - \mathbf{A}_t^{(i)})$ . Consequently

$$\begin{aligned} \frac{d}{dt} \nabla^2 f_t(x_t) &= \frac{d}{dt} \left( \sum_{i \in [n]} \frac{t^2}{1 + g_t^{(i)}} (\mathbf{I} - \mathbf{A}_t^{(i)}) \right) \\ &= 2t \left(\frac{1}{t^2}\right) \nabla^2 f_t(x_t) + t^2 \sum_{i \in [n]} \frac{-\frac{d}{dt}g_t^{(i)}}{(1 + g_t^{(i)})^2} (\mathbf{I} - \mathbf{A}_t^{(i)}) - \sum_{i \in [n]} \frac{t^2}{1 + g_t^{(i)}} \frac{d}{dt} \mathbf{A}_t^{(i)} \end{aligned}$$

Now, since  $\mathbf{0} \preceq \mathbf{A}_t^{(i)} \preceq \mathbf{I}$  we have  $0 \preceq \nabla^2 f_t(x_t) \preceq t^2 w_t \mathbf{I}$ . For all unit vectors  $v$ , using Lemma 7.3.2

yields

$$\begin{aligned} \left| v^\top \left( \frac{d}{dt} \nabla^2 f_t(x_t) \right) v \right| &\leq 2t \cdot w_t \cdot \|v\|_2^2 + t^2 \sum_{i \in [n]} \frac{\left| \frac{d}{dt} g_t^{(i)} \right|}{(1 + g_t^{(i)})^2} \|v\|_2^2 + \sum_{i \in [n]} \frac{t^2}{1 + g_t^{(i)}} \left| v^\top \left( \frac{d}{dt} \mathbf{A}_t^{(i)} \right) v \right| \\ &\leq 4t \cdot w_t + \sum_{i \in [n]} \frac{t^2}{1 + g_t^{(i)}} \left| v^\top \left( \frac{d}{dt} \mathbf{A}_t^{(i)} \right) v \right|. \end{aligned}$$

Next

$$\begin{aligned} \frac{d}{dt} \mathbf{A}_t^{(i)} &= 2t \left( \frac{1}{t^2} \right) \mathbf{A}_t^{(i)} - \left( \frac{t}{(1 + g_t^{(i)}) g_t^{(i)}} \right)^2 \left[ (1 + g_t^{(i)}) \frac{d}{dt} g_t^{(i)} + g_t^{(i)} \frac{d}{dt} g_t^{(i)} \right] (x_t - a^{(i)}) (x_t - a^{(i)})^\top \\ &\quad + \frac{t^2}{(1 + g_t^{(i)}) g_t^{(i)}} \left[ (x_t - a^{(i)}) \left( \frac{d}{dt} x_t \right)^\top + \left( \frac{d}{dt} x_t \right) (x_t - a^{(i)})^\top \right], \end{aligned}$$

and therefore by Lemma 7.3.2 and the that  $t \|x_t - a^{(i)}\|_2 \leq g_t^{(i)}$  we have

$$\begin{aligned} \left| v^\top \left( \frac{d}{dt} \mathbf{A}_t^{(i)} \right) v \right| &\leq \left( \frac{2}{t} + \frac{2t^2 \left| \frac{d}{dt} g_t^{(i)} \right|}{(1 + g_t^{(i)}) (g_t^{(i)})^2} \|x_t - a^{(i)}\|_2^2 + \frac{2t^2 \|x_t - a^{(i)}\|_2 \left\| \frac{d}{dt} x_t \right\|_2}{(1 + g_t^{(i)}) g_t^{(i)}} \right) \|v\|_2^2 \\ &\leq \frac{2}{t} + \frac{2}{t} \cdot \frac{g_t^{(i)} + \bar{g}_t}{1 + g_t^{(i)}} + \frac{2}{t} \cdot \frac{\bar{g}_t}{1 + g_t^{(i)}} \leq \frac{4}{t} + \frac{4}{t} \frac{\bar{g}_t}{1 + g_t^{(i)}}. \end{aligned}$$

Consequently, we have

$$\left| v^\top \left( \frac{d}{dt} \nabla^2 f_t(x_t) \right) v \right| \leq 8t \cdot w_t + 4t \sum_{i \in [n]} \frac{\bar{g}_t}{(1 + g_t^{(i)})^2} \leq 12t \cdot w_t$$

which completes the proof of (7.3). To prove (7.4), let  $v$  be any unit vector and note that

$$\begin{aligned} \left| v^\top (\nabla^2 f_{t(1+\beta)}(x) - \nabla^2 f_t(x)) v \right| &= \left| \int_t^{t(1+\beta)} v^\top \frac{d}{dt} [\nabla^2 f_\alpha(x_\alpha)] v \cdot d\alpha \right| \leq 12 \int_t^{t(1+\beta)} \alpha \cdot w_\alpha d\alpha \\ &\leq 12 \int_t^{t(1+\beta)} \alpha \left( \frac{\alpha}{t} \right)^2 w_t d\alpha \leq \frac{12}{t^2} \left( \frac{1}{4} [t(1+\beta)]^4 - \frac{1}{4} t^4 \right) w_t \\ &= 3t^2 [(1+\beta)^4 - 1] w_t \leq 15t^2 \beta w_t \end{aligned}$$

where we used Lemma 7.3.3 and  $0 \leq \beta \leq \frac{1}{8}$  at the last line.  $\square$

#### ■ 7.5.4 Where is the next Optimal Point?

*Lemma 7.5.4.* For all  $t$  we have  $\frac{1}{2} [t^2 \cdot w_t \mathbf{I} - (t^2 \cdot w_t - \mu_t) v_t v_t^\top] \preceq \nabla^2 f_t(x_t) \preceq t^2 \cdot w_t \mathbf{I} - (t^2 \cdot w_t - \mu_t) v_t v_t^\top$

*Proof of Lemma 7.5.4.* This follows immediately from Lemma 7.5.2, regarding the hessian of the penalized objective function, and Lemma 7.9.1, regarding the sum of PSD matrices expressed as the identity matrix minus a rank 1 matrix.  $\square$

*Lemma 7.3.5 (The Central Path is Almost Straight).* For all  $t \geq 0$ ,  $\beta \in [0, \frac{1}{600}]$ , and any unit vector  $y$  with  $|\langle y, v_t \rangle| \leq \frac{1}{t^{2-\kappa}}$  where  $\kappa = \max_{\delta \in [t, (1+\beta)t]} \frac{w_\delta}{\mu_\delta}$ , we have  $y^\top (x_{(1+\beta)t} - x_t) \leq \frac{6\beta}{t}$ .

*Proof of Lemma 7.3.5.* Clearly

$$\begin{aligned} y^\top(x_{(1+\beta)t} - x_t) &= \int_t^{(1+\beta)t} y^\top \frac{d}{d\alpha} x_\alpha d\alpha \leq \int_\beta^{(1+\beta)t} \left| y^\top \frac{d}{d\alpha} x_\alpha \right| d\alpha \\ &\leq \int_t^{(1+\beta)t} \left| y^\top (\nabla^2 f_\alpha(x_\alpha))^{-1} \sum_{i \in [n]} \frac{\alpha}{(1 + g_\alpha^{(i)}) g_\alpha^{(i)}} (x_\alpha - a^{(i)}) \right| d\alpha \\ &\leq \int_t^{(1+\beta)t} \|(\nabla^2 f_\alpha(x_\alpha))^{-1} y\|_2 \cdot \left\| \sum_{i \in [n]} \frac{\alpha}{(1 + g_\alpha^{(i)}) g_\alpha^{(i)}} (x_\alpha - a^{(i)}) \right\|_2 d\alpha \end{aligned}$$

Now since clearly  $\alpha \|x_\alpha - a^{(i)}\|_2 \leq g_\alpha^{(i)}$ , invoking Lemma 7.3.2 yields that

$$\left\| \sum_{i \in [n]} \frac{\alpha (x_\alpha - a^{(i)})}{(1 + g_\alpha^{(i)}) g_\alpha^{(i)}} \right\|_2 \leq \sum_{i \in [n]} \frac{1}{1 + g_\alpha^{(i)}} = w_\alpha \leq \left(\frac{\alpha}{t}\right)^2 w_t.$$

Now by invoking Lemma 7.3.3 and the Lemma 7.5.4, we have that

$$\nabla^2 f_\alpha(x_\alpha) \succeq \nabla^2 f_t(x_t) - 15\beta t^2 w_t \mathbf{I} \succeq \frac{1}{2} \left[ t^2 \cdot w_t \mathbf{I} - (t^2 \cdot w_t - \mu_t) v_t v_t^\top \right] - 15\beta t^2 w_t \mathbf{I}.$$

For notational convenience let  $\mathbf{H}_t \stackrel{\text{def}}{=} \nabla^2 f_t(x_t)$  for all  $t > 0$ . Then Lemma 7.3.3 shows that  $\mathbf{H}_\alpha = \mathbf{H}_t + \Delta_\alpha$  where  $\|\Delta_\alpha\|_2 \leq 15\beta t^2 w_t$ . Using  $\mathbf{H}_\alpha^2 = \mathbf{H}_t^2 + \Delta_\alpha \mathbf{H}_t + \mathbf{H}_t \Delta_\alpha + \Delta_\alpha^2$ , we have

$$\begin{aligned} \|\mathbf{H}_\alpha^2 - \mathbf{H}_t^2\|_2 &\leq \|\Delta_\alpha \mathbf{H}_t\|_2 + \|\mathbf{H}_t \Delta_\alpha\|_2 + \|\Delta_\alpha^2\|_2 \\ &\leq 2\|\Delta_\alpha\|_2 \|\mathbf{H}_t\|_2 + \|\Delta_\alpha\|_2^2 \leq 40\beta t^4 w_t^2. \end{aligned}$$

Let  $S$  be the subspace orthogonal to  $v_t$ . Then, Lemma 7.5.4 shows that  $\mathbf{H}_t \succeq \frac{1}{2} t^2 w_t \mathbf{I}$  on  $S$  and hence  $\mathbf{H}_t^2 \succeq \frac{1}{4} t^4 w_t^2 \mathbf{I}$  on  $S$ .<sup>4</sup> Since  $\|\mathbf{H}_\alpha^2 - \mathbf{H}_t^2\|_2 \leq 40\beta t^4 w_t^2$ , we have that

$$\mathbf{H}_\alpha^2 \succeq \frac{1}{4} t^4 w_t^2 - 40\beta t^4 w_t^2 \mathbf{I} \text{ on } S$$

and hence

$$\mathbf{H}_\alpha^{-2} \preceq \left( \frac{1}{4} t^4 w_t^2 - 40\beta t^4 w_t^2 \right)^{-1} \mathbf{I} \text{ on } S.$$

Therefore, for any  $z \in S$ , we have

$$\left\| (\nabla^2 f_\alpha(x_\alpha))^{-1} z \right\|_2 = \|\mathbf{H}_\alpha^{-1} z\|_2 \leq \frac{\|z\|_2}{\sqrt{\frac{1}{4} t^4 w_t^2 - 40\beta t^4 w_t^2}}.$$

Now, we split  $y = z + \langle y, v_t \rangle v_t$  where  $z \in S$ . Then, we have that

$$\begin{aligned} \left\| (\nabla^2 f_\alpha(x_\alpha))^{-1} y \right\|_2 &\leq \left\| (\nabla^2 f_\alpha(x_\alpha))^{-1} z \right\|_2 + |\langle y, v_t \rangle| \left\| (\nabla^2 f_\alpha(x_\alpha))^{-1} v_t \right\|_2 \\ &\leq \frac{1}{\sqrt{\frac{1}{4} t^4 w_t^2 - 40\beta t^4 w_t^2}} + \frac{1}{t^2 \cdot \kappa} \left\| (\nabla^2 f_\alpha(x_\alpha))^{-1} v_t \right\|_2. \end{aligned}$$

Note that, we also know that  $\lambda_{\min}(\nabla^2 f_\alpha(x_\alpha)) \geq \mu_\alpha$  and hence  $\lambda_{\max}(\nabla^2 f_\alpha(x_\alpha)^{-2}) \leq \mu_\alpha^{-2}$ . Therefore,

<sup>4</sup>By  $\mathbf{A} \preceq \mathbf{B}$  on  $S$  we mean that for all  $x \in S$  we have  $x^\top \mathbf{A} x \leq x^\top \mathbf{B} x$ . The meaning of  $\mathbf{A} \succeq \mathbf{B}$  on  $S$  is analogous.

we have

$$\left\| (\nabla^2 f_\alpha(x_\alpha))^{-1} y \right\|_2 \leq \frac{1}{t^2 w_t \sqrt{\frac{1}{4} - 40\beta}} + \frac{1}{t^2} \frac{\mu_\alpha}{w_\alpha} \frac{1}{\mu_\alpha} \leq \frac{1}{t^2 w_t} \left( 2 + \frac{1}{\sqrt{\frac{1}{4} - 40\beta}} \right) \leq \frac{5}{t^2 w_t}.$$

Combining these and using that  $\beta \in [0, 1/600]$  yields that

$$\begin{aligned} y^\top (x_{(1+\beta)t} - x_t) &\leq \int_t^{(1+\beta)t} \frac{5}{t^2 w_t} \left( \frac{\alpha}{t} \right)^2 w_t d\alpha \leq \frac{5}{t^4} \left( \frac{1}{3} (1+\beta)^3 t^3 - \frac{1}{3} t^3 \right) \\ &\leq \frac{5}{3t} [(1+\beta)^3 - 1] \leq \frac{6\beta}{t}. \end{aligned}$$

□

### ■ 7.5.5 Where is the End?

*Lemma 7.3.6.*  $f(x_t) - f(x_*) \leq \frac{2n}{t}$  for all  $t > 0$ .

*Proof of Lemma 7.3.6.* Clearly,  $\nabla f_t(x_t) = 0$  by definition of  $x_t$ . Consequently  $\frac{1}{t} \nabla f_t(x_t)^\top (x_t - x_*) = 0$  and using Lemma 7.5.1 to give the formula for  $\nabla f_t(x_t)$  yields

$$0 = \sum_{i \in [n]} \frac{t(x_t - a^{(i)})^\top (x_t - x_*)}{1 + g_t^{(i)}(x_t)} = \sum_{i \in [n]} \frac{t \|x_t - a^{(i)}\|_2^2 + t(x_t - a^{(i)})^\top (a^{(i)} - x_*)}{1 + g_t^{(i)}(x_t)}.$$

Therefore, by Cauchy Schwarz and the fact that  $t \|x_t - a^{(i)}\|_2 \leq g_t^{(i)}(x_t) \leq 1 + g_t^{(i)}$

$$\sum_{i \in [n]} \frac{t \|x_t - a^{(i)}\|_2^2}{1 + g_t^{(i)}(x_t)} \leq \sum_{i \in [n]} \frac{t \|x_t - a^{(i)}\|_2 \|a^{(i)} - x_*\|_2}{1 + g_t^{(i)}(x_t)} \leq f(x_*).$$

Furthermore, since  $1 + g_t^{(i)}(x_t) \leq 2 + t \|x_t - a^{(i)}\|_2$  we have

$$\sum_{i \in [n]} \frac{t \|x_t - a^{(i)}\|_2^2}{1 + g_t^{(i)}(x_t)} \geq \sum_{i \in [n]} \|x_t - a^{(i)}\|_2 - \sum_{i \in [n]} \frac{2 \|x_t - a^{(i)}\|_2}{1 + g_t^{(i)}(x_t)} \geq f(x_t) - \frac{2n}{t}.$$

Combining yields the result. □

### ■ 7.5.6 Simple Lemmas

Here we provide various small technical Lemmas that we will use to bound the accuracy with which we need to carry out various operations in our algorithm. Here we use some notation from Section 7.4 to simplify our bounds and make them more readily applied.

*Lemma 7.5.5.* For any  $x$ , we have that  $\|x - x_t\|_2 \leq f(x)$ .

*Proof of Lemma 7.5.5.* Since  $\sum_{i \in [n]} \|x - a^{(i)}\|_2 = f(x)$ , we have that  $\|x - a^{(i)}\|_2 \leq f(x)$  for all  $i \in [n]$ . Since  $\nabla f(x_t) = 0$  by Lemma 7.5.1 we see that  $x_t$  is a convex combination of the  $a^{(i)}$  and therefore  $\|x - x_t\|_2 \leq f(x)$  by convexity. □

*Lemma 7.5.6.*  $x^{(0)} = \frac{1}{n} \sum_{i \in [n]} a^{(i)}$  is a 2-approximate geometric median, i.e.  $\tilde{f}_* \leq 2 \cdot f(x_*)$ .

*Proof.* For all  $x \in \mathbb{R}^d$  we have

$$\|x^{(0)} - x\|_2 = \left\| \frac{1}{n} \sum_{i \in [n]} a^{(i)} - \frac{1}{n} \sum_{i \in [n]} x \right\|_2 \leq \frac{1}{n} \sum_{i \in [n]} \|a^{(i)} - x\|_2 \leq \frac{f(x)}{n}.$$

Consequently,

$$f(x^{(0)}) \leq \sum_{i \in [n]} \|x^{(0)} - a^{(i)}\|_2 \leq \sum_{i \in [n]} \left( \|x^{(0)} - x_*\|_2 + \|x_* - a^{(i)}\|_2 \right) \leq 2 \cdot f(x_*)$$

□

*Lemma 7.5.7.* For all  $t \geq 0$ , we have

$$1 \leq \frac{t^2 \cdot w_t(x)}{\mu_t(x)} \leq \bar{g}_t(x) \leq \max_{i \in [n]} g_t^{(i)}(x) \leq 1 + t \cdot f(x).$$

In particular, if  $\|x - x_t\|_2 \leq \frac{1}{t} + \alpha$  and  $t \leq \frac{2n}{\tilde{\varepsilon}_* \cdot f_*}$  then  $g_t^{(i)}(x) \leq 6n\varepsilon_*^{-1} + t \cdot \alpha$ .

*Proof of Lemma 7.5.7.* The first claim  $1 \leq \frac{t^2 \cdot w_t(x)}{\mu_t(x)} \leq \bar{g}_t(x)$ , follows from  $\mu_t(x) \geq \sum_{i \in [n]} \frac{t^2}{g_t^{(i)}(x)(1+g_t^{(i)}(x))}$  and the fact that the largest eigenvalue of  $\nabla^2 f_t(x)$  is at most  $t^2 \cdot w_t(x)$ . The second follows from the fact that  $\bar{g}_t(x)$  is a weighted harmonic mean of  $g_t^{(i)}(x)$  and therefore

$$\bar{g}_t(x) \leq \max_{i \in [n]} g_t^{(i)}(x) \leq 1 + t \cdot \max_{i \in [n]} \|x - a^{(i)}\|_2 \leq 1 + t \cdot f(x).$$

The final inequality comes from the fact that  $\|x - a^{(i)}\|_2 \leq \|x - x_t\|_2 + \|x_t - a^{(i)}\|_2 \leq \|x - x_t\|_2 + f(x_t)$  and the fact that  $f(x_t) \leq f(x_*) + \frac{2n}{t}$  by Lemma 7.3.6. □

*Lemma 7.5.8.* For all  $x \in \mathbb{R}^d$  and  $t > 0$  such that  $\|x - x_t\|_2 \leq \frac{1}{t} + \alpha$  and  $t \leq \frac{2n}{\tilde{\varepsilon}_* \cdot f_*}$  we have

$$n \left( \frac{t}{7n\tilde{\varepsilon}_*^{-1} + t \cdot \alpha} \right)^2 \|x - x_t\|_2^2 \leq f_t(x) - f_t(x_t) \leq \frac{nt^2}{2} \|x - x_t\|_2^2$$

*Proof of Lemma 7.5.8.* For the first inequality, note that  $\nabla^2 f_t(x) \preceq \sum_{i \in [n]} \frac{t^2}{1+g_t^{(i)}(x)} \mathbf{I} \preceq n \cdot t^2 \mathbf{I}$ . Consequently, if we let  $n \cdot t^2 \mathbf{I} = \mathbf{H}$  in Lemma 7.9.5, we have that

$$f_t(x) - f_t(x_t) \leq \frac{1}{2} \|x - x_t\|_{\mathbf{H}}^2 \leq \frac{nt^2}{2} \|x - x_t\|_2^2.$$

For the second inequality, note that Lemma 7.5.2 and Lemma 7.5.7 yields that

$$\nabla^2 f_t(x) \succeq \sum_{i \in [n]} \frac{t^2}{(1+g_t^{(i)}(x))g_t^{(i)}(x)} \mathbf{I} \succeq n \left( \frac{t}{7n\tilde{\varepsilon}_*^{-1} + t \cdot \alpha} \right)^2 \mathbf{I}.$$

Consequently, if we let  $S$  denote the set of all  $x \in \mathbb{R}^n$  with  $\|x - x_t\|_2 \leq \frac{1}{t} + \alpha$  then we see that  $S$  is convex and contains the minimizer of  $f_t$  and therefore by applying 7.9.5 again yields the lower bound. □

## ■ 7.6 Analysis of the Algorithm

Here we provide proofs, algorithms, and technical lemmas from Section 7.4.

### ■ 7.6.1 Eigenvector Computation and Hessian Approximation

Below we prove that the power method can be used to compute an  $\varepsilon$ -approximate top eigenvector of a symmetric PSD matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  with a non-zero eigenvalue gap  $g = \frac{\lambda_1(\mathbf{A}) - \lambda_2(\mathbf{A})}{\lambda_1(\mathbf{A})}$ . While it is well known that this can be done by applying  $\mathbf{A}$  to a random initial vector  $O(\frac{\alpha}{g} \log(\frac{d}{\varepsilon}))$  times in the following theorem we provide a slightly less known refinement that the dimension  $d$  can be replaced with the stable rank of  $\mathbf{A}$ ,  $s = \sum_{i \in [d]} \frac{\lambda_i(\mathbf{A})}{\lambda_1(\mathbf{A})}$ . We use this fact to avoid a dependence on  $d$  in our logarithmic factors.

---

#### Algorithm 17: PowerMethod( $\mathbf{A}, k$ )

---

**Input:** symmetric PSD matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  and a number of iterations  $k \geq 1$ .

Let  $x \sim \mathcal{N}(0, 1)$  be drawn from a  $d$  dimensional normal distribution.

Let  $y = \mathbf{A}^k x$

**Output:**  $u = y / \|y\|_2$

---

*Lemma 7.6.1 (Power Method).* Let  $\mathbf{A} \in \mathbb{R}^{d \times d}$  be a symmetric PSD matrix, let  $g \stackrel{\text{def}}{=} \frac{\lambda_1(\mathbf{A}) - \lambda_2(\mathbf{A})}{\lambda_1(\mathbf{A})}$ ,  $s = \sum_{i \in [d]} \frac{\lambda_i}{\lambda_1}$ , and let  $\varepsilon > 0$  and  $k \geq \frac{\alpha}{g} \log(\frac{n}{\varepsilon})$  for large enough constant  $\alpha$ . In time  $O(\text{nnz}(\mathbf{A}) \cdot \log(\frac{n}{\varepsilon}))$ , the algorithm PowerMethod( $\mathbf{A}, k$ ) outputs a vector  $u$  such that  $\langle v_1(\mathbf{A}), u \rangle^2 \geq 1 - \varepsilon$  and  $u^\top \mathbf{A} u \geq (1 - \varepsilon)\lambda_1(\mathbf{A})$  with high probability in  $n/\varepsilon$ .

*Proof.* We write  $u = \sum_{i \in [d]} \alpha_i v_i(\mathbf{A})$ . Then, we have

$$\langle v_1(\mathbf{A}), u \rangle^2 = \left\langle v_1(\mathbf{A}), \frac{\sum_{i \in [d]} \alpha_i \lambda_i(\mathbf{A})^k v_i(\mathbf{A})}{\sqrt{\sum_{i \in [d]} \alpha_i^2 \lambda_i(\mathbf{A})^{2k}}} \right\rangle^2 = \frac{\alpha_1^2}{\alpha_1^2 + \sum_{j \neq 1} \alpha_j^2 \left(\frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})}\right)^{2k}} \geq 1 - \sum_{j \neq 1} \frac{\alpha_j^2}{\alpha_1^2} \left(\frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})}\right)^{2k}$$

Re arranging terms we have

$$\begin{aligned} 1 - \langle v_1(\mathbf{A}), u \rangle^2 &\leq \sum_{j \neq 1} \frac{\alpha_j^2}{\alpha_1^2} \left(\frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})}\right) \left(\frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})}\right)^{2k-1} \leq \sum_{j \neq 1} \frac{\alpha_j^2}{\alpha_1^2} \left(\frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})}\right) \left(\frac{\lambda_2(\mathbf{A})}{\lambda_1(\mathbf{A})}\right)^{2k-1} \\ &= \sum_{j \neq 1} \frac{\alpha_j^2}{\alpha_1^2} \cdot \left(\frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})}\right) \cdot (1 - g)^{2k-1} \leq \sum_{j \neq 1} \frac{\alpha_j^2}{\alpha_1^2} \cdot \left(\frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})}\right) \cdot \exp(-(2k - 1)g) \end{aligned}$$

where we used that  $\frac{\lambda_2}{\lambda_1} = 1 - g \leq e^{-g}$ .

Now with high probability in  $n$  we have that  $\alpha_1^2 \geq \frac{1}{O(\text{poly}(n/\varepsilon))}$  by known properties of the chi-squared distribution. All that remains is to upper bound  $\sum_{j \neq 1} \alpha_j^2 \cdot \left(\frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})}\right)$ . To bound this consider  $h(\alpha) \stackrel{\text{def}}{=} \sqrt{\sum_{j \neq 1} \alpha_j^2 (\lambda_j(\mathbf{A})/\lambda_1(\mathbf{A}))}$ . Note that

$$\|\nabla h(\alpha)\|_2 = \left\| \frac{\sum_{j \neq 1} \vec{1}_j \cdot \alpha_j \left(\frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})}\right)}{\sqrt{\sum_{j \neq 1} \alpha_j^2 \left(\frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})}\right)}} \right\|_2 = \sqrt{\frac{\sum_{j \neq 1} \alpha_j^2 \left(\frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})}\right)^2}{\sum_{j \neq 1} \alpha_j^2 \left(\frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})}\right)}} \leq 1.$$

where  $\vec{1}_j$  is the indicator vector for coordinate  $j$ . Consequently  $h$  is 1-Lipschitz and by Gaussian concentration for Lipschitz functions we know there are absolute constants  $C$  and  $c$  such that

$$\Pr[h(\alpha) \geq \mathbf{E}h(\alpha) + \lambda] \leq C \exp(-c\lambda^2).$$

By the concavity of square root and the expected value of the chi-squared distribution we have

$$\mathbf{E}h(\alpha) \leq \sqrt{\mathbf{E} \sum_{j \neq i} \alpha_j^2 \cdot \left( \frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})} \right)} = \sqrt{\sum_{j \neq i} \left( \frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})} \right)} \leq \sqrt{s}.$$

Consequently, since  $s \geq 1$  we have that  $\Pr[h(\alpha) \geq (1 + \sqrt{t})^2 \cdot s] \leq C \exp(-c \cdot t)$  for  $t \geq 1$  and that  $\sum_{j \neq 1} \alpha_j \cdot \left( \frac{\lambda_j(\mathbf{A})}{\lambda_1(\mathbf{A})} \right) = O(\text{poly}(n/\varepsilon) \cdot s)$  with high probability in  $n/\varepsilon$ . Since  $k = \Omega(\frac{1}{g} \log(\frac{ns}{\varepsilon}))$ , we have  $\langle v_1(\mathbf{A}), u \rangle^2 \geq 1 - \varepsilon$  with high probability in  $n$ . Furthermore, this implies that

$$u^\top \mathbf{A} u = u^\top \left( \sum_{i \in [d]} \lambda_i(\mathbf{A}) v_i(\mathbf{A}) v_i(\mathbf{A})^\top \right) u \geq \lambda_1(\mathbf{A}) \langle v_1(\mathbf{A}), u \rangle^2 \geq (1 - \varepsilon) \lambda_1(\mathbf{A}).$$

□

*Lemma 7.4.1* (Computing Hessian Approximation). Let  $x \in \mathbb{R}^d$ ,  $t > 0$ , and  $\varepsilon \in (0, \frac{1}{4})$ . The algorithm `ApproxMinEig`( $x, t, \varepsilon$ ) outputs  $(\lambda, u)$  in  $O(nd \log \frac{n}{\varepsilon})$  time such that if  $\mu_t(x) \leq \frac{1}{4} t^2 w_t(x)$  then  $\langle v_t(x), u \rangle^2 \geq 1 - \varepsilon$  with high probability in  $n$ . Furthermore, if  $\varepsilon \leq \frac{\mu_t(x)}{8t^2 \cdot w_t(x)}$  then  $\frac{1}{4} \mathbf{Q} \preceq \nabla^2 f_t(x) \preceq 4\mathbf{Q}$  with high probability in  $n$  where  $\mathbf{Q} \stackrel{\text{def}}{=} t^2 \cdot w_t(x) - (t^2 \cdot w_t(x) - \lambda) uu^\top$ .

*Proof of Lemma 7.4.1.* By Lemma 7.5.4 we know that  $\frac{1}{2} \mathbf{Z} \preceq \nabla^2 f_t(x) \preceq \mathbf{Z}$  where

$$\mathbf{Z} = t^2 \cdot w_t(x) - (t^2 \cdot w_t(x) - \mu_t(x)) v_t(x) v_t(x)^\top.$$

Consequently, if  $\mu_t(x) \leq \frac{1}{4} t^2 w_t(x)$ , then for all unit vectors  $w \perp v_t(x)$ , we have that

$$w^\top \nabla^2 f_t(x) w \geq \frac{1}{2} w^\top \mathbf{Z} w \geq \frac{1}{2} t^2 w_t(x).$$

Since  $\nabla^2 f_t(x) = t^2 \cdot w_t(x) - \mathbf{A}$ , for  $\mathbf{A}$  in the definition of `ApproxMinEig` (Algorithm 14) this implies that  $v_t(x)^\top \mathbf{A} v_t(x) \geq \frac{3}{4} t^2 \cdot w_t(x)$  and  $w^\top \mathbf{A} w \leq \frac{1}{2} t^2 w_t(x)$ . Furthermore, we see that

$$\sum_{i \in [d]} \lambda_i(\mathbf{A}) = \text{Tr}(\mathbf{A}) = \sum_{i \in [n]} \frac{t^4 \|x - a^{(i)}\|_2^2}{(1 + g_t^{(i)}(x))^2 g_t^{(i)}(x)} \leq t^2 \cdot w_t(x)$$

Therefore, in this case,  $\mathbf{A}$  has a constant multiplicative gap between its top two eigenvectors stable rank at most a constant (i.e.  $g = O(1)$  and  $s = O(1)$  in Theorem 7.6.1). Consequently, by Theorem 7.6.1 we have  $\langle v_t(x), u \rangle^2 \geq 1 - \varepsilon$ .

For the second claim, we note that

$$t^2 \cdot w_t(x) - \mu_t(x) \geq u^\top \mathbf{A} u \geq (1 - \varepsilon) \lambda_1(\mathbf{A}) = (1 - \varepsilon)(t^2 \cdot w_t(x) - \mu_t(x))$$

Therefore, since  $\lambda = u^\top \nabla^2 f_t(x) u = t^2 \cdot w_t(x) - u^\top \mathbf{A} u$ , we have

$$(1 - \varepsilon) \mu_t(x) - \varepsilon \cdot t^2 w_t(x) \leq \lambda \leq \mu_t(x). \quad (7.7)$$



On the other hand, by Lemma 7.9.2, we have that

$$\sqrt{\varepsilon}\mathbf{I} \preceq v_t(x)v_t(x)^\top - uu^\top \preceq \sqrt{\varepsilon}\mathbf{I}. \quad (7.8)$$

Combining (7.7) and (7.8), we have  $\frac{1}{2}\mathbf{Z} \preceq \mathbf{Q} \preceq 2\mathbf{Z}$  if  $\varepsilon \leq \frac{\mu_t(x)}{8t^2 \cdot w_t(x)}$  and  $\frac{1}{4}\mathbf{Q} \preceq \nabla^2 f_t(x) \preceq 4\mathbf{Q}$  follows. On the other hand, when  $\mu_t(x) > \frac{1}{4}t^2 w_t(x)$ , it is the case that  $\frac{1}{4}t^2 \cdot w_t(x)\mathbf{I} \preceq \nabla^2 f_t(x) \preceq t^2 \cdot w_t(x)\mathbf{I}$  and  $\frac{1}{4}t^2 \cdot w_t(x)\mathbf{I} \preceq \mathbf{Q} \preceq t^2 \cdot w_t(x)\mathbf{I}$  again yielding  $\frac{1}{4}\mathbf{Q} \preceq \nabla^2 f_t(x) \preceq 4\mathbf{Q}$ .  $\square$

*Lemma 7.4.2.* Let  $u = \text{ApproxMinEig}(x, t, \varepsilon_v)$  for  $\varepsilon_v < \frac{1}{8}$  and  $x$  such that  $\|x - x_t\|_2 \leq \frac{\varepsilon_c}{t}$  for  $\varepsilon_c \leq (\frac{\varepsilon_v}{36})^{\frac{3}{2}}$ .  $\mu_t \leq \frac{1}{4}t^2 \cdot w_t$ . For all unit vectors  $y \perp u$ , we have  $\langle y, v_t \rangle^2 \leq 8\varepsilon_v$ .

*Proof of Lemma 7.4.2.* By Lemma 7.4.1 we know that  $\langle v_t(x), u \rangle^2 \geq 1 - \varepsilon_v$ . Since clearly  $\|x - x_t\|_2 \leq \frac{1}{20t}$ , by assumption, Lemma 7.3.1 shows

$$(1 - 6\varepsilon_c^{2/3})\nabla^2 f_t(x_t) \preceq \nabla^2 f_t(x) \preceq (1 + 6\varepsilon_c^{2/3})\nabla^2 f_t(x_t).$$

Furthermore, since  $\mu_t \leq \frac{1}{4}t^2 \cdot w_t$ , as in Lemma 7.4.1 we know that the largest eigenvalue of  $\mathbf{A}$  defined in  $\text{ApproxMinEig}(x, t, \varepsilon)$  is at least  $\frac{3}{4}t^2 \cdot w_t$  while the second largest eigenvalue is at most  $\frac{1}{2}t^2 \cdot w_t$ . Consequently, the eigenvalue gap,  $g$ , defined in Lemma 7.9.3 is at least  $\frac{1}{3}$  and this lemma shows that  $\langle v_t(x), v_t \rangle^2 \geq 1 - 36\varepsilon_c^{2/3} \geq 1 - \varepsilon_v$ . Consequently, by Lemma 7.9.4, we have that  $\langle u, v_t \rangle^2 \geq 1 - 4\varepsilon_v$ .

To prove the final claim, we write  $u = \alpha v_t + \beta w$  for an unit vector  $w \perp v_t$ . Since  $y \perp u$ , we have that  $0 = \alpha \langle v_t, y \rangle + \beta \langle w, y \rangle$ . Then, either  $\langle v_t, y \rangle = 0$  and the result follows or  $\alpha^2 \langle v_t, y \rangle^2 = \beta^2 \langle w, y \rangle^2$  and since  $\alpha^2 + \beta^2 = 1$ , we have

$$\langle v_t, y \rangle^2 \leq \frac{\beta^2 \langle w, y \rangle^2}{\alpha^2} \leq \frac{1 - \alpha^2}{\alpha^2} \leq 2(1 - \alpha^2) \leq 8\varepsilon_v$$

where in the last line we used that  $\alpha^2 \geq 1 - 4\varepsilon_v > \frac{1}{2}$  since  $\varepsilon_v \leq \frac{1}{8}$ .  $\square$

*Lemma 7.4.3.* Suppose  $\mu_t \geq \frac{1}{4}t^2 \cdot w_t$  and let  $t' \in [t, (1 + \frac{1}{600})t]$  then  $\|x_{t'} - x_t\|_2 \leq \frac{1}{100t}$ .

*Proof of Lemma 7.4.3.* Note that  $t' = (1 + \beta)t$  where  $\beta \in [0, \frac{1}{600}]$ . Since  $\frac{1}{4}t^2 \cdot w_t \mathbf{I} \preceq \mu_t \mathbf{I} \preceq \nabla^2 f(x_t)$  applying Lemma 7.3.3 then yields that for all  $s \in [t, t']$

$$\nabla^2 f(x_s) \succeq \nabla^2 f(x_t) - 15\beta t^2 w_t \mathbf{I} \succeq \left(\frac{1}{4} - 15\beta\right) t^2 \cdot w_t \mathbf{I} \succeq \frac{t^2 \cdot w_t}{6} \mathbf{I}.$$

Consequently, by Lemma 7.5.1, the fact that  $t\|x_t - a^{(i)}\|_2 \leq g_t^{(i)}$ , and Lemma 7.3.2 we have

$$\begin{aligned} \|x_{t'} - x_t\|_2 &\leq \int_t^{t'} \left\| \frac{d}{ds} x_s \right\|_2 ds = \int_t^{t'} \left\| \left( \nabla^2 f_s(x_s) \right)^{-1} \sum_{i \in [n]} \frac{s}{(1 + g_s^{(i)}) g_s^{(i)}} (x_s - a^{(i)}) \right\|_2 ds \\ &\leq \int_t^{t'} \frac{6}{t^2 \cdot w_t} \sum_{i \in [n]} \frac{s \|x_s - a^{(i)}\|_2}{(1 + g_s^{(i)}) g_s^{(i)}} ds \leq \int_t^{t'} \frac{6w_s}{t^2 \cdot w_t} ds \leq \int_t^{t'} \frac{6}{t^2} \cdot \left(\frac{s}{t}\right)^2 ds \\ &= \frac{6}{3t^4} [(t')^2 - (t)^3] = \frac{6}{3t} [(1 + \beta)^3 - 1] \leq \frac{6\beta}{t} \leq \frac{1}{100t} \end{aligned}$$

where in the last line we used that for  $\alpha \geq 0$  we have  $1 + 3\alpha \leq (1 + \alpha)^3$ .  $\square$

## ■ 7.6.2 Line Searching

Here we prove the main results we use on centering, Lemma 7.4.4, and line searching Lemma 7.4.5. These results are our main tools for computing approximations to the central path. To prove Lemma 7.4.5 we also include here two preliminary lemmas, Lemma 7.6.2 and Lemma 7.6.3, on the structure of  $g_{t,y,v}$  defined in (7.5).

*Lemma 7.4.4.* Given some  $y \in \mathbb{R}^d$ ,  $t > 0$  and  $\varepsilon \in (0, \min\{\frac{1}{4}, \frac{\mu_t(x)}{8t^2 \cdot w_t(x)}\})$ . In  $O(nd \log(\frac{n}{\varepsilon}))$  time  $\text{LocalCenter}(y, t, \varepsilon)$  computes  $x^{(k)}$  such that with high probability in  $n/\varepsilon$ .

$$f_t(x^{(k)}) - \min_{\|x-y\|_2 \leq \frac{1}{49t}} f_t(x) \leq \varepsilon \left( f_t(y) - \min_{\|x-y\|_2 \leq \frac{1}{49t}} f_t(x) \right).$$

*Proof of Lemma 7.4.4.* By Lemma 7.4.1 we know that  $\frac{1}{4}\mathbf{Q} \preceq \nabla^2 f_t(y) \preceq 4\mathbf{Q}$  with high probability in  $n/\varepsilon$ . Furthermore for  $x$  such that  $\|x-y\|_2 \leq \frac{1}{50t}$  Lemma 7.3.1 shows that  $\frac{1}{2}\nabla^2 f_t(x) \preceq \nabla^2 f_t(y) \preceq 2\nabla^2 f_t(x)$ . Combining these we have that  $\frac{1}{8}\mathbf{Q} \preceq \nabla^2 f_t(x) \preceq 8\mathbf{Q}$  for all  $x$  with  $\|x-y\|_2 \leq \frac{1}{50t}$ . Therefore, Lemma 7.9.5 shows that

$$f_t(x^{(k)}) - \min_{\|x-y\|_2 \leq \frac{1}{49t}} f_t(x) \leq \left(1 - \frac{1}{64}\right)^k \left( f_t(x^{(0)}) - \min_{\|x-y\|_2 \leq \frac{1}{49t}} f_t(x) \right).$$

The guarantee on  $x^{(k)}$  then follows from our choice of  $k$ .

For the running time, Lemma 7.4.1 showed the cost of  $\text{ApproxMinEig}$  is  $O(nd \log(\frac{n}{\varepsilon}))$ . Using Lemma 7.9.7 we see that the cost per iteration is  $O(nd)$  and therefore, the total cost of the  $k$  iterations is  $O(nd \log(\frac{1}{\varepsilon}))$ . Combining yields the running time.  $\square$

*Lemma 7.6.2.* For  $t > 0$ ,  $y \in \mathbb{R}^d$ , and unit vector  $v \in \mathbb{R}^d$ , the function  $g_{t,y,v} : \mathbb{R} \rightarrow \mathbb{R}$  defined by (7.5) is convex and  $nt$ -Lipschitz.

*Proof.* Changing variables yields  $g_{t,y,v}(\alpha) = \min_{z \in S} f_t(z + \alpha v)$  for  $S \stackrel{\text{def}}{=} \{z \in \mathbb{R}^d : \|z-y\|_2 \leq \frac{1}{49t}\}$ . Since  $f_t$  is convex and  $S$  is a convex set, by Lemma 7.9.6 we have that  $g_{t,y,v}$  is convex.

Next, by Lemma 7.5.1, triangle inequality, and the fact that  $t\|x - a^{(i)}\|_2 \leq g_t^{(i)}(x)$  we have

$$\|\nabla f_t(x)\|_2 = \left\| \sum_{i \in [n]} \frac{t^2(x - a^{(i)})}{1 + g_t^{(i)}(x)} \right\|_2 \leq \sum_{i \in [n]} \frac{t^2\|x - a^{(i)}\|_2}{1 + g_t^{(i)}(x)} \leq tn. \quad (7.9)$$

Consequently,  $f_t(x)$  is  $nt$ -Lipschitz and for all  $x, y \in \mathbb{R}^n$  we have  $|f_t(x) - f_t(y)| \leq nt\|x-y\|_2$ . Now if we consider the set  $S_\alpha \stackrel{\text{def}}{=} \{x \in \mathbb{R}^d : \|x - (y + \alpha v)\|_2 \leq \frac{1}{49t}\}$  then we see that for all  $\alpha, \beta \in \mathbb{R}$  there is a bijection from  $S_\alpha$  to  $S_\beta$  where every point in the set moves by at most  $\|(\alpha - \beta)v\|_2 \leq |\alpha - \beta|$ . Consequently, since  $g_{t,y,v}(\alpha)$  simply minimizes  $f_t$  over  $S_\alpha$  we have that  $|g_{t,y,v}(\alpha) - g_{t,y,v}(\beta)| \leq nt|\alpha - \beta|$  for all  $\alpha, \beta \in \mathbb{R}$  we have that  $g_{t,y,v}$  is  $nt$ -Lipschitz as desired.  $\square$

*Lemma 7.6.3.* Let  $\frac{1}{400f(x_*)} \leq t \leq t' \leq (1 + \frac{1}{600})t \leq \frac{2n}{\varepsilon \cdot f_*}$  and let  $u = \text{ApproxMinEig}(y, t, \varepsilon_v)$  for  $\varepsilon_v \leq \frac{1}{8}(\frac{\varepsilon_*}{7n})^2$  and  $y \in \mathbb{R}^d$  such that  $\|y - x_t\|_2 \leq \frac{1}{t}(\frac{\varepsilon_v}{36})^{\frac{3}{2}}$ . The function  $g_{t',y,v} : \mathbb{R} \rightarrow \mathbb{R}$  defined in (7.5) satisfies  $g_{t',y,v}(\alpha_*) = \min_\alpha g_{t',y,v}(\alpha) = f_{t'}(x_{t'})$  for some  $\alpha_* \in [-6f(x_*), 6f(x_*)]$ .

*Proof of Lemma 7.4.5.* Let  $z \in \mathbb{R}^d$  be an arbitrary unit vector. By Lemma 7.3.5 we know if  $|\langle z, v_t \rangle| \leq \min_{\delta \in [t, (1+\beta)t]} \frac{\mu_\delta}{t^2 \cdot w_\delta}$  for  $\beta \in [0, \frac{1}{600}]$  then  $z^\top(x_{(1+\beta)t} - x_t) \leq \frac{6\beta}{t}$ . Now by Lemma 7.5.7 and our bound on  $\beta$  we know that  $\frac{\delta^2 \cdot w_\delta}{\mu_\delta} \leq 7n\varepsilon_*^{-1}$  and consequently  $\frac{\varepsilon_*}{6n} \leq \min_{\delta \in [t, (1+\beta)t]} \frac{w_\delta}{t^2 \cdot \mu_\delta}$ . Consequently, if  $\mu_t \leq$

$\frac{1}{4}t^2 \cdot w_t$  then by Lemma 7.4.2 and our choice of  $\varepsilon_v$  we have that if  $z \perp u$  then  $|\langle z, v_t \rangle|^2 \leq (\frac{\tilde{\varepsilon}_*}{6n})^2$  and  $z^\top(x_{t'} - x_t) \leq \frac{6}{600t} = \frac{1}{100t}$ . Otherwise,  $\mu_t \geq \frac{1}{4}t^2 \cdot w_t$  and by Lemma 7.4.3 we have  $\|x_{t'} - x_t\|_2 \leq \frac{1}{100t}$ . In either case, since  $\|y - x_t\|_2 \leq \frac{1}{100t}$ , we can reach  $x_{t'}$  from  $y$  by first moving an Euclidean distance of  $\frac{1}{100t}$  to go from  $x$  to  $x_t$ , then adding some multiple of  $v$ , then moving an Euclidean distance of  $\frac{1}{100t}$  in a direction perpendicular to  $v$ . Since the total movement perpendicular to  $v$  is  $\frac{1}{100t} + \frac{1}{100t} \leq \frac{1}{49t'}$  we have that  $\min_{\alpha} g_{t',y,v}(\alpha) = f_{t'}(x_{t'})$  as desired.

All that remains is to show that there is a minimizer of  $g_{t,y,v}$  in the range  $[-6f(x_*), 6f(x_*)]$ . However, by Lemma 7.3.6 and Lemma 7.5.5 we know that

$$\|y - x_{t'}\|_2 \leq \|y - x_t\|_2 + \|x_t - x_*\|_2 + \|x_* - x_{t'}\|_2 \leq \frac{1}{100t} + f(x_*) + f(x_*) \leq 6f(x_*).$$

Consequently,  $\alpha_* \in [-6f(x_*), 6f(x_*)]$  as desired.

*Lemma 7.4.5.* Let  $\frac{1}{400f(x_*)} \leq t \leq t' \leq (1 + \frac{1}{600})t \leq \frac{2n}{\tilde{\varepsilon} \cdot \tilde{f}_*}$  and let  $u = \text{ApproxMinEig}(y, t, \varepsilon_v)$  for  $\varepsilon_v \leq \frac{1}{8}(\frac{\tilde{\varepsilon}_*}{7n})^2$  and  $y \in \mathbb{R}^d$  such that  $\|y - x_t\|_2 \leq \frac{1}{t}(\frac{\varepsilon_v}{36})^{\frac{3}{2}}$ . In  $O(nd \log^2(\frac{n}{\tilde{\varepsilon} \cdot \varepsilon_v}))$  time,  $\text{LineSearch}(y, t, t', u, \varepsilon)$  outputs  $x'$  such that  $\|x' - x_{t'}\|_2 \leq \frac{\varepsilon}{t'}$  with high probability in  $n$ .

By (7.9) we know that  $f_{t'}$  is  $nt'$  Lipschitz and therefore

$$f_t(y) - \min_{\|x-y\|_2 \leq \frac{1}{49t'}} f_{t'}(x) \leq \frac{nt'}{49t'} = \frac{n}{49}.$$

Furthermore, for  $\alpha \in [-6\tilde{f}_*, 6\tilde{f}_*]$  we know that by Lemma 7.5.5

$$\|y + \alpha u - x_{t'}\|_2 \leq \|y - x_t\|_2 + |\alpha| + \|x_t - x_*\|_2 + \|x_{t'} - x_*\|_2 \leq \frac{1}{t'} + 14f(x_*)$$

consequently by Lemma 7.5.7 we have

$$\frac{(t')^2 \cdot w_{t'}(y + \alpha u)}{\mu_{t'}(y + \alpha u)} \leq 6n\tilde{\varepsilon}^{-1} + t' \cdot 14f(x_*) \leq 34n\tilde{\varepsilon}^{-1}.$$

Since  $\varepsilon_O \leq \frac{1}{8} \cdot \frac{\tilde{\varepsilon}_*}{34n}$  invoking Lemma 7.4.4 yields that  $|q(\alpha) - g_{t',y,u}(\alpha)| \leq \frac{n\varepsilon_O}{49}$  with high probability in  $n/\tilde{\varepsilon}_*$  and each call to  $\text{LocalCenter}$  takes  $O(nd \log \frac{n}{\varepsilon_O})$  time. Furthermore, by Lemma 7.6.2 we have that  $g_{t',y,u}$  is a  $nt'$ -Lipschitz convex function and by Lemma 7.6.3 we have that the minimizer has value  $f_{t'}(x_{t'})$  and is achieved in the range  $[-12\tilde{f}_*, 12\tilde{f}_*]$ . Consequently, combining all these facts and invoking Lemma 7.9.3, i.e. our result on one dimensional function minimization, we have  $f_{t'}(x') - f_{t'}(x_{t'}) \leq \frac{\varepsilon_O}{49t'}$  using only  $O(\log \frac{nt'f(x_*)}{\varepsilon_O})$  calls to  $\text{LocalCenter}$ .

Finally, by Lemma 7.5.8 and Lemma 7.3.6 we have

$$\|x - x_{t'}\| \leq \frac{1}{\sqrt{n}} \cdot \left( \frac{7n\tilde{\varepsilon}^{-1} + t'\alpha}{t'} \right) \cdot \sqrt{\frac{n\varepsilon_O}{49}} \leq \frac{34n}{7 \cdot \tilde{\varepsilon}_* \cdot t'} \sqrt{\varepsilon_O} \leq \frac{6n}{\tilde{\varepsilon}_* \cdot t'} \sqrt{\varepsilon_O}.$$

Since  $\varepsilon_O \leq \frac{\varepsilon_*^2 \cdot \tilde{\varepsilon}_*^2}{36n^2}$  we have that  $\|x - x_{t'}\|_2 \leq \frac{\varepsilon}{t'}$  as desired.  $\square$

*Lemma 7.4.6.* Let  $\frac{1}{400f(x_*)} \leq t \leq t' \leq (1 + \frac{1}{600})t \leq \frac{2n}{\tilde{\varepsilon}_* \cdot \tilde{f}_*}$  and let  $x \in \mathbb{R}^d$  satisfy  $\|x - x_t\|_2 \leq \frac{1}{100t}$ . Then, in  $O(nd \log^2(\frac{n}{\tilde{\varepsilon} \cdot \tilde{\varepsilon}_*}))$  time,  $\text{LineSearch}(x, t, t, u, \varepsilon)$  output  $y$  such that  $\|y - x_t\|_2 \leq \frac{\varepsilon}{t}$  for any vector  $u \in \mathbb{R}^d$ .

*Proof of Lemma 7.4.6.* The proof is strictly easier than the proof of Lemma 7.4.5 as  $\|x - \alpha^*u - x_t\|_2 \leq \frac{1}{100t}$  is satisfied automatically for  $\alpha^* = 0$ . Note that this lemma assume less for the initial point.  $\square$

### ■ 7.6.3 Putting It All Together

**Theorem 7.4.1.** *In  $O(nd \log^3(\frac{n}{\varepsilon}))$  time, Algorithm 13 outputs an  $(1 + \varepsilon)$ -approximate geometric median with constant probability.*

*Proof of Theorem 7.4.1.* By Lemma 7.5.6 we know that  $x^{(0)}$  is a 2-approximate geometric median and therefore  $f(x^{(0)}) = \tilde{f}_* \leq 2 \cdot f(x_*)$ . Furthermore, since  $\|x^{(0)} - x_{t_1}\|_2 \leq f(x^{(0)})$  by Lemma 7.5.5 and  $t_1 = \frac{1}{400\tilde{f}_*}$  we have  $\|x^{(0)} - x_{t_1}\|_2 \leq \frac{1}{400t_1}$ . Hence, by Lemma 7.4.6, we have  $\|x^{(1)} - x_{t_1}\|_2 \leq \frac{\varepsilon_c}{t_1}$  with high probability in  $n/\varepsilon$ . Consequently, by Lemma 7.4.5 we have that  $\|x^{(k)} - x_{t_i}\|_2 \leq \frac{\varepsilon_c}{t_i}$  for all  $i$  with high probability in  $n/\varepsilon$ .

Now, Lemma 7.3.6 shows that

$$f(x_{t_k}) - f(x_*) \leq \frac{2n}{t_k} \leq \frac{2n}{\tilde{t}_*} \left(1 + \frac{1}{600}\right) \leq \tilde{\varepsilon}_* \cdot \tilde{f}_* \left(1 + \frac{1}{600}\right) \leq \frac{2}{3} \left(1 + \frac{1}{600}\right) \varepsilon \cdot f(x_*).$$

Since  $\|x^{(k)} - x_{t_k}\|_2 \leq \frac{\varepsilon_c}{t_k} \leq 400 \cdot \tilde{f}_* \cdot \varepsilon_c$  we have that  $f(x_k) \leq f(x_{t_k}) + 400n \cdot \tilde{f}_* \cdot \varepsilon_c$  by triangle inequality. Combining these facts and using that  $\varepsilon_c$  is sufficiently small yields that  $f(x^{(k)}) \leq (1 + \varepsilon)f(x_*)$  as desired.

To bound the running time, Lemma 7.4.1 shows `ApproxMinEvec` takes  $O(nd \log(\frac{n}{\varepsilon}))$  per iteration and Lemma 7.4.5 shows `LineSearch` takes  $O(nd \log^2(\frac{n}{\varepsilon}))$  time per iteration, using that  $\varepsilon_v$  and  $\varepsilon_c$  are  $O(\Omega(\varepsilon/n))$ . Since for  $l = \Omega(\log \frac{n}{\varepsilon})$  we have that  $t_l > \tilde{t}_*$  we have that  $k = O(\log \frac{n}{\varepsilon})$ .  $t_{i+1} \leq \frac{1}{400}$ . Since there are  $O(\log(\frac{n}{\varepsilon}))$  iterations taking time  $O(nd \log^2(\frac{n}{\varepsilon}))$  the running time follows.  $\square$

## ■ 7.7 Pseudo Polynomial Time Algorithm

Here we provide a self-contained result on computing a  $1 + \varepsilon$  approximate geometric median in  $O(d\varepsilon^{-2})$  time. Note that it is impossible to achieve such approximation for the mean,  $\min_{x \in \mathbb{R}^d} \sum_{i \in [n]} \|x - a^{(i)}\|_2^2$ , because the mean can be changed arbitrarily by changing only 1 point. However, [172] showed that the geometric median is far more stable. In Section 7.7.1, we show how this stability property allows us to get an constant approximate in  $O(d)$  time. In Section 7.7.2, we show how to use stochastic subgradient descent to then improve the accuracy.

### ■ 7.7.1 A Constant Approximation of Geometric Median

We first prove that the geometric median is stable even if we are allowed to modify up to half of the points. The following lemma is a strengthening of the robustness result in [172].

*Lemma 7.7.1.* Let  $x_*$  be a geometric median of  $\{a^{(i)}\}_{i \in [n]}$  and let  $S \subseteq [n]$  with  $\|S\| < \frac{n}{2}$ . For all  $x$

$$\|x_* - x\|_2 \leq \left( \frac{2n - 2|S|}{n - 2|S|} \right) \max_{i \notin S} \|a^{(i)} - x\|_2.$$

*Proof.* For notational convenience let  $r = \|x_* - x\|_2$  and let  $M = \max_{i \notin S} \|a^{(i)} - x\|_2$ .

For all  $i \notin S$ , we have that  $\|x - a^{(i)}\|_2 \leq M$ , hence, we have

$$\begin{aligned} \|x_* - a^{(i)}\|_2 &\geq r - \|x - a^{(i)}\|_2 \\ &\geq r - 2M + \|x - a^{(i)}\|_2. \end{aligned}$$

Furthermore, by triangle inequality for all  $i \in S$ , we have

$$\|x_* - a^{(i)}\|_2 \geq \|x - a^{(i)}\|_2 - r.$$

Hence, we have that

$$\sum_{i \in [n]} \|x_* - a^{(i)}\|_2 \geq \sum_{i \in [n]} \|x - a^{(i)}\|_2 + (n - |S|)(r - 2M) - |S|r.$$

Since  $x^*$  is a minimizer of  $\sum_{i \in [n]} \|x_* - a^{(i)}\|_2$ , we have that

$$(n - |S|)(r - 2M) - |S|r \leq 0.$$

Hence, we have

$$\|x^* - x\|_2 = r \leq \frac{2n - 2|S|}{n - 2|S|}M.$$

□

Now, we use Lemma 7.7.1 to show that the algorithm `CrudeApproximate` outputs a constant approximation of the geometric median with high probability.

---

**Algorithm 18:** `CrudeApproximateK`

---

**Input:**  $a^{(1)}, a^{(2)}, \dots, a^{(n)} \in \mathbb{R}^d$ .

Sample two independent random subset of  $[n]$  of size  $K$ . Call them  $S_1$  and  $S_2$ .

Let  $i^* \in \arg \min_{i \in S_2} \alpha_i$  where  $\alpha_i$  is the 65 percentile of the numbers  $\{\|a^{(i)} - a^{(j)}\|_2\}_{j \in S_1}$ .

**Output:** Output  $a^{(i^*)}$  and  $\alpha_{i^*}$ .

---

*Lemma 7.7.2.* Let  $x_*$  be a geometric median of  $\{a^{(i)}\}_{i \in [n]}$  and  $(\tilde{x}, \lambda)$  be the output of `CrudeApproximateK`. We define  $d_T^k(x)$  be the  $k$ -percentile of  $\{\|x - a^{(i)}\|\}_{i \in T}$ . Then, we have that  $\|x_* - \tilde{x}\|_2 \leq 6d_{[n]}^{60}(\tilde{x})$ . Furthermore, with probability  $1 - e^{-\Theta(K)}$ , we have

$$d_{[n]}^{60}(\tilde{x}) \leq \lambda = d_{S_1}^{60}(\tilde{x}) \leq 2d_{[n]}^{70}(x_*).$$

*Proof.* Lemma 7.7.1 shows that for all  $x$  and  $T \subseteq [n]$  with  $|T| \leq \frac{n}{2}$

$$\|x_* - x\|_2 \leq \left( \frac{2n - 2|T|}{n - 2|T|} \right) \max_{i \notin T} \|a^{(i)} - x\|_2.$$

Picking  $T$  to be the indices of largest 40% of  $\|a^{(i)} - \tilde{x}\|_2$ , we have

$$\|x_* - \tilde{x}\|_2 \leq \left( \frac{2n - 0.8n}{n - 0.8n} \right) d_{[n]}^{60}(\tilde{x}) = 6d_{[n]}^{60}(\tilde{x}). \quad (7.10)$$

For any point  $x$ , we have that  $d_{[n]}^{60}(x) \leq d_{S_1}^{65}(x)$  with probability  $1 - e^{-\Theta(K)}$  because  $S_1$  is a random subset of  $[n]$  with size  $K$ . Taking union bound over elements on  $S_2$ , with probability  $1 - Ke^{-\Theta(K)} = 1 - e^{-\Theta(K)}$ , for all points  $x \in S_2$

$$d_{[n]}^{60}(x) \leq d_{S_1}^{65}(x). \quad (7.11)$$

yielding that  $d_{[n]}^{60}(\tilde{x}) \leq \lambda$ .

Next, for any  $i \in S_2$ , we have

$$\|a^{(i)} - a^{(j)}\|_2 \leq \|a^{(i)} - x_*\|_2 + \|x_* - a^{(j)}\|_2.$$

and hence

$$d_{[n]}^{70}(a^{(i)}) \leq \|a^{(i)} - x_*\|_2 + d_{[n]}^{70}(x_*).$$

Again, since  $S_1$  is a random subset of  $[n]$  with size  $K$ , we have that  $d_{S_1}^{65}(a^{(i)}) \leq d_{[n]}^{70}(a^{(i)})$  with probability  $1 - Ke^{-\Theta(K)} = 1 - e^{-\Theta(K)}$ . Therefore,

$$d_{S_1}^{65}(a^{(i)}) \leq \|a^{(i)} - x_*\|_2 + d_{[n]}^{70}(x_*).$$

Since  $S_2$  is an independent random subset, with probability  $1 - e^{-\Theta(K)}$ , there is  $i \in S_2$  such that  $\|a^{(i)} - x_*\|_2 \leq d_{[n]}^{70}(x_*)$ . In this case, we have

$$d_{S_1}^{65}(a^{(i)}) \leq 2d_{[n]}^{70}(x_*).$$

Since  $i^*$  minimize  $d_{S_1}^{65}(a^{(i)})$  over all  $i \in S_2$ , we have that

$$\lambda \stackrel{\text{def}}{=} d_{S_1}^{65}(\tilde{x}) \stackrel{\text{def}}{=} d_{S_1}^{65}(a^{(i^*)}) \leq d_{S_1}^{65}(a^{(i)}) \leq 2d_{[n]}^{70}(x_*).$$

□

### ■ 7.7.2 A $1 + \varepsilon$ Approximation of Geometric Median

Here we show how to improve the constant approximation in the previous section to a  $1 + \varepsilon$  approximation. Our algorithm is essentially stochastic subgradient where we use the information from the previous section to bound the domain in which we need to search for a geometric median.

---

#### Algorithm 19: ApproximateMedian( $\varepsilon$ )

---

**Input:**  $a^{(1)}, a^{(2)}, \dots, a^{(n)} \in \mathbb{R}^d$ .

Let  $T = (120/\varepsilon)^2$  and let  $\eta = \frac{6\lambda}{n} \sqrt{\frac{2}{T}}$ .

Let  $(x^{(1)}, \lambda) = \text{CrudeApproximate}_{\sqrt{T}}(a^{(1)}, a^{(2)}, \dots, a^{(n)})$ .

**for**  $k \leftarrow 1, 2, \dots, T$  **do**

    Sample  $i_k$  from  $[n]$  and let

$$g^{(k)} = \begin{cases} n(x^{(k)} - a^{(i_k)}) / \|x^{(k)} - a^{(i_k)}\|_2 & \text{if } x^{(i)} \neq a^{(i_k)} \\ 0 & \text{otherwise} \end{cases}$$

    Let  $x^{(k+1)} = \arg \min_{\|x - x^{(1)}\|_2 \leq 6\lambda} \eta \langle g^{(k)}, x - x^{(k)} \rangle + \frac{1}{2} \|x - x^{(k)}\|_2^2$ .

**end**

**Output:** Output  $\frac{1}{T} \sum_{i=1}^T x^{(k)}$ .

---

**Theorem 7.7.1.** *Let  $x$  be the output of ApproximateMedian( $\varepsilon$ ). With probability  $1 - e^{-\Theta(1/\varepsilon)}$ , we have*

$$\mathbb{E}f(x) \leq (1 + \varepsilon) \min_{x \in \mathbb{R}^d} f(x).$$

Furthermore, the algorithm takes  $O(d/\varepsilon^2)$  time.

*Proof.* After computing  $x^{(1)}$  and  $\lambda$  the remainder of our algorithm is the stochastic subgradient descent method applied to  $f(x)$ . It is routine to check that  $\mathbf{E}_{i^{(k)}} g^{(k)}$  is a subgradient of  $f$  at  $x^{(k)}$ . Furthermore, since the diameter of the domain,  $\{x : \|x - x^{(1)}\|_2 \leq 6\lambda\}$ , is clearly  $\lambda$  and the norm of sampled gradient,  $g^{(k)}$ , is at most  $n$ , we have that

$$\mathbb{E}f\left(\frac{1}{T} \sum_{i=1}^T x^{(k)}\right) - \min_{\|x - x^{(1)}\|_2 \leq 6\lambda} f(x) \leq 6n\lambda \sqrt{\frac{2}{T}}$$

(see [41, Thm 6.1]). Lemma 7.7.2 shows that  $\|x^* - x^{(1)}\|_2 \leq 6\lambda$  and  $\lambda \leq 2d_{[n]}^{70}(x^*)$  with probability  $1 - \sqrt{T}e^{-\Theta(\sqrt{T})}$ . In this case, we have

$$\mathbb{E}f\left(\frac{1}{T}\sum_{i=1}^T x^{(k)}\right) - f(x^*) \leq \frac{12\sqrt{2}nd_{[n]}^{70}(x_*)}{\sqrt{T}}.$$

Since  $d_{[n]}^{70}(x^*) \leq \frac{1}{0.3n}f(x^*)$ , we have

$$\mathbb{E}f\left(\frac{1}{T}\sum_{i=1}^T x^{(k)}\right) \leq \left(1 + \frac{60}{\sqrt{T}}\right)f(x_*) \leq \left(1 + \frac{\varepsilon}{2}\right)f(x_*).$$

□

## ■ 7.8 Derivation of Penalty Function

Here we derive our penalized objective function. Consider the following optimization problem:

$$\min_{x \in \mathbb{R}^d, \alpha \geq 0 \in \mathbb{R}^n} f_t(x, \alpha) \quad \text{where} \quad t \cdot 1^T \alpha + \sum_{i \in [n]} -\ln\left(\alpha_i^2 - \|x - a^{(i)}\|_2^2\right).$$

Since  $p_i(\alpha, x) \stackrel{\text{def}}{=} -\ln\left(\alpha_i^2 - \|x - a^{(i)}\|_2^2\right)$  is a barrier function for the set  $\alpha_i^2 \geq \|x - a^{(i)}\|_2^2$ , i.e. as  $\alpha_i \rightarrow \|x - a^{(i)}\|_2$  we have  $p_i(\alpha, x) \rightarrow \infty$ , we see that as we minimize  $f_t(x, \alpha)$  for increasing values of  $t$  the  $x$  values converge to a solution to the geometric median problem. Our penalized objective function,  $f_t(x)$ , is obtain simply by minizing the  $\alpha_i$  in the above formula and dropping terms that do not affect the minimizing  $x$ . In the remainder of this section we show this formally.

Fix some  $x \in \mathbb{R}^d$  and  $t > 0$ . Note that for all  $j \in [n]$  we have

$$\frac{\partial}{\partial \alpha_j} f_t(x, \alpha) = t - \left(\frac{1}{\alpha_j^2 - \|x - a^{(j)}\|_2^2}\right) 2\alpha_j.$$

Since  $f(x, \alpha)$  is convex in  $\alpha$ , the minimum  $\alpha_j^*$  must satisfy

$$t \left( (\alpha_j^*)^2 - \|x - a^{(j)}\|_2^2 \right) - 2\alpha_j^* = 0. \quad (7.12)$$

Solving for such  $\alpha_j^*$  under the restriction  $\alpha_j^* \geq 0$  we obtain

$$\alpha_j^* = \frac{2 + \sqrt{4 + 4t^2\|x - a^{(j)}\|_2^2}}{2t} = \frac{1}{t} \left[ 1 + \sqrt{1 + t^2\|x - a^{(j)}\|_2^2} \right]. \quad (7.13)$$

Using (7.12) and (7.13) we have that

$$\min_{\alpha \geq 0 \in \mathbb{R}^n} f_t(x, \alpha) = \sum_{i \in [n]} \left[ 1 + \sqrt{1 + t^2\|x - a^{(i)}\|_2^2} - \ln \left[ \frac{2}{t^2} \left( 1 + \sqrt{1 + t^2\|x - a^{(i)}\|_2^2} \right) \right] \right].$$

If we drop the terms that do not affect the minimizing  $x$  we obtain our penalty function  $f_t$ :

$$f_t(x) = \sum_{i \in [n]} \left[ \sqrt{1 + t^2\|x - a^{(i)}\|_2^2} - \ln \left( 1 + \sqrt{1 + t^2\|x - a^{(i)}\|_2^2} \right) \right].$$

## ■ 7.9 Technical Lemmas

Here we provide various technical lemmas we use through this chapter.

### ■ 7.9.1 Linear Algebra

First we provide the following lemma that shows that any matrix obtained as a non-negative linear combination of the identity minus a rank 1 matrix less than the identity results in a matrix that is well approximated spectrally by the identity minus a rank 1 matrix. We use this lemma to characterize the Hessian of our penalized objective function and thereby imply that it is possible to apply the inverse of the Hessian to a vector with high precision.

*Lemma 7.9.1.* Let  $\mathbf{A} = \sum_i (\alpha_i \mathbf{I} - \beta_i a_i a_i^\top) \in \mathbb{R}^{d \times d}$  where the  $a_i$  are unit vectors and  $0 \leq \beta_i \leq \alpha_i$  for all  $i$ . Let  $v$  denote a unit vector that is the maximum eigenvector of  $\sum_i \beta_i a_i a_i^\top$  and let  $\lambda$  denote the corresponding eigenvalue. Then,

$$\frac{1}{2} \left( \sum_i \alpha_i \mathbf{I} - \lambda v v^\top \right) \preceq \mathbf{A} \preceq \sum_i \alpha_i \mathbf{I} - \lambda v v^\top.$$

*Proof.* Let  $\alpha \stackrel{\text{def}}{=} \sum_i \alpha_i$ . Since clearly  $v^\top \mathbf{A} v = v^\top (\sum_i \alpha_i \mathbf{I} - \lambda v v^\top) v$  it suffices to show that for  $w \perp v$  it is the case that  $\frac{1}{2} \alpha \|w\|_2^2 \preceq w^\top \mathbf{A} w \preceq \alpha \|w\|_2^2$  or equivalently, that  $\lambda_i(\mathbf{A}) \in [\frac{1}{2} \alpha, \alpha]$  for  $i \neq d$ . However we know that  $\sum_{i \in [d]} \lambda_i(\mathbf{A}) = \text{Tr}(\mathbf{A}) = d\alpha - \sum_i \beta_i \geq (d-1)\alpha$  and  $\lambda_i(\mathbf{A}) \leq \alpha$  for all  $i \in [d]$ . Consequently, since  $\lambda_d(\mathbf{A}) \in [0, \lambda_{d-1}]$  we have

$$2 \cdot \lambda_{d-1}(\mathbf{A}) \geq (d-1)\alpha - \sum_{i=1}^{d-2} \lambda_i(\mathbf{A}) \geq (d-1)\alpha - (d-2)\alpha = \alpha.$$

Consequently,  $\lambda_{d-1}(\mathbf{A}) \in [\frac{\alpha}{2}, \alpha]$  and the result holds by the monotonicity of  $\lambda_i$ .  $\square$

Next we bound the spectral difference between the outer product of two unit vectors by their inner product. We use this lemma to bound the amount of precision required in our eigenvector computations.

*Lemma 7.9.2.* For unit vectors  $u_1$  and  $u_2$  we have

$$\|u_1 u_1^\top - u_2 u_2^\top\|_2^2 = 1 - (u_1^\top u_2)^2 \quad (7.14)$$

Consequently if  $(u_1^\top u_2)^2 \geq 1 - \varepsilon$  for  $\varepsilon \leq 1$  we have that

$$-\sqrt{\varepsilon} \mathbf{I} \preceq u_1 u_1^\top - u_2 u_2^\top \preceq \sqrt{\varepsilon} \mathbf{I}$$

*Proof.* Note that  $u_1 u_1^\top - u_2 u_2^\top$  is a symmetric matrix and all eigenvectors are either orthogonal to both  $u_1$  and  $u_2$  (with eigenvalue 0) or are of the form  $v = \alpha u_1 + \beta u_2$  where  $\alpha$  and  $\beta$  are real numbers that are not both 0. Thus, if  $v$  is an eigenvector of eigenvalue  $\lambda$  it must be that

$$\begin{aligned} \lambda(\alpha u_1 + \beta u_2) &= (u_1 u_1^\top - u_2 u_2^\top)(\alpha u_1 + \beta u_2) \\ &= (\alpha + \beta(u_1^\top u_2))u_1 - (\alpha(u_1^\top u_2) + \beta)u_2 \end{aligned}$$

or equivalently

$$\begin{pmatrix} (1-\lambda) & u_1^\top u_2 \\ -(u_1^\top u_2) & -(1+\lambda) \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$



By computing the determinant we see this has a solution only when

$$-(1 - \lambda^2) + (u_1^\top u_2)^2 = 0$$

Solving for  $\lambda$  then yields (7.14) and completes the proof.  $\square$

Next we show how the top eigenvectors of two spectrally similar matrices are related. We use this to bound the amount of spectral approximation we need to obtain accurate eigenvector approximations.

*Lemma 7.9.3.* Let  $\mathbf{A}$  and  $\mathbf{B}$  be symmetric PSD matrices such that  $(1 - \varepsilon)\mathbf{A} \preceq \mathbf{B} \preceq (1 + \varepsilon)\mathbf{A}$ . Then if  $g \stackrel{\text{def}}{=} \frac{\lambda_1(\mathbf{A}) - \lambda_2(\mathbf{A})}{\lambda_1(\mathbf{A})}$  satisfies  $g > 0$  we have  $[v_1(\mathbf{A})^\top v_1(\mathbf{B})]^2 \geq 1 - 2(\varepsilon/g)$ .

*Proof.* Without loss of generality  $v_1(\mathbf{B}) = \alpha v_1(\mathbf{A}) + \beta v$  for some unit vector  $v \perp v_1(\mathbf{A})$  and  $\alpha, \beta \in \mathbb{R}$  such that  $\alpha^2 + \beta^2 = 1$ . Now we know that

$$v_1(\mathbf{B})^\top \mathbf{B} v_1(\mathbf{B}) \leq (1 + \varepsilon) v_1(\mathbf{B})^\top \mathbf{A} v_1(\mathbf{B}) \leq (1 + \varepsilon) [\alpha^2 \lambda_1(\mathbf{A}) + \beta^2 \lambda_2(\mathbf{A})]$$

Furthermore, by the optimality of  $v_1(\mathbf{B})$  we have that

$$v_1(\mathbf{B})^\top \mathbf{B} v_1(\mathbf{B}) \geq (1 - \varepsilon) v_1(\mathbf{A})^\top \mathbf{A} v_1(\mathbf{A}) \geq (1 - \varepsilon) \lambda_1(\mathbf{A}).$$

Now since  $\beta^2 = 1 - \alpha^2$  combining these inequalities yields

$$(1 - \varepsilon) \lambda_1(\mathbf{A}) \leq (1 + \varepsilon) \alpha^2 (\lambda_1(\mathbf{A}) - \lambda_2(\mathbf{A})) + (1 + \varepsilon) \lambda_2(\mathbf{A}).$$

Rearranging terms, using the definition of  $g$ , and that  $g \in (0, 1]$  and  $\varepsilon \geq 0$  yields

$$\alpha^2 \geq \frac{\lambda_1(\mathbf{A}) - \lambda_2(\mathbf{A}) - \varepsilon(\lambda_1(\mathbf{A}) + \lambda_2(\mathbf{A}))}{(1 + \varepsilon)(\lambda_1(\mathbf{A}) - \lambda_2(\mathbf{A}))} = 1 - \frac{2\varepsilon \lambda_1(\mathbf{A})}{(1 + \varepsilon)(\lambda_1(\mathbf{A}) - \lambda_2(\mathbf{A}))} \geq 1 - 2(\varepsilon/g).$$

$\square$

Here we prove an approximate transitivity lemma for inner products of vectors. We use this to bound the accuracy need for certain eigenvector computations.

*Lemma 7.9.4.* Suppose that we have vectors  $v_1, v_2, v_3 \in \mathbb{R}^n$  such that  $\langle v_1, v_2 \rangle^2 \geq 1 - \varepsilon$  and  $\langle v_2, v_3 \rangle^2 \geq 1 - \varepsilon$  for  $0 < \varepsilon \leq \frac{1}{2}$  then  $\langle v_1, v_3 \rangle^2 \geq 1 - 4\varepsilon$ .

*Proof.* Without loss of generality, we can write  $v_1 = \alpha_1 v_2 + \beta_1 w_1$  for  $\alpha_1^2 + \beta_1^2 = 1$  and unit vector  $w_1 \perp v_2$ . Similarly we can write  $v_3 = \alpha_3 v_2 + \beta_3 w_3$  for  $\alpha_3^2 + \beta_3^2 = 1$  and unit vector  $w_3 \perp v_2$ . Now, by the inner products we know that  $\alpha_1^2 \geq 1 - \varepsilon$  and  $\alpha_3^2 \geq 1 - \varepsilon$  and therefore  $|\beta_1| \leq \sqrt{\varepsilon}$  and  $|\beta_3| \leq \sqrt{\varepsilon}$ . Consequently, since  $\varepsilon \leq \frac{1}{2}$ ,  $|\beta_1 \beta_3| \leq \varepsilon \leq 1 - \varepsilon \leq |\alpha_1 \alpha_3|$ , and we have

$$\begin{aligned} \langle v_1, v_3 \rangle^2 &\geq \langle \alpha_1 v_2 + \beta_1 w_1, \alpha_3 v_2 + \beta_3 w_3 \rangle^2 \geq (|\alpha_1 \alpha_3| - |\beta_1 \beta_3|)^2 \\ &\geq (1 - \varepsilon - \varepsilon)^2 = (1 - 2\varepsilon)^2 \geq 1 - 4\varepsilon. \end{aligned}$$

$\square$

## ■ 7.9.2 Convex Optimization

First we provide a single general lemma about about first order methods for convex optimization. We use this lemma for multiple purposes including bounding errors and quickly compute approximations to the central path.

*Lemma 7.9.5.* [206] Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a twice differentiable function, let  $B \subseteq \mathbb{R}^n$  be a convex set, and let  $x_*$  be a point that achieves the minimum value of  $f$  restricted to  $B$ . Further suppose that for a symmetric positive definite matrix  $\mathbf{H} \in \mathbb{R}^{n \times n}$  we have that  $\mu \mathbf{H} \preceq \nabla^2 f(y) \preceq L \mathbf{H}$  for all  $y \in B$ . Then for all  $x \in B$  we have

$$\frac{1}{2L} \|\nabla f(x)\|_{\mathbf{H}^{-1}}^2 \leq f(x) - f(x_*) \leq \frac{L}{2} \|x - x_*\|_{\mathbf{H}}^2$$

and

$$\frac{\mu}{2} \|x - x_*\|_{\mathbf{H}}^2 \leq f(x) - f(x_*) \leq \frac{1}{2\mu} \|\nabla f(x)\|_{\mathbf{H}^{-1}}^2.$$

Furthermore, if

$$x_1 = \arg \min_{x \in B} \left[ f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + \frac{L}{2} \|x_0 - x\|_{\mathbf{H}}^2 \right]$$

then

$$f(x_1) - f(x_*) \leq \left(1 - \frac{\mu}{L}\right) (f(x_0) - f(x_*)). \quad (7.15)$$

Next we provide a short technical lemma about the convexity of functions that arises naturally in our line searching procedure.

*Lemma 7.9.6.* Let  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  be a convex function and let  $g(\alpha) \stackrel{\text{def}}{=} \min_{x \in S} f(x + \alpha d)$  for any convex set  $S$  and  $d \in \mathbb{R}^n$ . Then  $g$  is convex.

*Proof.* Let  $\alpha, \beta \in \mathbb{R}$  and define  $x_\alpha = \arg \min_{x \in S} f(x + \alpha d)$  and  $x_\beta = \arg \min_{x \in S} f(x + \beta d)$ . For any  $t \in [0, 1]$  we have

$$\begin{aligned} g(t\alpha + (1-t)\beta) &= \min_{x \in S} f(x + (t\alpha + (1-t)\beta)d) \\ &\leq f(tx_\alpha + (1-t)x_\beta + (t\alpha + (1-t)\beta)d) && \text{(Convexity of } S\text{)} \\ &\leq t \cdot f(x_\alpha + \alpha d) + (1-t) \cdot f(x_\beta + \beta \cdot d) && \text{(Convexity of } f\text{)} \\ &= t \cdot g(\alpha) + (1-t) \cdot g(\beta) \end{aligned}$$

□

*Lemma 7.9.7.* For any vectors  $y, z, v \in \mathbb{R}^d$  and scalar  $\alpha$ , we can minimize  $\min_{\|x-y\|_2^2 \leq \alpha} \|x-z\|_{\mathbf{I}-vv^\top}^2$  exactly in time  $O(d)$ .

*Proof.* Let  $x^*$  be the solution of this problem. If  $\|x^* - y\|_2^2 < \alpha$ , then  $x^* = z$ . Otherwise, there is  $\lambda > 0$  such that  $x^*$  is the minimizer of

$$\min_{x \in \mathbb{R}^d} \|x - z\|_{\mathbf{I}-vv^\top}^2 + \lambda \|x - y\|_2^2.$$

Let  $\mathbf{Q} = \mathbf{I} - vv^\top$ . Then, the optimality condition of the above equation shows that

$$\mathbf{Q}(x^* - z) + \lambda(x^* - y) = \vec{0}.$$

Therefore,

$$x^* = (\mathbf{Q} + \lambda \mathbf{I})^{-1} (\mathbf{Q}z + \lambda y). \quad (7.16)$$

Hence,

$$\alpha = \|x - y\|_2^2 = (z - y)^\top \mathbf{Q} (\mathbf{Q} + \lambda \mathbf{I})^{-2} \mathbf{Q} (z - y).$$

Let  $\eta = 1 + \lambda$ , then we have  $(\mathbf{Q} + \lambda \mathbf{I}) = \eta \mathbf{I} - vv^\top$  and hence Sherman–Morrison formula shows that

$$(\mathbf{Q} + \lambda \mathbf{I})^{-1} = \eta^{-1} \mathbf{I} + \frac{\eta^{-2} vv^\top}{1 - \|v\|^2 \eta^{-1}} = \eta^{-1} \left( \mathbf{I} + \frac{vv^\top}{\eta - \|v\|^2} \right).$$

Hence, we have

$$(\mathbf{Q} + \lambda \mathbf{I})^{-2} = \eta^{-2} \left( \mathbf{I} + \frac{2vv^\top}{\eta - \|v\|^2} + \frac{vv^\top \|v\|^2}{(\eta - \|v\|^2)^2} \right) = \eta^{-2} \left( \mathbf{I} + \frac{2\eta - \|v\|^2}{(\eta - \|v\|^2)^2} vv^\top \right).$$

Let  $c_1 = \|\mathbf{Q}(z - y)\|_2^2$  and  $c_2 = (v^\top \mathbf{Q}(z - y))^2$ , then we have

$$\alpha = \eta^{-2} \left( c_1 + \frac{2\eta - \|v\|^2}{(\eta - \|v\|^2)^2} c_2 \right).$$

Hence, we have

$$\alpha \eta^2 (\eta - \|v\|^2)^2 = c_1 (\eta - \|v\|^2)^2 + c_2 (2\eta - \|v\|^2)^2.$$

Note that this is a polynomial of degree 4 in  $\eta$  and all coefficients can be computed in  $O(d)$  time. Solving this by explicit formula, one can test all 4 possible  $\eta$ 's into the formula (7.16) of  $x$ . Together with trivial case  $x^* = z$ , we simply need to check among 5 cases to check which is the solution.  $\square$

### ■ 7.9.3 Noisy One Dimensional Convex Optimization

Here we show how to minimize a one dimensional convex function giving a noisy oracle for evaluating the function. While this could possibly be done using general results on convex optimization with a membership oracle, the proof in one dimension is much simpler and we provide it here for completeness.

---

#### Algorithm 20: OneDimMinimizer( $\ell, u, \varepsilon, g, L$ )

---

**Input:** Interval  $[\ell, u] \subseteq \mathbb{R}$  and target additive error  $\varepsilon \in \mathbb{R}$

**Input:** noisy additive evaluation oracle  $g : \mathbb{R} \rightarrow \mathbb{R}$  and Lipschitz bound  $L > 0$

Let  $x = \ell$ ,  $y_\ell^{(0)} = \ell$ ,  $y_u^{(0)} = u$

**for**  $i = 1, \dots, \left\lceil \log_{3/2} \left( \frac{L(u-\ell)}{\varepsilon} \right) \right\rceil$  **do**

Let  $z_\ell^{(i)} = \frac{2y_\ell^{(i-1)} + y_u^{(i-1)}}{3}$  and  $z_u^{(i)} = \frac{y_\ell^{(i-1)} + 2y_u^{(i-1)}}{3}$

**if**  $g(z_\ell^{(i)}) \leq g(z_u^{(i)})$  **then**

Let  $(y_\ell^{(i)}, y_u^{(i)}) = (y_\ell^{(i-1)}, z_u^{(i)})$ .

If  $g(z_\ell^{(i)}) \leq g(x^{(i-1)})$  update  $x := z_\ell^{(i)}$  ..

**else if**  $g(z_\ell^{(i)}) > g(z_u^{(i)})$  **then**

Let  $(y_\ell^{(i)}, y_u^{(i)}) = (z_\ell^{(i)}, y_u^{(i-1)})$ .

If  $g(z_u^{(i)}) \leq g(x^{(i-1)})$  update  $x := z_u^{(i)}$ .

**end**

**end**

**Output:**  $x$

---

*Lemma 7.9.8.* Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be an  $L$ -Lipschitz convex function defined on the  $[\ell, u]$  interval and let  $g : \mathbb{R} \rightarrow \mathbb{R}$  be an oracle such that  $|g(y) - f(y)| \leq \varepsilon$  for all  $y$ . In  $O(\log(\frac{L(u-\ell)}{\varepsilon}))$  time and with  $O(\log(\frac{L(u-\ell)}{\varepsilon}))$  calls to  $g$ , the algorithm `OneDimMinimizer`( $\ell, u, \varepsilon, g, L$ ) outputs a point  $x$  such that

$$f(x) - \min_{y \in [\ell, u]} f(y) \leq 4\varepsilon.$$

*Proof.* First, note that for any  $y, y' \in \mathbb{R}$  if  $f(y) < f(y') - 2\varepsilon$  then  $g(y) < g(y')$ . This directly follows from our assumption on  $g$ . Second, note that the output of the algorithm,  $x$ , is simply the point

queried by the algorithm (i.e.  $\ell$  and the  $z_\ell^i$  and  $z_u^i$ ) with the smallest value of  $g$ . Combining these facts implies that  $f(x)$  is within  $2\varepsilon$  of the minimum value of  $f$  among the points queried. It thus suffices to show that the algorithm queries some point within  $2\varepsilon$  of optimal.

To do this, we break into two cases. First, consider the case where the intervals  $[y_\ell^{(i)}, y_u^{(i)}]$  all contain a minimizer of  $f$ . In this case, the final interval contains an optimum, and is of size at most  $\frac{\varepsilon}{L}$ . Thus, by the Lipschitz property, all points in the interval are within  $\varepsilon \leq 2\varepsilon$  of optimal, and at least one endpoint of the interval must have been queried by the algorithm.

For the other case, consider the last  $i$  for which this interval does contain an optimum of  $f$ . This means that  $g(z_\ell^{(i)}) \leq g(z_u^{(i)})$  while a minimizer  $x^*$  is to the right of  $z_u^{(i)}$ , or the symmetric case with a minimizer is to the left of  $z_\ell^{(i)}$ . Without loss of generality, we assume the former. We then have  $z_\ell^{(i)} \leq z_u^{(i)} \leq x^*$  and  $x^* - z_u^{(i)} \leq z_u^{(i)} - z_\ell^{(i)}$ . Consequently  $z_u^{(i)} = \alpha z_\ell^{(i)} + (1 - \alpha)x^*$  where  $\alpha \in [0, \frac{1}{2}]$  and the convexity of  $f$  implies  $f(z_u^{(i)}) \leq \frac{1}{2}f(z_\ell^{(i)}) + \frac{1}{2}f(x^*)$  or equivalently  $f(z_u^{(i)}) - f(x^*) \leq f(z_\ell^{(i)}) - f(z_u^{(i)})$ . But  $f(z_\ell^{(i)}) - f(z_u^{(i)}) \leq 2\varepsilon$  since  $g(z_\ell^{(i)}) \leq g(z_u^{(i)})$ . Thus,  $f(z_u^{(i)}) - f(x^*) \leq 2\varepsilon$ , and  $z_u^{(i)}$  is queried by the algorithm, as desired.  $\square$

## ■ 7.10 Weighted Geometric Median

In this section, we show how to extend our results to the *weighted geometric median* problem, also known as the Weber problem: given a set of  $n$  points in  $d$  dimensions,  $a^{(1)}, \dots, a^{(n)} \in \mathbb{R}^d$ , with corresponding weights  $w^{(1)}, \dots, w^{(n)} \in \mathbb{R}_{>0}$ , find a point  $x_* \in \mathbb{R}^d$  that minimizes the weighted sum of Euclidean distances to them:

$$x_* \in \arg \min_{x \in \mathbb{R}^d} f(x) \quad \text{where} \quad f(x) \stackrel{\text{def}}{=} \sum_{i \in [n]} w^{(i)} \|x - a^{(i)}\|_2.$$

As in the unweighted problem, our goal is to compute  $(1 + \varepsilon)$ -approximate solution, i.e.  $x \in \mathbb{R}^d$  with  $f(x) \leq (1 + \varepsilon)f(x_*)$ .

First, we show that it suffices to consider the case where the weights are integers with bounded sum (Lemma 7.10.1). Then, we show that such an instance of the weighted geometric median problem can be solved using the algorithms developed for the unweighted problem.

*Lemma 7.10.1.* Given points  $a^{(1)}, a^{(2)}, \dots, a^{(n)} \in \mathbb{R}^d$ , non-negative weights  $w^{(1)}, w^{(2)}, \dots, w^{(n)} \in \mathbb{R}_{>0}$ , and  $\varepsilon \in (0, 1)$ , we can compute in linear time weights  $w_1^{(1)}, w_1^{(2)}, \dots, w_1^{(n)}$  such that:

- Any  $(1 + \varepsilon/5)$ -approximate weighted geometric median of  $a^{(1)}, \dots, a^{(n)}$  with the weights  $w_1^{(1)}, \dots, w_1^{(n)}$  is also a  $(1 + \varepsilon)$ -approximate weighted geometric median of  $a^{(1)}, \dots, a^{(n)}$  with the weights  $w^{(1)}, \dots, w^{(n)}$ , and
- $w_1^{(1)}, \dots, w_1^{(n)}$  are nonnegative integers and  $\sum_{i \in [n]} w_1^{(i)} \leq 5n\varepsilon^{-1}$ .

*Proof.* Let

$$f(x) = \sum_{i \in [n]} w^{(i)} \|a^{(i)} - x\|$$

and  $W = \sum_{i \in [n]} w^{(i)}$ . Furthermore, let  $\varepsilon' = \varepsilon/5$  and for each  $i \in [n]$ , define  $w_0^{(i)} = \frac{n}{\varepsilon' W} w^{(i)}$ ,  $w_1^{(i)} = \lfloor w_0^{(i)} \rfloor$ , and  $w_2^{(i)} = w_0^{(i)} - w_1^{(i)}$ . We also define  $f_0, f_1, f_2, W_0, W_1, W_2$  analogously to  $f$  and  $W$ .

Now, assume  $f_1(x) \leq (1 + \varepsilon')f_1(x_*)$ , where  $x_*$  is the minimizer of  $f$  and  $f_0$ . Then:

$$f_0(x) = f_1(x) + f_2(x) \leq f_1(x) + f_2(x_*) + W_2 \|x - x_*\|_2$$

and

$$\begin{aligned} W_2 \|x - x_*\|_2 &= \frac{W_2}{W_1} \sum_{i \in [n]} w_1^{(i)} \|x - x_*\|_2 \leq \frac{W_2}{W_1} \sum_{i \in [n]} w_1^{(i)} \left( \|x - a^{(i)}\|_2 + \|a^{(i)} - x_*\| \right) \\ &\leq \frac{W_2}{W_1} (f_1(x) + f_1(x_*)). \end{aligned}$$

Now, since  $W_0 = \frac{n}{\varepsilon'}$  and  $W_1 \geq W_0 - n$  we have

$$\frac{W_2}{W_1} = \frac{W_0 - W_1}{W_1} = \frac{W_0}{W_1} - 1 \leq \frac{W_0}{W_0 - n} - \frac{W_0 - n}{W_0 - n} = \frac{n}{\frac{n}{\varepsilon'} - n} = \frac{\varepsilon'}{1 - \varepsilon'}.$$

Combining these yields that

$$\begin{aligned} f_0(x) &\leq f_1(x) + f_2(x_*) + \frac{\varepsilon'}{1 - \varepsilon'} (f_1(x) + f_1(x_*)) \\ &\leq \left( 1 + \frac{\varepsilon'}{1 - \varepsilon'} \right) (1 + \varepsilon') f_1(x_*) + \frac{\varepsilon'}{1 - \varepsilon'} f_1(x_*) + f_2(x_*) \\ &\leq (1 + 5\varepsilon') f_0(x_*) = (1 + \varepsilon) f_0(x_*). \end{aligned}$$

□

We now proceed to show the main result of this section.

*Lemma 7.10.2.* A  $(1 + \varepsilon)$ -approximate weighted geometric median of  $n$  points in  $\mathbb{R}^d$  can be computed in  $O(nd \log^3 \varepsilon^{-1})$  time.

*Proof.* By applying Lemma 7.10.1, we can assume that the weights are integer and their sum does not exceed  $n\varepsilon^{-1}$ . Note that computing the weighted geometric median with such weights is equivalent to computing an unweighted geometric median of  $O(n\varepsilon^{-1})$  points (where each point of the original input is repeated with the appropriate multiplicity). We now show how to simulate the behavior of our unweighted geometric median algorithms on such a set of points without computing it explicitly.

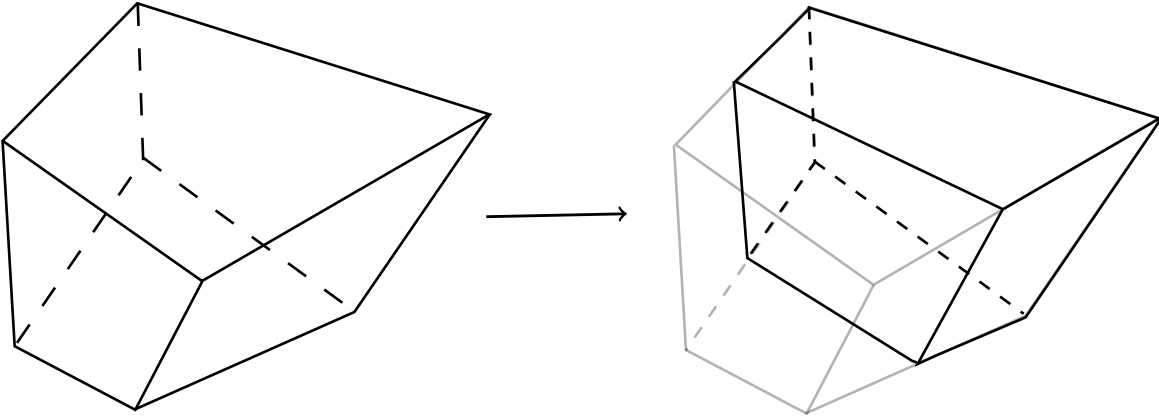
If  $\varepsilon > n^{-1/2}$ , we will apply the algorithm `ApproximateMedian`( $\varepsilon$ ), achieving a runtime of  $O(d\varepsilon^{-2}) = O(nd)$ . It is only necessary to check that we can implement weighted sampling from our points with  $O(n)$  preprocessing and  $O(1)$  time per sample. This is achieved by the alias method [154].

Now assume  $\varepsilon < n^{-1/2}$ . We will employ the algorithm `AccurateMedian`( $\varepsilon$ ). Note that the initialization using `ApproximateMedian`(2) can be performed as described in the previous paragraph. It remains to check that we can implement the subroutines `LineSearch` and `ApproxMinEig` on the implicitly represented multiset of  $O(n\varepsilon^{-1})$  points. It is enough to observe only  $n$  of the points are distinct, and all computations performed by these subroutines are identical for identical points. The total runtime will thus be  $O(nd \log^3(n/\varepsilon^2)) = O(nd \log^3 \varepsilon^{-1})$ . □

# Part II

---

# Cutting



---

---

## Convex Minimization In Nearly-Cubic Time

### ■ 8.1 Introduction

Throughout this chapter, we study the following *feasibility problem*:

**Definition 8.1.1** (Feasibility Problem). Given a separation oracle (see Definition 2.3.5) for a set  $K \subseteq \mathbb{R}^n$  contained in a box of radius  $R$  either find a point  $\mathbf{x} \in K$  or prove that  $K$  does not contain a ball of radius  $\varepsilon$ .

This feasibility problem is one of the most fundamental and classic problems in optimization. Since the celebrated result of Shor [236], Yudin and Nemirovski [274] and Khachiyan [142] in 1970s essentially proving that it can be solved in time  $O(\text{poly}(n) \cdot \text{SO} \cdot \log(R/\varepsilon))$ , this problem has served as one of the key primitives for solving numerous problems in both combinatorial and convex optimization.

Despite the prevalence of this feasibility problem, the best known running time for solving this problem has not been improved in over 25 years. In a seminal paper of Vaidya in 1989 [253], he showed how to solve the problem in  $\tilde{O}(n \cdot \text{SO} \cdot \log(nR/\varepsilon) + n^{\omega+1} \log(nR/\varepsilon))$  time. Despite interesting generalizations and practical improvements [11, 226, 101, 20, 103, 204, 272, 102], the best theoretical guarantees for solving this problem have not been improved since.

In this chapter, we show how to improve upon Vaidya’s running time in certain regimes. We provide a cutting plane algorithm which achieves an expected running time of  $O(n \cdot \text{SO} \cdot \log(nR/\varepsilon) + n^3 \log^{O(1)}(nR/\varepsilon))$ , improving upon the previous best known running time for the current known value of  $\omega < 2.373$  [269, 99] when  $R/\varepsilon = O(\text{poly}(n))$ .

We achieve our results by the combination of multiple techniques. First we show how to use techniques from the work of Vaidya and Atkinson to modify Vaidya’s scheme so that it is able to tolerate random noise in the computation in each iteration. We then show how to use known numerical machinery [257, 242, 164] in combination with some new techniques (Section 8.4.1 and Section 8.4.2) to implement each of these relaxed iterations efficiently. We hope that both these numerical techniques as well as our scheme for approximating complicated methods, such as Vaidya’s, may find further applications.

#### ■ 8.1.1 Previous Work

Throughout this chapter, we restrict our attention to algorithms for the feasibility problem that have a polynomial dependence on  $\text{SO}$ ,  $n$ , and  $\log(R/\varepsilon)$ . Such “efficient” algorithms typically follow the following iterative framework. First, they compute some trivial region  $\Omega$  that contains  $K$ . Then, they call the separation oracle at some point  $\mathbf{x} \in \Omega$ . If  $\mathbf{x} \in K$ , the algorithm successfully solved the problem. If  $\mathbf{x} \notin K$  then the separation oracle must return a half-space containing  $K$ . The algorithm then uses this half-space to shrink the region  $\Omega$  while maintaining the invariant that  $K \subseteq \Omega$ . The algorithm then repeats this process until it finds a point  $\mathbf{x} \in K$  or the region  $\Omega$  becomes too small to contain a ball with radius  $\varepsilon$ .

Year	Algorithm	Complexity
1979	Ellipsoid Method [236, 274, 142]	$O(n^2 \text{SO} \log \kappa + n^4 \log \kappa)$
1988	Inscribed Ellipsoid [144, 213]	$O(n \text{SO} \log \kappa + (n \log \kappa)^{4.5})$
1989	Volumetric Center [253]	$O(n \text{SO} \log \kappa + n^{1+\omega} \log \kappa)$
1995	Analytic Center [20]	$O\left(\begin{array}{l} n \text{SO} \log^2 \kappa + n^{\omega+1} \log^2 \kappa \\ + (n \log \kappa)^{2+\omega/2} \end{array}\right)$
2004	Random Walk [36]	$\rightarrow O(n \text{SO} \log \kappa + n^7 \log \kappa)$
-	This Chapter	$O(n \text{SO} \log \kappa + n^3 \log^{O(1)} \kappa)$

**Table 8.1:** Algorithms for the Feasibility Problem.  $\kappa$  indicates  $nR/\varepsilon$ . The arrow,  $\rightarrow$ , indicates that it solves a more general problem where only a membership oracle is given.

Previous works on efficient algorithms for the feasibility problem all follow this iterative framework. They vary in terms of what set  $\Omega$  they maintain, how they compute the center to query the separation oracle, and how they update the set. In Table 8.1, we list the previous running times for solving the feasibility problem. As usual SO indicates the cost of the separation oracle.

The first efficient algorithm for the feasibility problem is the *ellipsoid* method, due to Shor [236], Nemirovskii and Yudin [274], and Khachiyan [142]. The ellipsoid method maintains an ellipsoid as  $\Omega$  and uses the center of the ellipsoid as the next query point. It takes  $\Theta(n^2 \log(nR/\varepsilon))$  calls of oracle which is far from the lower bound  $\Omega(n \log(R/\varepsilon))$  calls [201].

To alleviate the problem, the algorithm could maintain all the information from the oracle, i.e., the polytope created from the intersection of all half-spaces obtained. The center of gravity method [168] achieves the optimal oracle complexity using this polytope and the center of gravity of this polytope as the next point. However, computing center of gravity is computationally expensive and hence we do not list its running time in Table 8.1. The Inscribed Ellipsoid Method [144] also achieved an optimal oracle complexity using this polytope as  $\Omega$  but instead using the center of the maximal inscribed ellipsoid in the polytope to query the separation oracle. We listed it as occurring in year 1988 in Table 8.1 because it was [213] that yielded the first polynomial time algorithm to actually compute this maximal inscribed ellipsoid for polytope.

Vaidya [253] obtained a faster algorithm by maintaining an approximation of this polytope and using a different center, namely the volumetric center. Although the oracle complexity of this volumetric center method is very good, the algorithm is not extremely efficient as each iteration involves matrix inversion. Atkinson and Vaidya [20] showed how to avoid this computation in certain settings. However, they were unable to achieve the desired convergence rate from their method.

Bertsimas and Vempala [36] also gives an algorithm that avoids these expensive linear algebra operations while maintaining the optimal convergence rate by using techniques in sampling convex sets. Even better, this result works for a much weaker oracle, the membership oracle. However, the additional cost of this algorithm is relatively high. We remark that while there are considerable improvements on the sampling techniques [180, 133, 164], the additional cost is still quite high compared to standard linear algebra.

### ■ 8.1.2 Challenges in Improving Previous Work

Our algorithm builds upon the previous fastest algorithm of Vaidya [253]. Ignoring implementation details and analysis, Vaidya's algorithm is quite simple. This algorithm simply maintains a polytope  $P^{(k)} = \{x \in \mathbb{R}^n : \mathbf{A}x - \mathbf{b} > \mathbf{0}\}$  as the current  $\Omega$  and uses the *volumetric center*, the minimizer of the



following *volumetric barrier function*

$$\arg \min_x \frac{1}{2} \log \det (\mathbf{A}^T \mathbf{S}_x^{-2} \mathbf{A}) \quad \text{where} \quad \mathbf{S}_x \stackrel{\text{def}}{=} \text{Diag}(\mathbf{A}\mathbf{x} - \mathbf{b}) \quad (8.1)$$

as the point at which to query the separation oracle. The polytope is then updated by adding shifts of the half-spaces returned by the separation oracle and dropping unimportant constraints. By choosing the appropriate shift, picking the right rule for dropping constraints, and using Newton's method to compute the volumetric center he achieved a running time of  $O(n \cdot SO \cdot \log \kappa + n^{1+\omega} \log \kappa)$ .

While Vaidya's algorithm's dependence on  $SO$  is essentially optimal, the additional per-iteration costs of his algorithm could possibly be improved. The computational bottleneck in each iteration of Vaidya's algorithm is computing the gradient of  $\log \det$  which in turn involves computing the leverage scores  $\boldsymbol{\sigma}(\mathbf{x}) \stackrel{\text{def}}{=} \text{diag}(\mathbf{S}_x^{-1} \mathbf{A} (\mathbf{A}^T \mathbf{S}_x^{-2} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{S}_x^{-1})$ , a key concept in part I. As the best known algorithms for computing leverage scores exactly in this setting take time  $O(n^\omega)$ , directly improving the running time of Vaidya's algorithm seems challenging.

However, similar to previous chapters, these leverage scores can be computed up to a multiplicative  $(1 \pm \varepsilon)$  error by solving  $O(\varepsilon^{-2} \log n)$  linear systems involving  $\mathbf{A}^T \mathbf{S}_x^{-2} \mathbf{A}$ . While in general this still takes time  $O(\varepsilon^{-2} n^\omega)$ , there are known techniques for efficiently maintaining the inverse of a matrix so that solving linear systems take amortized  $O(n^2)$  time [257, 164]. Consequently if it could be shown that computing *approximate* leverage scores sufficed, this would potentially decrease the amortized cost per iteration of Vaidya's method.

Unfortunately, Vaidya's method does not seem to tolerate this type of multiplicative error. If we compute compute gradients for (8.1) using this approximation of leverage scores, it seems that the point computed would be far from the true center. Moreover, without being fairly close to the true volumetric center, it is difficult to argue that such a cutting plane method would make sufficient progress.

To overcome this issue, it is tempting to directly use the linear program result in Chapter 6. In this chapter, we faced a similar issue where a volumetric, i.e.  $\log \det$ , potential function had the right analytic and geometric properties, however was computationally expensive to minimize. To overcome this issue the authors instead computed a weighted analytic center:

$$\arg \min_x - \sum_{i \in [m]} w_i \log s_i(\mathbf{x}) \quad \text{where} \quad \mathbf{s}(\mathbf{x}) \stackrel{\text{def}}{=} \mathbf{A}\mathbf{x} - \mathbf{b} \quad .$$

For carefully chosen weights this center provides the same convergence guarantees as the volumetric potential function, while each step can be computed by solving few linear systems (rather than forming the matrix inverse).

Unfortunately, it is unclear how to directly extend the work in Chapter 6 on solving an explicit linear program to the feasibility problem specified by a separation oracle. While it is possible to approximate the volumetric barrier by a weighted analytic center in many respects, proving that this approximation suffices for fast convergence remains open. In fact, the volumetric barrier function as used in Vaidya's algorithm is well approximated simply by the standard analytic center

$$\arg \min_x - \sum_{i \in [m]} \log s_i(\mathbf{x}) \quad \text{where} \quad \mathbf{s}(\mathbf{x}) \stackrel{\text{def}}{=} \mathbf{A}\mathbf{x} - \mathbf{b} \quad .$$

as all the unimportant constraints are dropped during the algorithm. However, despite decades of research, the best running times known for solving the feasibility problem using the analytic center are Vaidya and Atkinson algorithm from 1995 [20]. While the running time of this algorithm could possibly be improved using approximate leverage score computations and amortized efficient linear

system solvers, unfortunately at best, without further insight this would yield an algorithm which requires a suboptimal  $O(n \log^{O(1)} \kappa)$  queries to the separation oracle.

As pointed out in [20], the primary difficulty in using any sort of analytic center is quantifying the amount of progress made in each step. We still believe providing direct near-optimal analysis of weighted analytic center is a tantalizing open question warranting further investigation. However, rather than directly address the question of the performance of weighted analytic centers for the feasibility problem, we take a slightly different approach that side-steps this issue. We provide a partial answer that still sheds some light on the performance of the weighted analytic center while still providing our desired running time improvements.

### ■ 8.1.3 Our Approach

To overcome the shortcoming of the volumetric and analytic centers we instead consider a hybrid barrier function

$$\arg \min_{\mathbf{x}} - \sum_{i \in [m]} w_i \log s_i(\mathbf{x}) + \log \det(\mathbf{A}^T \mathbf{S}_x^{-1} \mathbf{A}) \quad \text{where} \quad \mathbf{s}(\mathbf{x}) \stackrel{\text{def}}{=} \mathbf{A}\mathbf{x} - \mathbf{b} \quad .$$

for carefully chosen weights. Our key observation is that for a correct choice of weights, we can compute the gradient of this potential function. In particular if we let  $\mathbf{w} = \boldsymbol{\tau} - \boldsymbol{\sigma}(\mathbf{x})$  then the gradient of this potential function is the same as the gradients of  $\sum_{i \in [m]} \tau_i \log s_i(\mathbf{x})$ , which we can compute efficiently. Moreover, since we are using  $\log \det$ , we can use analysis similar to Vaidya's algorithm [253] to analyze the convergence rate of this algorithm.

Unfortunately, this is a simple observation and does not immediately change the problem substantially. It simply pushes the problem of computing gradients of  $\log \det$  to computing  $\mathbf{w}$ . Therefore, for this scheme to work, we would need to ensure that the weights do not change too much and that when they change, they do not significantly hurt the progress of our algorithm. In other words, for this scheme to work, we would still need very precise estimates of leverage scores.

However, we note that the leverage scores  $\boldsymbol{\sigma}(\mathbf{x})$  do not change too much between iterations. Moreover, we provide what we believe is an interesting technical result that unbiased estimates to the changes in leverage scores can be computed using linear system solvers such that the *total error* of the estimate is bounded by the total change of the leverage scores (See Section 8.4.1). Using this result our scheme simply follows Vaidya's basic scheme in [253], however instead of minimizing the hybrid barrier function directly we alternate between taking Newton steps we can compute, changing the weights so that we can still compute Newton steps, and computing accurate unbiased estimates of the changes in the leverage scores so that the weights do not change adversarially by too much.

To make this scheme work, there are two additional details that need to be dealt with. First, we cannot let the weights vary too much as this might ultimately hurt the rate of progress of our algorithm. Therefore, in every iteration we compute a single leverage score to high precision to control the value of  $w_i$  and we show that by careful choice of the index we can ensure that no weight gets too large (See Section 8.4.2).

Second, we need to show that changing weights does not affect our progress by much more than the progress we make with respect to  $\log \det$ . To do this, we need to show the slacks are bounded above and below. We enforce this by adding regularization terms and consider the potential function

$$p_e(\mathbf{x}) = - \sum_{i \in [m]} w_i \log s_i(\mathbf{x}) + \frac{1}{2} \log \det (\mathbf{A}^T \mathbf{S}_x^{-2} \mathbf{A} + \lambda \mathbf{I}) + \frac{\lambda}{2} \|\mathbf{x}\|_2^2$$

This allows us to ensure that the entries of  $\mathbf{s}(\mathbf{x})$  do not get too large or too small and therefore

changing the weighting of the analytic center cannot affect the function value too much.

Third, we need to make sure our potential function is convex. If we simply take  $\mathbf{w} = \boldsymbol{\tau} - \boldsymbol{\sigma}(\mathbf{x})$  with  $\boldsymbol{\tau}$  as an estimator of  $\boldsymbol{\sigma}(\mathbf{x})$ ,  $\mathbf{w}$  can be negative and the potential function could be non-convex. To circumvent this issue, we use  $\mathbf{w} = c_e + \boldsymbol{\tau} - \boldsymbol{\sigma}(\mathbf{x})$  and make sure  $\|\boldsymbol{\tau} - \boldsymbol{\sigma}(\mathbf{x})\|_\infty < c_e$ .

Combining these insights, using efficient algorithms for solving a sequence of slowly changing linear systems [257, 164], and providing careful analysis ultimately allows us to achieve a running time of  $O(nSO \log \kappa + n^3 \log^{O(1)} \kappa)$  for the feasibility problem. Furthermore, in the case that  $K$  does not contain a ball of radius  $\varepsilon$ , our algorithm provides a proof that the polytope does not contain a ball of radius  $\varepsilon$ . This proof ultimately allows us to achieve running time improvements for strongly polynomial submodular minimization in Chapter 10.

### ■ 8.1.4 Organization

The rest of Chapter 8 is organized as follows. In Section 8.2 we provide some preliminary information and notation we use throughout Chapter 8. In Section 8.3 we then provide and analyze our cutting plane method. In Section 8.4 we provide key technical tools which may be of independent interest.

## ■ 8.2 Preliminaries

Here we introduce some notation and concepts we use throughout this chapter.

### ■ 8.2.1 Hybrid Barrier Function

As explained in Section 8.1.3 our cutting plane method maintains a polytope  $P = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} \geq \mathbf{b}\}$  for  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{b} \in \mathbb{R}^m$  that contains some target set  $K$ . We then maintain a minimizer of the following hybrid barrier function:

$$p_e(\mathbf{x}) \stackrel{\text{def}}{=} - \sum_{i \in [m]} (c_e + e_i) \log(s_i(\mathbf{x})/R) + \frac{1}{2} \log \det(R^2 (\mathbf{A}^T \mathbf{S}_x^{-2} \mathbf{A} + \lambda \mathbf{I})) + \frac{\lambda}{2} \|\mathbf{x}\|_2^2$$

where  $\mathbf{e} \in \mathbb{R}^m$  is a variable we maintain,  $c_e \geq 0$  and  $\lambda \geq 0$  are constants we fix later,  $\mathbf{s}(\mathbf{x}) \stackrel{\text{def}}{=} \mathbf{A}\mathbf{x} - \mathbf{b}$ ,  $\mathbf{S}_x = \mathbf{Diag}(\mathbf{s}(\mathbf{x}))$ , and  $R$  is the radius of the box containing  $P$ . When the meaning is clear from context we often use the shorthand  $\mathbf{A}_x \stackrel{\text{def}}{=} \mathbf{S}_x^{-1} \mathbf{A}$ . The  $R$  term is only used in the proof and can be ignored in the algorithm because it only affects a constant term in the potential function.

Rather than maintaining  $\mathbf{e}$  explicitly, we instead maintain a vector  $\boldsymbol{\tau} \in \mathbb{R}^m$  that approximates the leverage score

$$\boldsymbol{\psi}(\mathbf{x}) \stackrel{\text{def}}{=} \mathbf{diag} \left( \mathbf{A}_x (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1} \mathbf{A}_x^T \right) \quad .$$

and pick  $\mathbf{e}$  using the function  $e_P(\boldsymbol{\tau}, \mathbf{x}) \stackrel{\text{def}}{=} \boldsymbol{\tau} - \boldsymbol{\psi}(\mathbf{x})$ . Note that  $\boldsymbol{\psi}(\mathbf{x})$  is simply the leverage scores of certain rows of the matrix

$$\begin{bmatrix} \mathbf{A}_x \\ \sqrt{\lambda} \mathbf{I} \end{bmatrix}.$$

and therefore the usual properties of leverage scores hold, i.e.  $\psi_i(\mathbf{x}) \in (0, 1)$  and  $\|\boldsymbol{\psi}(\mathbf{x})\|_1 \leq n$ . We write  $\boldsymbol{\psi}(\mathbf{x})$  equivalently as  $\boldsymbol{\psi}_x$  or  $\boldsymbol{\psi}_P$  when we want the matrix to be clear. Furthermore, we let  $\boldsymbol{\Psi}_x \stackrel{\text{def}}{=} \mathbf{Diag}(\boldsymbol{\psi}(\mathbf{x}))$  and  $\mu(\mathbf{x}) \stackrel{\text{def}}{=} \min_i \psi_i(\mathbf{x})$ . Again, we use the subscripts of  $x$  and  $P$  interchangeably and often drop them when the meaning is clear from context.

We remark that the last term  $\frac{\lambda}{2} \|\mathbf{x}\|_2^2$  ensures that our point is always within a certain region (Lemma 8.3.18) and hence the term  $(c_e + e_i) \log s_i(\mathbf{x})_i$  never gets too large. However, this  $\ell^2$  term

changes the Hessian of the potential function and hence we need to put a  $\lambda \mathbf{I}$  term inside both the log det and the leverage score to reflect this. This is the reason why we use  $\psi$  instead of the standard leverage score.

## ■ 8.3 Our Cutting Plane Method

In this section we develop and prove the correctness of our cutting plane method. We break the presentation and proof of correctness of our cutting plane methods into multiple parts. First in Section 8.3.1 we describe how we maintain a center of the hybrid barrier function  $p_e$  and analyze this procedure. Then, in Section 8.3.2 we carefully analyze the effect of changing constraints on the hybrid barrier function and in Section 8.3.3 we prove properties of an approximate center of hybrid barrier function, which we call a hybrid center. In Section 8.3.4 we then provide our cutting plane method and in Section 8.3.5 we prove that the cutting plane method solves the feasibility problem as desired.

### ■ 8.3.1 Centering

In this section we show how to compute approximate *centers* or minimizers of the hybrid barrier function for the current polytope  $P = \{\mathbf{x} : \mathbf{A}\mathbf{x} \geq \mathbf{b}\}$ . We split this proof up into multiple parts. First we simply bound the gradient and Hessian of the hybrid barrier function,  $p_e$ , as follows.

**Lemma 8.3.1.** For  $f(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{2} \log \det (\mathbf{A}^T \mathbf{S}_x^{-2} \mathbf{A} + \lambda \mathbf{I})$ , we have that

$$\nabla f(\mathbf{x}) = -\mathbf{A}_x^T \psi(\mathbf{x}) \quad \text{and} \quad \mathbf{A}_x^T \Psi(\mathbf{x}) \mathbf{A}_x \preceq \nabla^2 f(\mathbf{x}) \preceq 3\mathbf{A}_x^T \Psi(\mathbf{x}) \mathbf{A}_x \quad .$$

*Proof.* Our proof is similar to [10, Appendix] which proved the statement when  $\lambda = 0$ . This case does not change the derivation significantly, however for completeness we include the proof below.

We take derivatives on  $\mathbf{s}$  first and then apply chain rule. Let  $f(\mathbf{s}) = \frac{1}{2} \log \det (\mathbf{A}^T \mathbf{S}^{-2} \mathbf{A} + \lambda \mathbf{I})$ . We use the notation  $Df(\mathbf{x})[\mathbf{h}]$  to denote the directional derivative of  $f$  along the direction  $\mathbf{h}$  at the point  $\mathbf{x}$ . Using the standard formula for the derivative of log det, i.e.  $\frac{d}{dt} \log \det \mathbf{B}_t = \text{Tr}((\mathbf{B}_t)^{-1}(\frac{d\mathbf{B}_t}{dt}))$ , we have

$$\begin{aligned} Df(\mathbf{s})[\mathbf{h}] &= \frac{1}{2} \text{Tr}((\mathbf{A}^T \mathbf{S}^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} (\mathbf{A}^T (-2) \mathbf{S}^{-3} \mathbf{H} \mathbf{A})) \\ &= -\sum_i \frac{h_i}{s_i} \mathbf{1}_i^T \mathbf{S}^{-1} \mathbf{A} (\mathbf{A}^T \mathbf{S}^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A} \mathbf{S}^{-1} \mathbf{1}_i = -\sum_i \frac{\psi_i h_i}{s_i} \quad . \end{aligned} \quad (8.2)$$

Applying chain rules, we have  $\nabla f(\mathbf{x}) = -\mathbf{A}_x^T \psi$ . Now let  $\mathbf{P} \stackrel{\text{def}}{=} \mathbf{S}^{-1} \mathbf{A} (\mathbf{A}^T \mathbf{S}^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{S}^{-1}$ . Taking the derivative of (8.2) again and using the cyclic property of trace, we have

$$\begin{aligned} D^2 f(\mathbf{s})[\mathbf{h}_1, \mathbf{h}_2] &= \text{Tr} \left( (\mathbf{A}^T \mathbf{S}^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} (\mathbf{A}^T (-2) \mathbf{S}^{-3} \mathbf{H}_2 \mathbf{A}) (\mathbf{A}^T \mathbf{S}^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} (\mathbf{A}^T \mathbf{S}^{-3} \mathbf{H}_1 \mathbf{A}) \right) \\ &\quad - \text{Tr} \left( (\mathbf{A}^T \mathbf{S}^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} (\mathbf{A}^T (-3) \mathbf{S}^{-4} \mathbf{H}_2 \mathbf{H}_1 \mathbf{A}) \right) \\ &= 3 \text{Tr} (\mathbf{P} \mathbf{S}^{-2} \mathbf{H}_2 \mathbf{H}_1) - 2 \text{Tr} (\mathbf{P} \mathbf{S}^{-1} \mathbf{H}_2 \mathbf{P} \mathbf{S}^{-1} \mathbf{H}_1) \\ &= 3 \sum_i P_{ii} \frac{\mathbf{h}_1(i) \mathbf{h}_2(i)}{s_i^2} - 2 \sum_{ij} P_{ij} \frac{\mathbf{h}_2(j)}{s_j} P_{ji} \frac{\mathbf{h}_2(i)}{s_i} \\ &= 3 \sum_i \psi_i \frac{\mathbf{h}_1(i) \mathbf{h}_2(i)}{s_i^2} - 2 \sum_{ij} P_{ij}^2 \frac{\mathbf{h}_2(j)}{s_j} \frac{\mathbf{h}_2(i)}{s_i} \quad . \end{aligned}$$

Consequently,  $D^2f(\mathbf{x})[\mathbf{1}_i, \mathbf{1}_j] = [\mathbf{S}^{-1} (3\Psi - 2\mathbf{P}^{(2)}) \mathbf{S}^{-1}]_{ij}$  where  $\mathbf{P}^{(2)}$  is the Schur product of  $\mathbf{P}$  with itself.

Now note that

$$\begin{aligned} \sum_i P_{ij}^2 &= \mathbf{1}_j \mathbf{S}^{-1} \mathbf{A} (\mathbf{A}^T \mathbf{S}^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{S}^{-2} \mathbf{A} (\mathbf{A}^T \mathbf{S}^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{S}^{-1} \mathbf{1}_j \\ &\leq \mathbf{1}_j \mathbf{S}^{-1} \mathbf{A} (\mathbf{A}^T \mathbf{S}^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{S}^{-1} \mathbf{1}_j = P_{jj} = \Psi_{jj} \quad . \end{aligned}$$

Hence, the Gershgorin circle theorem shows that the eigenvalues of  $\Psi - \mathbf{P}^{(2)}$  are lies in union of the interval  $[0, 2\psi_j]$  over all  $j$ . Hence,  $\Psi - \mathbf{P}^{(2)} \succeq \mathbf{0}$ . On the other hand, Schur product theorem shows that  $\mathbf{P}^{(2)} \succeq \mathbf{0}$  as  $\mathbf{P} \succeq \mathbf{0}$ . Hence, the result follows by chain rule.  $\square$

Lemma 8.3.1 immediately shows that under our choice of  $\mathbf{e} = \mathbf{e}_P(\mathbf{x}, \boldsymbol{\tau})$  we can compute the gradient of the hybrid barrier function,  $p_e(\mathbf{x})$  efficiently. Formally, Lemma 8.3.1 immediately implies the following:

**Lemma 8.3.2** (Gradient). *For  $\mathbf{x} \in P = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{A}\mathbf{y} \geq \mathbf{b}\}$  and  $\mathbf{e} \in \mathbb{R}^m$  we have*

$$\nabla p_e(\mathbf{x}) = -\mathbf{A}_x^T (c_e \mathbf{1} + \mathbf{e} + \boldsymbol{\psi}_P(\mathbf{x})) + \lambda \mathbf{x}$$

and therefore for all  $\boldsymbol{\tau} \in \mathbb{R}^m$ , we have

$$\nabla p_{e(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x}) = -\mathbf{A}_x^T (c_e \mathbf{1} + \boldsymbol{\tau}) + \lambda \mathbf{x}.$$

*Remark 8.3.3.* To be clear, the vector  $\nabla p_{e(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x})$  is defined as the vector such that

$$[\nabla p_{e(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x})]_i = \lim_{t \rightarrow 0} \frac{1}{t} (p_{e(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x} + t\mathbf{1}_i) - p_{e(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x})) \quad .$$

In other words, we treat the parameter  $\mathbf{e}(\boldsymbol{\tau}, \mathbf{x})$  as fixed. This is the reason we denote it by subscript to emphasize that  $p_e(\mathbf{x})$  is a family of functions,  $p_{e(\boldsymbol{\tau}, \mathbf{x})}$  is one particular function, and  $\nabla p_{e(\boldsymbol{\tau}, \mathbf{x})}$  means taking gradient on that particular function.

Consequently, we can always compute  $\nabla p_{e(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x})$  efficiently. Now, we measure *centrality* or how close we are to the hybrid center as follows.

**Definition 8.3.4** (Centrality). For  $\mathbf{x} \in P = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{A}\mathbf{y} \geq \mathbf{b}\}$  and  $\mathbf{e} \in \mathbb{R}^m$ , we define the *centrality* of  $\mathbf{x}$  by

$$\delta_e(\mathbf{x}) \stackrel{\text{def}}{=} \|\nabla p_e(\mathbf{x})\|_{\mathbf{H}(\mathbf{x})^{-1}}$$

where  $\mathbf{H}(\mathbf{x}) \stackrel{\text{def}}{=} \mathbf{A}_x^T (c_e \mathbf{I} + \Psi(\mathbf{x})) \mathbf{A}_x + \lambda \mathbf{I}$ . Often, we use *weights*  $\mathbf{w} \in \mathbb{R}_{>0}^m$  to approximate this Hessian and consider  $\mathbf{Q}(\mathbf{x}, \mathbf{w}) \stackrel{\text{def}}{=} \mathbf{A}_x^T (c_e \mathbf{I} + \mathbf{W}) \mathbf{A}_x + \lambda \mathbf{I}$ .

Next, we bound how much slacks can change in a region close to a nearly central point.

**Lemma 8.3.5.** *Let  $\mathbf{x} \in P = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{A}\mathbf{y} \geq \mathbf{b}\}$  and  $\mathbf{y} \in \mathbb{R}^n$  such that  $\|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})} \leq \varepsilon \sqrt{c_e + \mu(\mathbf{x})}$  for  $\varepsilon < 1$ . Then  $\mathbf{y} \in P$  and  $(1 - \varepsilon)\mathbf{S}_x \preceq \mathbf{S}_y \preceq (1 + \varepsilon)\mathbf{S}_x$ .*

*Proof.* Direct calculation reveals the following:

$$\begin{aligned} \|\mathbf{S}_x^{-1}(\mathbf{s}_y - \mathbf{s}_x)\|_\infty &\leq \|\mathbf{A}_x(\mathbf{y} - \mathbf{x})\|_2 \leq \frac{1}{\sqrt{c_e + \mu(\mathbf{x})}} \|\mathbf{A}_x(\mathbf{y} - \mathbf{x})\|_{c_e \mathbf{I} + \Psi(\mathbf{x})} \\ &\leq \frac{1}{\sqrt{c_e + \mu(\mathbf{x})}} \|\mathbf{y} - \mathbf{x}\|_{\mathbf{H}(\mathbf{x})} \leq \varepsilon \quad . \end{aligned}$$

Consequently,  $(1 - \varepsilon)\mathbf{S}_x \preceq \mathbf{S}_y \preceq (1 + \varepsilon)\mathbf{S}_x$ . Since  $\mathbf{y} \in P$  if and only if  $\mathbf{S}_y \succeq \mathbf{0}$  the result follows.  $\square$

Combining the previous lemmas we obtain the following.

**Lemma 8.3.6.** *Let  $\mathbf{x} \in P = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{A}\mathbf{y} \geq \mathbf{b}\}$  and  $\mathbf{e}, \mathbf{w} \in \mathbb{R}^m$  such that  $\|\mathbf{e}\|_\infty \leq \frac{1}{2}c_e \leq 1$  and  $\Psi(\mathbf{x}) \preceq \mathbf{W} \preceq \frac{4}{3}\Psi(\mathbf{x})$ . If  $\mathbf{y} \in \mathbb{R}^n$  satisfies  $\|\mathbf{x} - \mathbf{y}\|_{\mathbf{Q}(\mathbf{x}, \mathbf{w})} \leq \frac{1}{10}\sqrt{c_e + \mu(\mathbf{x})}$ , then*

$$\frac{1}{4}\mathbf{Q}(\mathbf{x}, \mathbf{w}) \preceq \nabla^2 p_e(\mathbf{y}) \preceq 8\mathbf{Q}(\mathbf{x}, \mathbf{w}) \quad \text{and} \quad \frac{1}{2}\mathbf{H}(\mathbf{x}) \preceq \mathbf{H}(\mathbf{y}) \preceq 2\mathbf{H}(\mathbf{x}) \quad .$$

Also, we have that  $\mathbf{H}(\mathbf{x}) \preceq \mathbf{Q}(\mathbf{x}, \mathbf{w}) \preceq \frac{4}{3}\mathbf{H}(\mathbf{x})$ .

*Proof.* Lemma 8.3.1 shows that

$$\mathbf{A}_y^T (c_e \mathbf{I} + \mathbf{E} + \Psi(\mathbf{y})) \mathbf{A}_y + \lambda \mathbf{I} \preceq \nabla^2 p_e(\mathbf{y}) \preceq \mathbf{A}_y^T (c_e \mathbf{I} + \mathbf{E} + 3\Psi(\mathbf{y})) \mathbf{A}_y + \lambda \mathbf{I} \quad . \quad (8.3)$$

Since  $\mathbf{W} \succeq \Psi$ , we have that  $\mathbf{Q}(\mathbf{x}, \mathbf{w}) \succeq \mathbf{H}(\mathbf{x})$  and therefore  $\|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})} \leq \varepsilon \sqrt{c_e + \mu(\mathbf{x})}$  with  $\varepsilon = 0.1$ . Consequently, by Lemma 8.3.5 we have  $(1 - \varepsilon)\mathbf{S}_x \preceq \mathbf{S}_y \preceq (1 + \varepsilon)\mathbf{S}_x$  and therefore

$$\frac{(1 - \varepsilon)^2}{(1 + \varepsilon)^2} \Psi(\mathbf{x}) \preceq \Psi(\mathbf{y}) \preceq \frac{(1 + \varepsilon)^2}{(1 - \varepsilon)^2} \Psi(\mathbf{x})$$

and

$$\frac{1}{2}\mathbf{H}(\mathbf{x}) \preceq \frac{(1 - \varepsilon)^2}{(1 + \varepsilon)^4} \mathbf{H}(\mathbf{x}) \preceq \mathbf{H}(\mathbf{y}) \preceq \frac{(1 + \varepsilon)^2}{(1 - \varepsilon)^4} \mathbf{H}(\mathbf{x}) \preceq 2\mathbf{H}(\mathbf{x})$$

Furthermore, (8.3) shows that

$$\begin{aligned} \nabla^2 p_e(\mathbf{y}) &\preceq \mathbf{A}_y^T (c_e \mathbf{I} + \mathbf{E} + 3\Psi(\mathbf{y})) \mathbf{A}_y + \lambda \mathbf{I} \\ &\preceq \frac{(1 + \varepsilon)^2}{(1 - \varepsilon)^4} \mathbf{A}_x^T (c_e \mathbf{I} + \mathbf{E} + 3\Psi(\mathbf{x})) \mathbf{A}_x + \lambda \mathbf{I} \\ &\preceq 2\mathbf{A}_x^T \left( \frac{3}{2}c_e \mathbf{I} + 3\mathbf{W} \right) \mathbf{A}_x + \lambda \mathbf{I} \preceq 8\mathbf{Q}(\mathbf{x}, \mathbf{w}) \end{aligned}$$

and

$$\begin{aligned} \nabla^2 p_e(\mathbf{y}) &\succeq \mathbf{A}_y^T (c_e \mathbf{I} + \mathbf{E} + \Psi(\mathbf{y})) \mathbf{A}_y + \lambda \mathbf{I} \\ &\succeq \frac{(1 - \varepsilon)^4}{(1 + \varepsilon)^2} \mathbf{A}_x^T (c_e \mathbf{I} + \mathbf{E} + \Psi(\mathbf{x})) \mathbf{A}_x + \lambda \mathbf{I} \\ &\succeq \frac{1}{2} \mathbf{A}_x^T \left( \frac{1}{2}c_e \mathbf{I} + \frac{3}{4}\mathbf{W} \right) \mathbf{A}_x + \lambda \mathbf{I} \succeq \frac{1}{4}\mathbf{Q}(\mathbf{x}, \mathbf{w}). \end{aligned}$$

The last inequality follows from the definition of  $\mathbf{H}(\mathbf{x})$  and  $\mathbf{Q}(\mathbf{x}, \mathbf{w})$  and the fact  $\Psi(\mathbf{x}) \preceq \mathbf{W} \preceq \frac{4}{3}\Psi(\mathbf{x})$ .  $\square$

To analyze our centering scheme we use standard facts about gradient descent we prove in Lemma 8.3.7.

**Lemma 8.3.7** (Gradient Descent). *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be twice differentiable and  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  be positive definite. Let  $\mathbf{x}_0 \in \mathbb{R}^n$  and  $\mathbf{x}_1 \stackrel{\text{def}}{=} \mathbf{x}_0 - \frac{1}{L}\mathbf{Q}^{-1}\nabla f(\mathbf{x}_0)$ . Furthermore, let  $\mathbf{x}_\alpha = \mathbf{x}_0 + \alpha(\mathbf{x}_1 - \mathbf{x}_0)$  and suppose that  $\mu\mathbf{Q} \preceq \nabla^2 f(\mathbf{x}_\alpha) \preceq L\mathbf{Q}$  for all  $\alpha \in [0, 1]$ . Then,*

1.  $\|\nabla f(\mathbf{x}_1)\|_{\mathbf{Q}^{-1}} \leq \left(1 - \frac{\mu}{L}\right) \|\nabla f(\mathbf{x}_0)\|_{\mathbf{Q}^{-1}}$
2.  $f(\mathbf{x}_1) \geq f(\mathbf{x}_0) - \frac{1}{L} \|\nabla f(\mathbf{x}_0)\|_{\mathbf{Q}^{-1}}^2$

*Proof.* Integrating we have that

$$\nabla f(\mathbf{x}_1) = \nabla f(\mathbf{x}_0) + \int_0^1 \nabla^2 f(\mathbf{x}_\alpha)(\mathbf{x}_1 - \mathbf{x}_0) d\alpha = \int_0^1 \left( \mathbf{Q} - \frac{1}{L} \nabla^2 f(\mathbf{x}_\alpha) \right) \mathbf{Q}^{-1} \nabla f(\mathbf{x}_0) d\alpha$$

Consequently, by applying Jensen's inequality we have

$$\begin{aligned} \|\nabla f(\mathbf{x}_1)\|_{\mathbf{Q}^{-1}} &= \left\| \int_0^1 \left( \mathbf{Q} - \frac{1}{L} \nabla^2 f(\mathbf{x}_\alpha) \right) \mathbf{Q}^{-1} \nabla f(\mathbf{x}_0) d\alpha \right\|_{\mathbf{Q}^{-1}} \\ &\leq \int_0^1 \left\| \left( \mathbf{Q} - \frac{1}{L} \nabla^2 f(\mathbf{x}_\alpha) \right) \mathbf{Q}^{-1} \nabla f(\mathbf{x}_0) \right\|_{\mathbf{Q}^{-1}} d\alpha \\ &\leq \|\mathbf{Q}^{-1/2} \nabla f(\mathbf{x}_0)\|_{[\mathbf{Q}^{-1/2} (\mathbf{Q} - \frac{1}{L} \nabla^2 f(\mathbf{x}_\alpha)) \mathbf{Q}^{-1/2}]^2} \end{aligned}$$

Now we know that by assumption that

$$\mathbf{0} \preceq \mathbf{Q}^{-1/2} \left( \mathbf{Q} - \frac{1}{L} \nabla^2 f(\mathbf{x}_\alpha) \right) \mathbf{Q}^{-1/2} \preceq \left( 1 - \frac{\mu}{L} \right) \mathbf{I}$$

and therefore combining these (1) holds.

Using the convexity of  $f$ , we have

$$\begin{aligned} f(\mathbf{x}_1) &\geq f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x}_1 - \mathbf{x}_0 \rangle \\ &\geq f(\mathbf{x}_0) - \|\nabla f(\mathbf{x}_0)\|_{\mathbf{Q}^{-1}} \|\mathbf{x}_1 - \mathbf{x}_0\|_{\mathbf{Q}} \end{aligned}$$

and since  $\|\mathbf{x}_1 - \mathbf{x}_0\|_{\mathbf{Q}} = \frac{1}{L} \|\nabla f(\mathbf{x}_0)\|_{\mathbf{Q}^{-1}}$ , (2) holds as well.  $\square$

Next we bound the effect of changing  $\mathbf{e}$  on the hybrid barrier function  $p_e(\mathbf{x})$ .

**Lemma 8.3.8.** For  $\mathbf{x} \in P = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{A}\mathbf{y} \geq \mathbf{b}\}$ ,  $\mathbf{e}, \mathbf{f} \in \mathbb{R}^m$ , and  $\mathbf{w} \in \mathbb{R}_{>0}^m$  such that  $\mathbf{W} \succeq \Psi(\mathbf{x})$

$$\|\nabla p_{\mathbf{f}}(\mathbf{x})\|_{\mathbf{Q}(\mathbf{x}, \mathbf{w})^{-1}} \leq \|\nabla p_{\mathbf{e}}(\mathbf{x})\|_{\mathbf{Q}(\mathbf{x}, \mathbf{w})^{-1}} + \frac{1}{\sqrt{c_e + \mu(\mathbf{x})}} \|\mathbf{f} - \mathbf{e}\|_2$$

*Proof.* Direct calculation shows the following

$$\begin{aligned} \|\nabla p_{\mathbf{f}}(\mathbf{x})\|_{\mathbf{Q}(\mathbf{x}, \mathbf{w})^{-1}} &= \left\| -\mathbf{A}_x^T (c_e \mathbf{1} + \mathbf{f} + \boldsymbol{\psi}_P(\mathbf{x})) + \lambda \mathbf{x} \right\|_{\mathbf{Q}(\mathbf{x}, \mathbf{w})^{-1}} && \text{(Formula for } \nabla p_{\mathbf{f}}(\mathbf{x}) \text{)} \\ &\leq \|\nabla p_{\mathbf{e}}(\mathbf{x})\|_{\mathbf{Q}(\mathbf{x}, \mathbf{w})^{-1}} + \|\mathbf{A}_x^T (\mathbf{f} - \mathbf{e})\|_{\mathbf{Q}(\mathbf{x}, \mathbf{w})^{-1}} && \text{(Triangle inequality)} \\ &\leq \|\nabla p_{\mathbf{e}}(\mathbf{x})\|_{\mathbf{Q}(\mathbf{x}, \mathbf{w})^{-1}} + \frac{1}{\sqrt{c_e + \mu(\mathbf{x})}} \|\mathbf{A}_x^T (\mathbf{f} - \mathbf{e})\|_{(\mathbf{A}_x^T \mathbf{A}_x)^{-1}} && \text{(Bound on } \mathbf{Q}(\mathbf{x}, \mathbf{w}) \text{)} \\ &\leq \|\nabla p_{\mathbf{e}}(\mathbf{x})\|_{\mathbf{Q}(\mathbf{x}, \mathbf{w})^{-1}} + \frac{1}{\sqrt{c_e + \mu(\mathbf{x})}} \|\mathbf{f} - \mathbf{e}\|_2 && \text{(Property of projection matrix)} \end{aligned}$$

where in the second to third line we used  $\mathbf{Q}(\mathbf{x}, \mathbf{w}) \succeq \mathbf{H}(\mathbf{x}) \succeq (c_e + \mu(\mathbf{x})) \mathbf{A}_x^T \mathbf{A}_x$ .  $\square$

We now have everything we need to analyze our centering algorithm.

**Lemma 8.3.9.** Let  $\mathbf{x}^{(0)} \in P = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{A}\mathbf{y} \geq \mathbf{b}\}$  and let  $\boldsymbol{\tau}^{(0)} \in \mathbb{R}^m$  such that  $\|\mathbf{e}(\boldsymbol{\tau}^{(0)}, \mathbf{x}^{(0)})\|_\infty \leq \frac{1}{3} c_e \leq \frac{1}{3}$ . Assume that  $r$  is a positive integer,  $0 \leq c_\Delta \leq 0.01 c_e$  and  $\delta_{\mathbf{e}^{(0)}}(\mathbf{x}^{(0)}) \leq \frac{1}{100} \sqrt{c_e + \mu(\mathbf{x}^{(0)})}$ . With high probability in  $n$ , the algorithm `Centering`( $\mathbf{x}^{(0)}, \boldsymbol{\tau}^{(0)}, r, c_\Delta$ ) outputs  $(\mathbf{x}^{(r)}, \boldsymbol{\tau}^{(r)})$  such that

1.  $\delta_{\mathbf{e}^{(r)}}(\mathbf{x}^{(r)}) \leq 2 \left(1 - \frac{1}{64}\right)^r \delta_{\mathbf{e}^{(0)}}(\mathbf{x}^{(0)})$ .
2.  $\mathbf{E}[p_{\mathbf{e}^{(k)}}(\mathbf{x}^{(r)})] \geq p_{\mathbf{e}^{(0)}}(\mathbf{x}^{(0)}) - 8 (\delta_{\mathbf{e}^{(0)}}(\mathbf{x}^{(0)}))^2$ .

---

**Algorithm 21:**  $(\mathbf{x}^{(r)}, \boldsymbol{\tau}^{(r)}) = \text{Centering}(\mathbf{x}^{(0)}, \boldsymbol{\tau}^{(0)}, r, c_\Delta)$ 


---

**Input:** Initial point  $\mathbf{x}^{(0)} \in P = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{A}\mathbf{y} \geq \mathbf{b}\}$ , Estimator of leverage scores  $\boldsymbol{\tau}^{(0)} \in \mathbb{R}^n$ 
**Input:** Number of iterations  $r > 0$ , Accuracy of the estimator  $0 \leq c_\Delta \leq 0.01c_e$ .

**Given:**  $\|\mathbf{e}^{(0)}\|_\infty \leq \frac{1}{3}c_e \leq \frac{1}{3}$  where  $\mathbf{e}^{(0)} = \mathbf{e}(\boldsymbol{\tau}^{(0)}, \mathbf{x}^{(0)})$ .

**Given:**  $\delta_{e^{(0)}}(\mathbf{x}^{(0)}) = \|\nabla p_{e^{(0)}}(\mathbf{x}^{(0)})\|_{\mathbf{H}(\mathbf{x}^{(0)})^{-1}} \leq \frac{1}{100}\sqrt{c_e + \mu(\mathbf{x}^{(0)})}$ .

Compute  $\mathbf{w}$  such that  $\boldsymbol{\Psi}(\mathbf{x}^{(0)}) \preceq \mathbf{W} \preceq \frac{4}{3}\boldsymbol{\Psi}(\mathbf{x}^{(0)})$  (Similar to Lemma 2.3.4)

Let  $\mathbf{Q} \stackrel{\text{def}}{=} \mathbf{Q}(\mathbf{x}^{(0)}, \mathbf{w})$ .

**for**  $k = 1$  *to*  $r$  **do**

$$\mathbf{x}^{(k)} := \mathbf{x}^{(k-1)} - \frac{1}{8}\mathbf{Q}^{-1}\nabla p_{e^{(k-1)}}(\mathbf{x}^{(k-1)}).$$

Sample  $\boldsymbol{\Delta}^{(k)} \in \mathbb{R}^n$  s.t.

$$\mathbf{E}[\boldsymbol{\Delta}^{(k)}] = \boldsymbol{\psi}(\mathbf{x}^{(k)}) - \boldsymbol{\psi}(\mathbf{x}^{(k-1)}) \text{ and}$$

with high probability in  $n$ ,

$$\|\boldsymbol{\Delta}^{(k)} - (\boldsymbol{\psi}(\mathbf{x}^{(k)}) - \boldsymbol{\psi}(\mathbf{x}^{(k-1)}))\|_2 \leq c_\Delta \|\mathbf{S}_{\mathbf{x}^{(k-1)}}^{-1}(\mathbf{s}_{\mathbf{x}^{(k)}} - \mathbf{s}_{\mathbf{x}^{(k-1)}})\|_2 \text{ (See Section 8.4.1)}$$

$$\boldsymbol{\tau}^{(k)} := \boldsymbol{\tau}^{(k-1)} + \boldsymbol{\Delta}^{(k)}.$$

$$\mathbf{e}^{(k)} := \mathbf{e}(\boldsymbol{\tau}^{(k)}, \mathbf{x}^{(k)}).$$

**end**
**Output:**  $(\mathbf{x}^{(r)}, \boldsymbol{\tau}^{(r)})$ 


---

$$3. \mathbf{E}\mathbf{e}^{(r)} = \mathbf{e}^{(0)} \text{ and } \|\mathbf{e}^{(r)} - \mathbf{e}^{(0)}\|_2 \leq \frac{1}{10}c_\Delta.$$

$$4. \left\| \mathbf{S}_{\mathbf{x}^{(0)}}^{-1}(\mathbf{s}(\mathbf{x}^{(r)}) - \mathbf{s}(\mathbf{x}^{(0)})) \right\|_2 \leq \frac{1}{10}.$$

where  $\mathbf{e}^{(r)} = \mathbf{e}(\boldsymbol{\tau}^{(r)}, \mathbf{x}^{(r)})$ .

*Proof.* Let  $\eta = \|\nabla p_{e^{(0)}}(\mathbf{x}^{(0)})\|_{\mathbf{Q}^{-1}}$ . First, we use induction to prove that  $\|\mathbf{x}^{(r)} - \mathbf{x}^{(0)}\|_{\mathbf{Q}} \leq 8\eta$ ,  $\|\nabla p_{e^{(r)}}(\mathbf{x}^{(r)})\|_{\mathbf{Q}^{-1}} \leq (1 - \frac{1}{64})^r \eta$  and  $\|\mathbf{e}^{(r)} - \mathbf{e}^{(0)}\|_2 \leq \frac{1}{10}c_\Delta$  for all  $r$ .

Clearly the claims hold for  $r = 0$ . We now suppose they hold for all  $r \leq t$  and show that they hold for  $r = t + 1$ . Now, since  $\|\mathbf{x}^{(t)} - \mathbf{x}^{(0)}\|_{\mathbf{Q}} \leq 8\eta$ ,  $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \frac{1}{8}\mathbf{Q}^{-1}\nabla p_{e^{(t)}}(\mathbf{x}^{(t)})$ , and  $\|\nabla p_{e^{(t)}}(\mathbf{x}^{(t)})\|_{\mathbf{Q}^{-1}} \leq (1 - \frac{1}{64})^t \eta \leq \eta$ , we have

$$\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(0)}\|_{\mathbf{Q}} \leq \|\mathbf{x}^{(t)} - \mathbf{x}^{(0)}\|_{\mathbf{Q}} + \frac{1}{8}\|\nabla p_{e^{(t)}}(\mathbf{x}^{(t)})\|_{\mathbf{Q}^{-1}} \leq 9\eta.$$

We will improve this estimate later in the proof to finish the induction on  $\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(0)}\|_{\mathbf{Q}}$ , but using this,  $\eta \leq 0.01\sqrt{c_e + \mu(\mathbf{x}^{(0)})}$ , and  $\|\mathbf{e}^{(t)}\|_\infty \leq \|\mathbf{e}^{(t)} - \mathbf{e}^{(0)}\|_\infty + \|\mathbf{e}^{(0)}\|_\infty \leq \frac{c_e}{2}$ , we can invoke Lemma 8.3.6 and Lemma 8.3.7 and therefore

$$\|\nabla p_{e^{(t)}}(\mathbf{x}^{(t+1)})\|_{\mathbf{Q}^{-1}} \leq \left(1 - \frac{1}{32}\right) \|\nabla p_{e^{(t)}}(\mathbf{x}^{(t)})\|_{\mathbf{Q}^{-1}}.$$

By Lemma 8.3.8 we have

$$\|\nabla p_{e^{(t+1)}}(\mathbf{x}^{(t+1)})\|_{\mathbf{Q}^{-1}} \leq \left(1 - \frac{1}{32}\right) \|\nabla p_{e^{(t)}}(\mathbf{x}^{(t)})\|_{\mathbf{Q}^{-1}} + \frac{1}{\sqrt{c_e + \mu(\mathbf{x}^{(0)})}} \|\mathbf{e}^{(t+1)} - \mathbf{e}^{(t)}\|_2. \quad (8.4)$$



To bound  $\|\mathbf{e}^{(t+1)} - \mathbf{e}^{(t)}\|_2$ , we use the definition of  $\mathbf{\Delta}$  to shows that

$$\begin{aligned} \|\mathbf{e}^{(t+1)} - \mathbf{e}^{(t)}\|_2 &= \left\| \left( \boldsymbol{\tau}^{(t+1)} - \boldsymbol{\psi}(\mathbf{x}^{(t+1)}) \right) - \left( \boldsymbol{\tau}^{(t)} - \boldsymbol{\psi}(\mathbf{x}^{(t)}) \right) \right\|_2 \\ &= \left\| \mathbf{\Delta}^{(t+1)} - \left( \boldsymbol{\psi}(\mathbf{x}^{(t+1)}) - \boldsymbol{\psi}(\mathbf{x}^{(t)}) \right) \right\|_2 \\ &\leq c_{\Delta} \left\| \mathbf{S}_{\mathbf{x}^{(t)}}^{-1} (\mathbf{s}_{\mathbf{x}^{(t+1)}} - \mathbf{s}_{\mathbf{x}^{(t)}}) \right\|_2 \end{aligned}$$

with high probability in  $n$ . Now, we note that Lemma 8.3.5 and the induction hypothesis  $\|\mathbf{x}^{(t)} - \mathbf{x}^{(0)}\|_{\mathbf{H}(\mathbf{x}^{(0)})} \leq \|\mathbf{x}^{(t)} - \mathbf{x}^{(0)}\|_{\mathbf{Q}} \leq 8\eta$  shows that  $(1 - 0.1)\mathbf{S}_{\mathbf{x}^{(0)}} \preceq \mathbf{S}_{\mathbf{x}^{(t)}} \preceq (1 + 0.1)\mathbf{S}_{\mathbf{x}^{(0)}}$  and therefore

$$\begin{aligned} \|\mathbf{e}^{(t+1)} - \mathbf{e}^{(t)}\|_2 &\leq c_{\Delta} \left\| \mathbf{S}_{\mathbf{x}^{(t)}}^{-1} (\mathbf{s}_{\mathbf{x}^{(t)}} - \mathbf{s}_{\mathbf{x}^{(t+1)}}) \right\|_2 \\ &\leq \frac{c_{\Delta}}{1 - 0.1} \left\| \mathbf{S}_{\mathbf{x}^{(0)}}^{-1} \mathbf{A} \left( \mathbf{x}^{(t)} - \mathbf{x}^{(t+1)} \right) \right\|_2 \\ &= \frac{c_{\Delta}}{1 - 0.1} \left\| \frac{1}{8} \mathbf{Q}^{-1} \nabla p_{e^{(t)}}(\mathbf{x}^{(t)}) \right\|_{\mathbf{A}^T \mathbf{S}_{\mathbf{x}^{(0)}}^{-2} \mathbf{A}} \\ &\leq \frac{c_{\Delta}}{8(1 - 0.1) \sqrt{c_e + \mu(\mathbf{x}^{(0)})}} \|\nabla p_{e^{(t)}}(\mathbf{x}^{(t)})\|_{\mathbf{Q}^{-1}} \end{aligned} \quad (8.5)$$

where in the last line we used  $\min_{i \in [m]} w_i \geq \mu(\mathbf{x}^{(0)})$ . Since  $c_{\Delta} < 0.01c_e \leq 0.01$ , by (8.4), we have

$$\begin{aligned} \|\nabla p_{e^{(t+1)}}(\mathbf{x}^{(t+1)})\|_{\mathbf{Q}^{-1}} &\leq \left(1 - \frac{1}{32}\right) \|\nabla p_{e^{(t)}}(\mathbf{x}^{(t)})\|_{\mathbf{Q}^{-1}} + \frac{0.01}{8(1 - 0.1)} \|\nabla p_{e^{(t)}}(\mathbf{x}^{(t)})\|_{\mathbf{Q}^{-1}} \\ &\leq \left(1 - \frac{1}{64}\right) \|\nabla p_{e^{(t)}}(\mathbf{x}^{(t)})\|_{\mathbf{Q}^{-1}} \quad . \end{aligned}$$

Furthermore, this implies that

$$\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(0)}\|_{\mathbf{Q}} \leq \left\| \sum_{k=0}^t \frac{1}{8} \mathbf{Q}^{-1} \nabla p_{e^{(k)}}(\mathbf{x}^{(k)}) \right\|_{\mathbf{Q}^{-1}} \leq \frac{1}{8} \sum_{i=0}^{\infty} \left(1 - \frac{1}{64}\right)^k \eta \leq \frac{64}{8} \eta = 8\eta \quad .$$

Similarly, we have that

$$\begin{aligned} \|\mathbf{e}^{(t+1)} - \mathbf{e}^{(0)}\|_2 &\leq \sum_{k=0}^t \frac{c_{\Delta}}{8(1 - 0.1) \sqrt{c_e + \mu(\mathbf{x}^{(0)})}} \left(1 - \frac{1}{64}\right)^k \|\nabla p_{e^{(0)}}(\mathbf{x}^{(0)})\|_{\mathbf{Q}^{-1}} \\ &\leq \frac{8c_{\Delta}\eta}{(1 - 0.1) \sqrt{c_e + \mu(\mathbf{x}^{(0)})}} \leq \frac{8c_{\Delta}}{(1 - 0.1) \sqrt{c_e + \mu(\mathbf{x}^{(0)})}} \delta_{e^{(0)}}(\mathbf{x}^{(0)}) \leq \frac{1}{10} c_{\Delta} \end{aligned}$$

where we used  $\eta = \|\nabla p_{e^{(0)}}(\mathbf{x}^{(0)})\|_{\mathbf{Q}^{-1}} \leq \|\nabla p_{e^{(0)}}(\mathbf{x}^{(0)})\|_{\mathbf{H}^{-1}} = \delta_{e^{(0)}}(\mathbf{x}^{(0)})$  and this finishes the induction on  $\|\nabla p_{e^{(t)}}(\mathbf{x}^{(t)})\|_{\mathbf{Q}^{-1}}$ ,  $\|\mathbf{x}^{(t)} - \mathbf{x}^{(0)}\|_{\mathbf{Q}}$  and  $\|\mathbf{e}^{(t)} - \mathbf{e}^{(0)}\|_2$ .

Hence, for all  $r$ , Lemma 8.3.6 shows that

$$\begin{aligned} \delta_{e^{(r)}}(\mathbf{x}^{(r)}) &= \|\nabla p_{e^{(r)}}(\mathbf{x}^{(r)})\|_{\mathbf{H}(\mathbf{x}^{(r)})^{-1}} \leq \sqrt{2} \|\nabla p_{e^{(r)}}(\mathbf{x}^{(r)})\|_{\mathbf{H}(\mathbf{x}^{(0)})^{-1}} \\ &\leq \sqrt{\frac{8}{3}} \|\nabla p_{e^{(r)}}(\mathbf{x}^{(r)})\|_{\mathbf{Q}^{-1}} \leq \sqrt{\frac{8}{3}} \left(1 - \frac{1}{64}\right)^r \|\nabla p_{e^{(0)}}(\mathbf{x}^{(0)})\|_{\mathbf{Q}^{-1}} \\ &\leq 2 \left(1 - \frac{1}{64}\right)^r \delta_{e^{(0)}}(\mathbf{x}^{(0)}). \end{aligned}$$

Using that  $\mathbf{E}\mathbf{e}^{(t+1)} = \mathbf{e}^{(t)}$ , we see that the expected change in function value is only due to the change

while taking centering steps and therefore Lemma 8.3.7 shows that

$$\mathbf{E}[p_{e^{(r)}}(\mathbf{x}^{(r)})] \geq p_{e^{(0)}}(\mathbf{x}^{(0)}) - \frac{1}{8} \sum_{k=0}^{\infty} \left(1 - \frac{1}{64}\right)^{2k} \|\nabla p_{e^{(0)}}(\mathbf{x}^{(0)})\|_{\mathbf{Q}^{-1}}^2 \geq p_{e^{(0)}}(\mathbf{x}^{(0)}) - 8 \left(\delta_{e^{(0)}}(\mathbf{x}^{(0)})\right)^2.$$

Finally, for (4), we note that

$$\left\| \frac{s(\mathbf{x}^{(r)}) - s(\mathbf{x}^{(0)})}{s(\mathbf{x}^{(0)})} \right\|_2 = \left\| \mathbf{x}^{(r)} - \mathbf{x}^{(0)} \right\|_{\mathbf{A}^T \mathbf{S}_{\mathbf{x}^{(0)}}^{-2} \mathbf{A}} \leq \frac{1}{\sqrt{\mu(\mathbf{x}^{(0)}) + c_e}} \left\| \mathbf{x}^{(r)} - \mathbf{x}^{(0)} \right\|_{\mathbf{Q}^{-1}} \leq \frac{1}{10}.$$

□

### ■ 8.3.2 Changing Constraints

Here we bound the effect that adding or a removing a constraint has on the hybrid barrier function. Much of the analysis in this section follows from the following lemma which follows easily from the Sherman Morrison Formula.

**Lemma 8.3.10** (Sherman Morrison Formula Implications). *Let  $\mathbf{B} \in \mathbb{R}^{n \times n}$  be an invertible symmetric matrix and let  $\mathbf{a} \in \mathbb{R}^n$  be arbitrary vector satisfying  $\mathbf{a}^T \mathbf{B}^{-1} \mathbf{a} < 1$ . The following hold:*

1.  $(\mathbf{B} \pm \mathbf{a}\mathbf{a}^T)^{-1} = \mathbf{B}^{-1} \mp \frac{\mathbf{B}^{-1}\mathbf{a}\mathbf{a}^T\mathbf{B}^{-1}}{1 \pm \mathbf{a}^T\mathbf{B}^{-1}\mathbf{a}}$ .
2.  $\mathbf{0} \preceq \frac{\mathbf{B}^{-1}\mathbf{a}\mathbf{a}^T\mathbf{B}^{-1}}{1 \pm \mathbf{a}^T\mathbf{B}^{-1}\mathbf{a}} \preceq \frac{\mathbf{a}^T\mathbf{B}^{-1}\mathbf{a}}{1 \pm \mathbf{a}^T\mathbf{B}^{-1}\mathbf{a}} \mathbf{B}^{-1}$ .
3.  $\log \det(\mathbf{B} \pm \mathbf{a}\mathbf{a}^T) = \log \det \mathbf{B} + \log(1 \pm \mathbf{a}^T \mathbf{B}^{-1} \mathbf{a})$ .

*Proof.* (1) follows immediately from Sherman Morrison Formula (Lemma 2.3.7). (2) follows since  $\mathbf{a}\mathbf{a}^T$  is PSD,

$$\frac{\mathbf{B}^{-1}\mathbf{a}\mathbf{a}^T\mathbf{B}^{-1}}{1 \pm \mathbf{a}^T\mathbf{B}^{-1}\mathbf{a}} = \mathbf{B}^{-1/2} \left[ \frac{\mathbf{B}^{-1/2}\mathbf{a}\mathbf{a}^T\mathbf{B}^{-1/2}}{1 \pm \mathbf{a}^T\mathbf{B}^{-1}\mathbf{a}} \right] \mathbf{B}^{-1/2},$$

and  $\mathbf{y}\mathbf{y}^T \preceq \|\mathbf{y}\|_2^2 \mathbf{I}$  for any vector  $\mathbf{y}$ . (3) follows immediately from the Matrix Determinant Lemma. □

We also make use of the following technical helper lemma.

**Lemma 8.3.11.** *For  $\mathbf{A} \in \mathbb{R}^{n \times m}$  and all  $\mathbf{a} \in \mathbb{R}^n$  we have*

$$\sum_{i \in [m]} \frac{1}{\psi_{\mathbf{A}}[i]} \left( \mathbf{1}_i^T (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a} \right)_i^4 \leq \left( \mathbf{a}^T (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a} \right)^2$$

*Proof.* We have by Cauchy Schwarz that

$$\left( \mathbf{1}_i^T (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a} \right)^2 \leq \psi_{\mathbf{A}}[i] \cdot \mathbf{a}^T (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a}$$

and consequently

$$\sum_{i \in [m]} \frac{\left( \mathbf{1}_i^T (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a} \right)^4}{\psi_{\mathbf{A}}[i]} \leq \left( \mathbf{a}^T (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a} \right) \sum_{i \in [m]} \left( \mathbf{1}_i (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a} \right)^2.$$

Since

$$\begin{aligned} \sum_{i \in [m]} \left( \mathbf{1}_i^T \mathbf{A} (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a} \right)^2 &= \mathbf{a}^T (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{A} (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a} \\ &\leq \mathbf{a}^T (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a}, \end{aligned}$$

we have the desired result.  $\square$

We now bound the effect of adding a constraint.

**Lemma 8.3.12.** *Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $\boldsymbol{\tau} \in \mathbb{R}^m$ , and  $\mathbf{x} \in P \stackrel{\text{def}}{=} \{\mathbf{y} \in \mathbb{R}^n : \mathbf{A}\mathbf{y} \geq \mathbf{b}\}$ . Let  $\overline{\mathbf{A}} \in \mathbb{R}^{(m+1) \times n}$  be  $\mathbf{A}$  with a row  $\mathbf{a}_{m+1}$  added, let  $\overline{\mathbf{b}} \in \mathbb{R}^{m+1}$  be the vector  $\mathbf{b}$  with an entry  $b_{m+1}$  added, and let  $\overline{P} \stackrel{\text{def}}{=} \{\mathbf{y} \in \mathbb{R}^n : \overline{\mathbf{A}}\mathbf{y} \geq \overline{\mathbf{b}}\}$ . Let  $s_{m+1} = \mathbf{a}_{m+1}^T \mathbf{x} - b_{m+1} > 0$ ,  $\psi_a = \frac{\mathbf{a}_{m+1}^T (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1} \mathbf{a}_{m+1}}{s_{m+1}^2}$ .*

Now, let  $\mathbf{v} \in \mathbb{R}^{m+1}$  be defined so that  $v_{m+1} = \frac{\psi_a}{1 + \psi_a}$  and for all  $i \in [m]$

$$v_i = \tau_i - \frac{1}{1 + \psi_a} \left[ \mathbf{A}_x (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1} \frac{\mathbf{a}_{m+1}}{s_{m+1}} \right]_i^2.$$

Then, the following hold

- [Leverage Score Estimation]  $e_{\overline{P}}(\vec{v}, \mathbf{x})_{m+1} = 0$  and  $e_{\overline{P}}(\vec{v}, \mathbf{x})_i = e_P(\boldsymbol{\tau}, \mathbf{x})_i$  for all  $i \in [m]$ .
- [Function Value Increase]  $p_{e_{\overline{P}}(\mathbf{v}, \mathbf{x})}(\mathbf{x}) = p_{e_P(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x}) - c_e \ln(s(\mathbf{x})_{m+1}/R) + \ln(1 + \psi_a)$ .
- [Centrality Increase]  $\delta_{e_{\overline{P}}(\mathbf{v}, \mathbf{x})}(\mathbf{x}) \leq \sqrt{1 + \psi_a} \delta_{e_P(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x}) + (c_e \sqrt{1 + \psi_a} + \psi_a) \sqrt{\frac{\psi_a}{\mu(\mathbf{x})}} + \psi_a$ .

*Proof.* By (1) in Lemma 8.3.10, we have that for all  $i \in [m]$

$$\psi_{\overline{P}}(\mathbf{x})_i = \psi_P(\mathbf{x})_i - \frac{1}{1 + \psi_a} \left[ \mathbf{A}_x (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1} \frac{\mathbf{a}_{m+1}}{s_{m+1}} \right]_i^2$$

and that

$$\psi_{\overline{P}}(\mathbf{x})_{m+1} = \psi_a - \frac{\psi_a^2}{1 + \psi_a} = \frac{\psi_a}{1 + \psi_a}.$$

Consequently [Leverage Score Estimation] holds.

By (3) in Lemma 8.3.10 we have that [Function Value Change] holds.

To bound the change in centrality, we first note that

$$\begin{aligned} \psi_{\overline{P}}(\mathbf{x})_i &\geq \psi_P(\mathbf{x})_i - \frac{1}{1 + \psi_a} \left[ \frac{\mathbf{a}_i^T}{s_i} (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1} \frac{\mathbf{a}_i}{s_i} \right] \left[ \frac{\mathbf{a}_{m+1}^T}{s_{m+1}} (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1} \frac{\mathbf{a}_{m+1}}{s_{m+1}} \right] \\ &= \frac{\psi_P(\mathbf{x})_i}{1 + \psi_a}. \end{aligned}$$

Therefore, we have that the approximate Hessian for  $\overline{P}$ , denoted  $\overline{\mathbf{H}}(\mathbf{x})$ , is bounded by  $\overline{\mathbf{H}}^{-1} \preceq (1 + \psi_a) \mathbf{H}^{-1}$ .

To bound the change in centrality note that by (2) in Lemma 8.3.10 we have that  $\overline{\mathbf{H}}^{-1} \preceq \mathbf{H}^{-1}$ . Therefore if let  $\mathbf{v}' \in \mathbb{R}^m$  be defined so that  $\mathbf{v}'_i = \mathbf{v}_i$  for all  $i \in [m]$  then by triangle inequality we have

$$\begin{aligned}
\delta_{e_{\bar{P}}(\mathbf{v}, \mathbf{x})}(\mathbf{x}) &= \|\bar{\mathbf{A}}_x^T(c_e \mathbf{1} + \mathbf{v})\|_{\mathbf{H}^{-1}} \leq \sqrt{1 + \psi_a} \|\bar{\mathbf{A}}_x^T(c_e \mathbf{1} + \mathbf{v})\|_{\mathbf{H}^{-1}} \\
&\leq \sqrt{1 + \psi_a} \left( \|\mathbf{A}_x^T(c_e \mathbf{1} + \boldsymbol{\tau})\|_{\mathbf{H}^{-1}} + \left\| \frac{\mathbf{a}_{m+1}}{s_{m+1}}(c_e + v_{m+1}) \right\|_{\mathbf{H}^{-1}} + \|\mathbf{A}_x^T(\mathbf{v}' - \boldsymbol{\tau})\|_{\mathbf{H}^{-1}} \right) \\
&= \sqrt{1 + \psi_a} \left( \delta_{e_P(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x}) + \left( c_e + \frac{\psi_a}{1 + \psi_a} \right) \left\| \frac{\mathbf{a}_{m+1}}{s_{m+1}} \right\|_{\mathbf{H}^{-1}} + \|\mathbf{A}_x^T(\mathbf{v}' - \boldsymbol{\tau})\|_{\mathbf{H}^{-1}} \right)
\end{aligned}$$

Now, since  $\mathbf{H}^{-1} \preceq \frac{1}{\mu(\mathbf{x})} (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1}$ , we have that

$$\left\| \frac{\mathbf{a}_{m+1}}{s_{m+1}} \right\|_{\mathbf{H}^{-1}} \leq \frac{1}{\sqrt{\mu(\mathbf{x})}} \left\| \frac{\mathbf{a}_{m+1}}{s_{m+1}} \right\|_{(\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1}} = \sqrt{\frac{\psi_a}{\mu(\mathbf{x})}}.$$

Since  $\boldsymbol{\Psi}^{1/2} \mathbf{A}_x (\mathbf{A}_x^T \boldsymbol{\Psi} \mathbf{A}_x)^{-1} \mathbf{A}_x^T \boldsymbol{\Psi}^{1/2}$  is a projection matrix, we have  $\boldsymbol{\Psi}^{-1} \succeq \mathbf{A}_x (\mathbf{A}_x^T \boldsymbol{\Psi} \mathbf{A}_x)^{-1} \mathbf{A}_x^T \succeq \mathbf{A}_x \mathbf{H}^{-1} \mathbf{A}_x^T$ . By Lemma 8.3.11, we have

$$\begin{aligned}
\|\mathbf{A}_x^T(\boldsymbol{\tau}' - \mathbf{v})\|_{\mathbf{H}^{-1}}^2 &\leq \|\boldsymbol{\tau}' - \mathbf{v}\|_{\boldsymbol{\Psi}^{-1}}^2 \\
&= \sum_{i \in [m]} \frac{1}{\psi(\mathbf{x})_i} \left( \frac{1}{1 + \psi_a} \left( \mathbf{1}_i \mathbf{A}_x (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1} \frac{\mathbf{a}_{m+1}}{s_{m+1}} \right)^2 \right)^2 \\
&\leq \left( \frac{1}{1 + \psi_a} \right)^2 \left( \frac{\mathbf{a}_{m+1}^T (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1} \mathbf{a}_{m+1}}{s_{m+1}^2} \right)^2 = \left( \frac{\psi_a}{1 + \psi_a} \right)^2
\end{aligned}$$

Combining, we have that

$$\begin{aligned}
\delta_{e_{\bar{P}}(\mathbf{v}, \mathbf{x})}(\mathbf{x}) &\leq \sqrt{1 + \psi_a} \left( \delta_{e_P(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x}) + \left( c_e + \frac{\psi_a}{1 + \psi_a} \right) \sqrt{\frac{\psi_a}{\mu(\mathbf{x})}} + \frac{\psi_a}{1 + \psi_a} \right) \\
&\leq \sqrt{1 + \psi_a} \delta_{e_P(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x}) + \left( c_e \sqrt{1 + \psi_a} + \psi_a \right) \sqrt{\frac{\psi_a}{\mu(\mathbf{x})}} + \psi_a.
\end{aligned}$$

□

We now bound the effect of removing a constraint.

**Lemma 8.3.13** (Removing a Constraint). *Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $\boldsymbol{\tau} \in \mathbb{R}^m$ , and  $\mathbf{x} \in P \stackrel{\text{def}}{=} \{\mathbf{y} \in \mathbb{R}^n : \mathbf{A}\mathbf{y} \geq \mathbf{b}\}$ . Let  $\bar{\mathbf{A}} \in \mathbb{R}^{(m-1) \times n}$  be  $\mathbf{A}$  with row  $m$  removed, let  $\bar{\mathbf{b}} \in \mathbb{R}^{m-1}$  denote the first  $m-1$  coordinates of  $\mathbf{b}$ , and let  $\bar{P} \stackrel{\text{def}}{=} \{\mathbf{y} \in \mathbb{R}^n : \bar{\mathbf{A}}\mathbf{y} \geq \bar{\mathbf{b}}\}$ . Let  $\psi_d = \psi_P(\mathbf{x})_m$ .*

*Now, let  $\mathbf{v} \in \mathbb{R}^{m-1}$  be defined so that for all  $i \in [m-1]$*

$$v_i = \tau_i + \frac{1}{1 - \psi_d} \left( \mathbf{A}_x (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1} \mathbf{A}_x^T \mathbf{1}_m \right)_i^2.$$

*Assume  $\psi_d \leq 1.1\mu(\mathbf{x}) \leq \frac{1}{10}$  and  $\|e_P(\boldsymbol{\tau}, \mathbf{x})\|_\infty \leq c_e \leq \frac{1}{2}$ , we have the following:*

- [Leverage Score Estimation]  $e_{\bar{P}}(\bar{\mathbf{v}}, \mathbf{x})_i = e_P(\boldsymbol{\tau}, \mathbf{x})_i$  for all  $i \in [m-1]$ .
- [Function Value Decrease]  $p_{e_{\bar{P}}(\mathbf{v}, \mathbf{x})}(\mathbf{x}) = p_{e_P(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x}) + [c_e + e_P(\boldsymbol{\tau}, \mathbf{x})_m] \ln(s(\mathbf{x})_m/R) + \ln(1 - \psi_d)$
- [Centrality Increase]  $\delta_{e_{\bar{P}}(\mathbf{v}, \mathbf{x})}(\mathbf{x}) \leq \frac{1}{\sqrt{1-2\mu(\mathbf{x})}} \delta_{e_P(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x}) + 3(c_e + \mu(\mathbf{x}))$ .

*Proof.* By (1) in Lemma 8.3.10, we have that for all  $i \in [m-1]$

$$\psi_{\bar{P}}(\mathbf{x})_i = \psi_P(\mathbf{x})_i + \frac{1}{1 - \psi_d} \left( \mathbf{1}_i^T \mathbf{A}_x (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1} \mathbf{A}_x^T \mathbf{1}_m \right)^2.$$

Consequently, [Leverage Score Estimation] holds. Furthermore, by (3) in Lemma 8.3.10, this then implies that [Function Value Change] holds.

To bound the change in centrality, we first note that by (1) and (2) in Lemma 8.3.10 and the fact  $\psi_{\bar{P}}(\mathbf{x})_i \geq \psi_P(\mathbf{x})_i$ , we have that the approximate Hessian for  $\bar{P}$ , denoted  $\bar{\mathbf{H}}(\mathbf{x})$ , is bounded by

$$\begin{aligned} \bar{\mathbf{H}}(\mathbf{x})^{-1} &\preceq \left( \mathbf{H}(\mathbf{x}) - \mathbf{A}_x^T (c_e \mathbf{I} + \Psi_x)^{1/2} \mathbf{1}_m \mathbf{1}_m^T (c_e \mathbf{I} + \Psi_x)^{1/2} \mathbf{A}_x \right)^{-1} \\ &\succeq \left( 1 + \frac{\alpha}{1 - \alpha} \right) \mathbf{H}(\mathbf{x})^{-1} = \left( \frac{1}{1 - \alpha} \right) \mathbf{H}(\mathbf{x})^{-1} \end{aligned}$$

where  $\alpha \stackrel{\text{def}}{=} \mathbf{1}_m^T (c_e \mathbf{I} + \Psi_x)^{1/2} \mathbf{A}_x \mathbf{H}(\mathbf{x})^{-1} \mathbf{A}_x^T (c_e \mathbf{I} + \Psi_x)^{1/2} \mathbf{1}_m$ . Using  $c_e + \mu(\mathbf{x}) \leq \frac{1}{2} + \frac{1}{10} \leq 1$ , we have

$$\mathbf{H}(\mathbf{x})^{-1} \preceq (\mathbf{A}_x^T (c_e + \mu(\mathbf{x})) \mathbf{A}_x + \lambda \mathbf{I})^{-1} \preceq (c_e + \mu(\mathbf{x}))^{-1} (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1}. \quad (8.6)$$

Using this, we have

$$\alpha \leq \left( \frac{c_e + \psi_d}{c_e + \mu(\mathbf{x})} \right) \mathbf{1}_m^T \mathbf{A}_x (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1} \mathbf{A}_x^T \mathbf{1}_m = \left( \frac{c_e + \psi_d}{c_e + \mu(\mathbf{x})} \right) \psi_d. \quad (8.7)$$

Now let  $\boldsymbol{\tau}' \in \mathbb{R}^{m-1}$  be defined so that  $\tau'_i = \tau_i$  for all  $i \in [m-1]$ . We have by above that

$$\delta_{e_{\bar{P}}(\mathbf{v}, \mathbf{x})}(\mathbf{x}) = \|\bar{\mathbf{A}}_x^T (c_e \mathbf{1} + \mathbf{v})\|_{\bar{\mathbf{H}}^{-1}} \leq \frac{1}{\sqrt{1 - \alpha}} \|\bar{\mathbf{A}}_x^T (c_e \mathbf{1} + \mathbf{v})\|_{\mathbf{H}^{-1}}$$

and therefore, by triangle inequality

$$\begin{aligned} \|\bar{\mathbf{A}}_x^T (c_e \mathbf{1} + \mathbf{v})\|_{\mathbf{H}^{-1}} &\leq \|\mathbf{A}_x^T (c_e \mathbf{1} + \boldsymbol{\tau})\|_{\mathbf{H}^{-1}} + \|\mathbf{A}_x^T \mathbf{1}_m (c_e + \tau_m)\|_{\mathbf{H}^{-1}} + \|\mathbf{A}_x^T (\boldsymbol{\tau}' - \mathbf{v})\|_{\mathbf{H}^{-1}} \\ &= \delta_{e_P(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x}) + (c_e + \tau_m) \|\mathbf{A}_x^T \mathbf{1}_m\|_{\mathbf{H}^{-1}} + \|\mathbf{A}_x^T (\boldsymbol{\tau}' - \mathbf{v})\|_{\mathbf{H}^{-1}}. \end{aligned}$$

Now, (8.6) shows that

$$\|\mathbf{A}_x^T \mathbf{1}_m\|_{\mathbf{H}^{-1}} \leq \frac{1}{\sqrt{c_e + \mu(\mathbf{x})}} \|\mathbf{A}_x^T \mathbf{1}_m\|_{(\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1}} \leq \sqrt{\frac{\psi_d}{c_e + \mu(\mathbf{x})}}$$

Furthermore, since  $\Psi^{-1} \succeq \mathbf{A}_x (\mathbf{A}_x^T \Psi \mathbf{A}_x)^{-1} \mathbf{A}_x^T \succeq \mathbf{A}_x \mathbf{H}^{-1} \mathbf{A}_x^T$ , by Lemma 8.3.11 we have

$$\begin{aligned} \|\mathbf{A}_x^T (\boldsymbol{\tau}' - \mathbf{v})\|_{\mathbf{H}^{-1}}^2 &\leq \|\boldsymbol{\tau}' - \mathbf{v}\|_{\Psi^{-1}}^2 \\ &= \sum_{i \in [m]} \frac{1}{\psi(\mathbf{x})_i} \left( \frac{1}{1 - \psi_d} \left( \mathbf{1}_i^T \mathbf{A}_x (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1} \mathbf{A}_x^T \mathbf{1}_m \right)^2 \right)^2 \\ &\leq \left( \frac{1}{1 - \psi_d} \right)^2 \left( \mathbf{1}_m^T \mathbf{A}_x (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1} \mathbf{A}_x^T \mathbf{1}_m \right)^2 = \left( \frac{\psi_d}{1 - \psi_d} \right)^2 \end{aligned}$$

Combining, we have that

$$\delta_{e_{\bar{P}}(\mathbf{v}, \mathbf{x})}(\mathbf{x}) \leq \frac{1}{\sqrt{1 - \alpha}} \left[ \delta_{e_P(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x}) + (c_e + \tau_m) \sqrt{\frac{\psi_d}{c_e + \mu(\mathbf{x})}} + \frac{\psi_d}{1 - \psi_d} \right].$$

Using the assumption  $\psi_d \leq 1.1\mu(\mathbf{x}) \leq \frac{1}{10}$ ,  $\|\mathbf{e}_P(\boldsymbol{\tau}, \mathbf{x})\|_\infty \leq c_e$  and (8.7), we have  $\alpha \leq 1.1\psi_d \leq 1.21\mu(\mathbf{x})$  and  $\tau_m \leq \psi_d + c_e$ , and

$$\begin{aligned} \delta_{e_{\bar{P}}(\mathbf{v}, \mathbf{x})}(\mathbf{x}) &\leq \frac{1}{\sqrt{1 - 1.3\mu(\mathbf{x})}} \left[ \delta_{e_P(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x}) + (c_e + \tau_m)\sqrt{1.1} + 1.2\psi_d \right] \\ &\leq \frac{1}{\sqrt{1 - 2\mu(\mathbf{x})}} \delta_{e_P(\boldsymbol{\tau}, \mathbf{x})}(\mathbf{x}) + \frac{1}{\sqrt{1 - \frac{1.3}{10}}} \left( \sqrt{1.1} \cdot 2c_e + 1.1(\sqrt{1.1} + 1.2)\mu(\mathbf{x}) \right) \end{aligned}$$

□

### ■ 8.3.3 Hybrid Center Properties

Here we prove properties of points near the hybrid center. First we bound the distance between points in the  $\mathbf{H}(\mathbf{x})$  norm in terms of the  $\ell_2$  norm of the points.

**Lemma 8.3.14.** *For  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{b} \in \mathbb{R}^m$  suppose that  $\mathbf{x} \in P = \{\mathbf{y} : \mathbf{A}\mathbf{y} \geq \mathbf{b}\}$  and  $\mathbf{e} \in \mathbb{R}^m$  such that  $\|\mathbf{e}\|_\infty \leq \frac{1}{2}c_e < \frac{1}{20}$  and  $\delta_e \leq 0.1\sqrt{c_e + \mu(\mathbf{x})}$ . Then for all  $\mathbf{y} \in P$  we have*

$$\|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})} \leq \frac{12mc_e + 6n + 2\lambda\|\mathbf{y}\|_2^2}{\sqrt{c_e + \mu(\mathbf{x})}} \quad (8.8)$$

and

$$\|\mathbf{x}\|_2^2 \leq 4\lambda^{-1}(mc_e + n) + 2\|\mathbf{y}\|_2^2 \quad .$$

*Proof.* For notational simplicity let  $\mathbf{t} \stackrel{\text{def}}{=} c_e \mathbf{1} + \mathbf{e} + \boldsymbol{\psi}_x$ ,  $\mathbf{T} \stackrel{\text{def}}{=} \mathbf{Diag}(\mathbf{t})$ , and  $\mathbf{M} \stackrel{\text{def}}{=} \mathbf{A}_x^T (c_e \mathbf{I} + \boldsymbol{\Psi}_x) \mathbf{A}_x$ . We have

$$\|\mathbf{x} - \mathbf{y}\|_{\mathbf{A}_x^T \mathbf{T} \mathbf{A}_x}^2 = \sum_{i \in [m]} t_i \frac{[\mathbf{s}_x - \mathbf{s}_y]_i^2}{[\mathbf{s}_x]_i^2} = \sum_{i \in [m]} t_i \left( 1 - 2 \frac{[\mathbf{s}_y]_i}{[\mathbf{s}_x]_i} + \frac{[\mathbf{s}_y]_i^2}{[\mathbf{s}_x]_i^2} \right) \quad (8.9)$$

and

$$\sum_{i \in [m]} t_i \frac{[\mathbf{s}_y]_i^2}{[\mathbf{s}_x]_i^2} \leq \left( \sum_{i \in [m]} t_i \frac{[\mathbf{s}_y]_i}{[\mathbf{s}_x]_i} \right) \max_{i \in [m]} \frac{[\mathbf{s}_y]_i}{[\mathbf{s}_x]_i} \leq \left( \sum_{i \in [m]} t_i \frac{[\mathbf{s}_y]_i}{[\mathbf{s}_x]_i} \right) (1 + \|\mathbf{S}_x^{-1}(\mathbf{s}_y - \mathbf{s}_x)\|_\infty) \quad (8.10)$$

and

$$\begin{aligned} \|\mathbf{S}_x^{-1}(\mathbf{s}_x - \mathbf{s}_y)\|_\infty &= \max_{i \in [m]} |\mathbf{1}_i \mathbf{S}_x^{-1} \mathbf{A}(\mathbf{x} - \mathbf{y})| \leq \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})} \sqrt{\max_{i \in [m]} [\mathbf{S}_x^{-1} \mathbf{A} \mathbf{H}(\mathbf{x})^{-1} \mathbf{A}^T \mathbf{S}_x^{-1}]_{ii}} \\ &\leq (c_e + \mu(\mathbf{x}))^{-1/2} \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})} \quad . \end{aligned} \quad (8.11)$$

Now, clearly  $\sum_{i \in [m]} t_i [\mathbf{s}_y]_i / [\mathbf{s}_x]_i$  is positive and since  $\|\mathbf{e}\|_\infty \leq \frac{1}{2}c_e$  we know that  $\frac{1}{2}\mathbf{M} \preceq \mathbf{A}_x^T \mathbf{T} \mathbf{A}_x$ . Therefore, by combining (8.9), (8.10), and (8.11) we have

$$\begin{aligned} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_{\mathbf{M}}^2 &\leq \|\mathbf{t}\|_1 - \sum_{i \in [m]} t_i \frac{[\mathbf{s}_y]_i}{[\mathbf{s}_x]_i} + \left( \sum_{i \in [m]} t_i \frac{[\mathbf{s}_y]_i}{[\mathbf{s}_x]_i} \right) \frac{\|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})}}{\sqrt{c_e + \mu(\mathbf{x})}} \\ &\leq \|\mathbf{t}\|_1 + \left( \sum_{i \in [m]} t_i \frac{[\mathbf{s}_y]_i}{[\mathbf{s}_x]_i} \right) \frac{\|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})}}{\sqrt{c_e + \mu(\mathbf{x})}} \end{aligned} \quad (8.12)$$

Now since  $\nabla p_e(\mathbf{x}) = -\mathbf{A}^T \mathbf{S}_x^{-1} \mathbf{T} \mathbf{1} + \lambda \mathbf{x}$  we have

$$\langle \mathbf{x} - \mathbf{y}, \nabla p_e(\mathbf{x}) \rangle = - \sum_{i \in [m]} t_i \frac{[\mathbf{s}_x - \mathbf{s}_y]_i}{[\mathbf{s}_x]_i} + \lambda \mathbf{x}^T (\mathbf{x} - \mathbf{y})$$

and therefore by Cauchy Schwarz and  $\mathbf{x}^T \mathbf{y} \leq \|\mathbf{x}\|_2^2 + \frac{1}{4} \|\mathbf{y}\|_2^2$ ,

$$\sum_{i \in [m]} t_i \frac{[\mathbf{s}_y]_i}{[\mathbf{s}_x]_i} = \|\mathbf{t}\|_1 - \lambda \|\mathbf{x}\|_2^2 + \lambda \mathbf{x}^T \mathbf{y} + \langle \mathbf{x} - \mathbf{y}, \nabla p_e(\mathbf{x}) \rangle \quad (8.13)$$

$$\leq \|\mathbf{t}\|_1 + \frac{\lambda}{4} \|\mathbf{y}\|_2^2 + \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})} \delta_e(\mathbf{x}) \quad . \quad (8.14)$$

Now, using (8.12), (8.14) and the definition of  $\mathbf{H}(\mathbf{x})$ , we have

$$\begin{aligned} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})}^2 &= \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_M^2 + \frac{\lambda}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \\ &\leq \|\mathbf{t}\|_1 + \left( \sum_{i \in [m]} t_i \frac{[\mathbf{s}_y]_i}{[\mathbf{s}_x]_i} \right) \frac{\|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})}}{\sqrt{c_e + \mu(\mathbf{x})}} + \frac{\lambda}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \\ &\leq \|\mathbf{t}\|_1 + \frac{\|\mathbf{t}\|_1 + \frac{\lambda}{4} \|\mathbf{y}\|_2^2}{\sqrt{c_e + \mu(\mathbf{x})}} \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})} + \delta_e(\mathbf{x}) \frac{\|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})}^2}{\sqrt{c_e + \mu(\mathbf{x})}} + \frac{\lambda}{2} \|\mathbf{x} - \mathbf{y}\|_2^2. \end{aligned}$$

Using the fact that  $\delta_e(\mathbf{x}) \leq 0.1 \sqrt{c_e + \mu(\mathbf{x})}$ , we have

$$\frac{1}{4} \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})}^2 \leq \|\mathbf{t}\|_1 + \frac{\lambda}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \frac{\|\mathbf{t}\|_1 + \frac{\lambda}{4} \|\mathbf{y}\|_2^2}{\sqrt{c_e + \mu(\mathbf{x})}} \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})}. \quad (8.15)$$

Furthermore, since  $\sum_{i \in [m]} t_i [\mathbf{s}_y]_i / [\mathbf{s}_x]_i$  is positive, (8.13) shows that

$$\lambda \mathbf{x}^T (\mathbf{x} - \mathbf{y}) = \lambda \|\mathbf{x}\|_2^2 - \lambda \mathbf{x}^T \mathbf{y} \leq \|\mathbf{t}\|_1 + \langle \mathbf{x} - \mathbf{y}, \nabla p_e(\mathbf{x}) \rangle \leq \|\mathbf{t}\|_1 + \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})} \delta_e(\mathbf{x})$$

and hence

$$\begin{aligned} \frac{\lambda}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 &\leq \frac{\lambda}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \|\mathbf{x}\|_2^2 = \lambda \mathbf{x}^T (\mathbf{x} - \mathbf{y}) + \frac{\lambda}{2} \|\mathbf{y}\|_2^2 \\ &\leq \|\mathbf{t}\|_1 + \frac{\lambda}{2} \|\mathbf{y}\|_2^2 + \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})} \delta_e(\mathbf{x}) \quad . \end{aligned} \quad (8.16)$$

Putting (8.16) into (8.15) and using the fact that  $\delta_e(\mathbf{x}) \leq 0.1 \sqrt{c_e + \mu(\mathbf{x})}$ , we have

$$\frac{1}{4} \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})}^2 \leq 2\|\mathbf{t}\|_1 + \frac{\lambda}{2} \|\mathbf{y}\|_2^2 + \left( 0.1 + \frac{\|\mathbf{t}\|_1 + \frac{\lambda}{4} \|\mathbf{y}\|_2^2}{\sqrt{c_e + \mu(\mathbf{x})}} \right) \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})}.$$

Now, using  $\|\mathbf{t}\|_1 \leq 2mc_e + n$  and  $\sqrt{c_e + \mu(\mathbf{x})} \leq 1.05$ , we have

$$\frac{1}{4} \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})}^2 \leq 2.2\alpha + (0.1 + \alpha) \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})} \quad \text{for} \quad \alpha = \frac{2mc_e + n + \frac{\lambda}{4} \|\mathbf{y}\|_2^2}{\sqrt{c_e + \mu(\mathbf{x})}} \quad .$$

Since  $\alpha \geq 1/\sqrt{1.1}$ , we have that

$$\|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})} \leq \frac{0.1 + \alpha + \sqrt{(\alpha + 0.1)^2 + 2.2\alpha}}{2 \cdot \frac{1}{4}} \leq 6\alpha$$

yielding (8.8).

We also have by (8.13) and the fact that  $\delta_e(\mathbf{x}) \leq 0.1\sqrt{c_e + \mu(\mathbf{x})}$ ,

$$\begin{aligned} \lambda\|\mathbf{x}\|_2^2 &= \|\mathbf{t}\|_1 + \lambda\mathbf{x}^T\mathbf{y} + \langle \mathbf{x} - \mathbf{y}, \nabla p_e(\mathbf{x}) \rangle - \sum_{i \in [m]} t_i \frac{[\mathbf{s}_y]_i}{[\mathbf{s}_x]_i} \\ &\leq \|\mathbf{t}\|_1 + \frac{\lambda}{2}\|\mathbf{x}\|_2^2 + \frac{\lambda}{2}\|\mathbf{y}\|_2^2 + \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})}\delta_e(\mathbf{x}) \\ &\leq \|\mathbf{t}\|_1 + \frac{\lambda}{2}\|\mathbf{x}\|_2^2 + \frac{\lambda}{2}\|\mathbf{y}\|_2^2 + 0.1\sqrt{c_e + \mu(\mathbf{x})}\|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})} \end{aligned}$$

Hence, using  $\|\mathbf{t}\|_1 \leq 2mc_e + n$  and  $\|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}(\mathbf{x})} \leq 6\alpha$ , we have

$$\begin{aligned} \frac{\lambda}{2}\|\mathbf{x}\|_2^2 &\leq \|\mathbf{t}\|_1 + \frac{\lambda}{2}\|\mathbf{y}\|_2^2 + 0.6 \left( 2mc_e + n + \frac{\lambda}{4}\|\mathbf{y}\|_2^2 \right) \\ &\leq \lambda\|\mathbf{y}\|_2^2 + 2(mc_e + n). \end{aligned}$$

□

In the following lemma we show how we can write one hyperplane in terms of the others provided that we are nearly centered and show there is a constraint that the central point is close to.

**Lemma 8.3.15.** *Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{b} \in \mathbb{R}^m$  such that  $\|a_i\|_2 = 1$  for all  $i$ . Suppose that  $\mathbf{x} \in P = \{\mathbf{y} : \mathbf{A}\mathbf{y} \geq \mathbf{b}\}$  and  $\mathbf{e} \in \mathbb{R}^m$  such that  $\|\mathbf{e}\|_\infty \leq \frac{1}{2}c_e \leq \frac{1}{2}$ . Furthermore, let  $\varepsilon = \min_{j \in [m]} s_j(\mathbf{x})$  and suppose that  $i = \arg \min_{j \in [m]} s_j(\mathbf{x})$  then*

$$\left\| \mathbf{a}_i + \sum_{j \neq i} \left( \frac{s(x)_i}{s(x)_j} \right) \left( \frac{c_e + e_j + \psi_j(\mathbf{x})}{c_e + e_i + \psi_i(\mathbf{x})} \right) \mathbf{a}_j \right\|_2 \leq \frac{2\varepsilon}{(c_e + \mu(\mathbf{x}))} \left[ \lambda\|\mathbf{x}\|_2 + \delta_e(\mathbf{x}) \sqrt{\frac{mc_e + n}{\varepsilon^2} + \lambda} \right].$$

*Proof.* We know that

$$\begin{aligned} \nabla p_e(\mathbf{x}) &= -\mathbf{A}^T \mathbf{S}_x^{-1} (c_e \mathbf{1} + \mathbf{e} + \boldsymbol{\psi}_x) + \lambda \mathbf{x} \\ &= \lambda \mathbf{x} - \sum_{i \in [m]} \frac{(c_e + e_i + \psi_i)}{s(\mathbf{x})_i} \mathbf{a}_i \end{aligned}$$

Consequently, by  $\|\mathbf{e}\|_\infty \leq \frac{1}{2}c_e$ , and  $\psi_i(\mathbf{x}) \geq \mu(\mathbf{x})$ , we have

$$\begin{aligned} \left\| \mathbf{a}_i + \sum_{j \neq i} \left( \frac{s(x)_i}{s(x)_j} \right) \left( \frac{c_e + e_j + \psi_j(\mathbf{x})}{c_e + e_i + \psi_i(\mathbf{x})} \right) \mathbf{a}_j \right\|_2 &= \frac{s_i(\mathbf{x})}{c_e + e_i + \psi_i(\mathbf{x})} \|\mathbf{A}^T \mathbf{S}_x^{-1} (c_e \mathbf{1} + \mathbf{e} + \boldsymbol{\psi}_x)\|_2 \\ &\leq \frac{2\varepsilon}{c_e + \mu(\mathbf{x})} [\lambda\|\mathbf{x}\|_2 + \|\nabla p_e(\mathbf{x})\|_2]. \end{aligned}$$

Using  $\|\mathbf{a}_i\| = 1$ ,  $\sum_i \psi_i \leq n$ , and  $s_i(\mathbf{x}) \geq \varepsilon$ , we have

$$\begin{aligned} \text{Tr}(\mathbf{A}_x^T (c_e \mathbf{I} + \boldsymbol{\Psi}_x) \mathbf{A}_x) &= \text{Tr}(\mathbf{A}_x \mathbf{A}_x^T (c_e \mathbf{I} + \boldsymbol{\Psi}_x)) \\ &= \sum_i (c_e + \psi_i) \frac{\|a_i\|_2^2}{s_i^2(\mathbf{x})} \leq \frac{mc_e + n}{\varepsilon^2}. \end{aligned} \tag{8.17}$$

Hence, we have  $\mathbf{H}(\mathbf{x}) \preceq \left( \frac{mc_e + n}{\varepsilon^2} + \lambda \right) \mathbf{I}$  and  $\|\nabla p_e(\mathbf{x})\|_2 \leq \delta_e(\mathbf{x}) \sqrt{\frac{mc_e + n}{\varepsilon^2} + \lambda}$  yielding the result. □



### ■ 8.3.4 The Algorithm

Here, we put all the results in the previous sections to get our ellipsoid algorithm. Below is a sketch of the pseudocode; we use  $c_a, c_d, c_e, c_\Delta$  to denote parameters we decide later.

---

#### Algorithm 22: Our Cutting Plane Method

---

**Input:**  $\mathbf{A}^{(0)} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b}^{(0)} \in \mathbb{R}^m$ ,  $\varepsilon > 0$ , and radius  $R > 0$ .  
**Input:** A separation oracle for a non-empty set  $K \subset B_\infty(R)$ .  
**Check:** Throughout the algorithm, if  $s_i(\mathbf{x}^{(k)}) < \varepsilon$  output  $P^{(k)}$ .  
**Check:** Throughout the algorithm, if  $\mathbf{x}^{(k)} \in K$ , output  $\mathbf{x}^{(k)}$ .  
Set  $P^{(0)} = B_\infty(R)$ .  
Set  $\mathbf{x}^{(0)} := \mathbf{0}$  and compute  $\tau_i^{(0)} = \psi_{P^{(0)}}(\mathbf{x}^{(0)})_i$  for all  $i \in [m]$  exactly.  
**for**  $k = 0$  **to**  $\infty$  **do**  
    Let  $m^{(k)}$  be the number of constraints in  $P^{(k)}$ .  
    Compute  $\mathbf{w}^{(k)}$  such that  $\Psi_{P^{(k)}}(\mathbf{x}^{(k)}) \preceq \mathbf{W}^{(k)} \preceq (1 + c_\Delta) \Psi_{P^{(k)}}(\mathbf{x}^{(k)})$ .  
    Let  $i^{(k)} \in \arg \max_{i \in [m^{(k)}]} |w_i^{(k)} - \tau_i^{(k)}|$ .  
    Set  $\tau_{i^{(k)}}^{(k+\frac{1}{3})} = \psi_{P^{(k)}}(\mathbf{x}^{(k)})_{i^{(k)}}$  and  $\tau_j^{(k+\frac{1}{3})} = \tau_j^{(k)}$  for all  $j \neq i^{(k)}$ .  
    **if**  $\min_{i \in [m^{(k)}]} w_i^{(k)} \leq c_d$  **then**  
        Remove constraint with minimum  $w_i^{(k)}$  yielding polytope  $P^{(k+1)}$ .  
        Update  $\tau$  according to Lemma 8.3.13 to get  $\tau_j^{(k+\frac{2}{3})}$ .  
    **else**  
        Use separation oracle at  $\mathbf{x}^{(k)}$  to get a constraint  $\{\mathbf{x} : \mathbf{a}^T \mathbf{x} \geq \mathbf{a}^T \mathbf{x}^{(k)}\}$  with  $\|\mathbf{a}\|_2 = 1$ .  
        Add constraint  $\{\mathbf{x} : \mathbf{a}^T \mathbf{x} \geq \mathbf{a}^T \mathbf{x}^{(k)} - c_a^{-1/2} \sqrt{\mathbf{a}^T (\mathbf{A}^T \mathbf{S}_{\mathbf{x}^{(k)}}^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a}}\}$  yielding polytope  $P^{(k+1)}$ .  
        Update  $\tau$  according to Lemma 8.3.12 to get  $\tau_j^{(k+\frac{2}{3})}$ .  
     $(\mathbf{x}^{(k+1)}, \tau^{(k+1)}) = \text{Centering}(\mathbf{x}^{(k)}, \tau^{(k+\frac{2}{3})}, 200, c_\Delta)$ .  
**end**

---

In the algorithm, there are two main invariants we maintain. First, we maintain that the centrality  $\delta_{P,e}(\mathbf{x})$ , which indicates how close  $\mathbf{x}$  is to the minimum point of  $p_e$ , is small. Second, we maintain that  $\|\mathbf{e}(\tau, \mathbf{x})\|_\infty$ , which indicates how accurate the leverage score estimate  $\tau$  is, is small. In the following lemma we show that we maintain both invariants throughout the algorithm.

**Lemma 8.3.16.** *Assume that  $c_e \leq c_d \leq \frac{1}{10^6}$ ,  $c_a \sqrt{c_a} \leq \frac{c_d}{10^3}$ ,  $c_d \leq c_a$ , and  $c_\Delta \leq C c_e / \log(n \log(R/\varepsilon))$  for some small enough universal constant  $C$ . During our cutting plane method, for all  $k$ , with high probability in  $n$ , we have*

1.  $\|\mathbf{e}(\tau^{(k+\frac{1}{3})}, \mathbf{x}^{(k)})\|_\infty \leq \frac{1}{1000} c_e$ ,  $\|\mathbf{e}(\tau^{(k+\frac{2}{3})}, \mathbf{x}^{(k)})\|_\infty \leq \frac{1}{1000} c_e$ ,  $\|\mathbf{e}(\tau^{(k+1)}, \mathbf{x}^{(k+1)})\|_\infty \leq \frac{1}{400} c_e$ .
2.  $\delta_{P^{(k)}, \mathbf{e}(\tau^{(k+\frac{2}{3})}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) \leq \frac{1}{100} \sqrt{c_e + \min(\mu(\mathbf{x}^{(k)}), c_d)}$ .
3.  $\delta_{P^{(k+1)}, \mathbf{e}(\tau^{(k+1)}, \mathbf{x}^{(k+1)})}(\mathbf{x}^{(k+1)}) \leq \frac{1}{400} \sqrt{c_e + \min(\mu(\mathbf{x}^{(k+1)}), c_d)}$ .

*Proof.* Some statements of the proof hold only with high probability in  $n$ ; we omit mentioning this for simplicity.

We prove by induction on  $k$ . Note that the claims are written in order consistent with the algorithm and proving the statement for  $k$  involves bounding centrality at the point  $\mathbf{x}^{(k+1)}$ . Trivially we define,  $\boldsymbol{\tau}^{(-1)} = \boldsymbol{\tau}^{(-\frac{2}{3})} = \boldsymbol{\tau}^{(-\frac{1}{3})} = \boldsymbol{\tau}^{(0)}$ ,  $\mathbf{x}^{(-1)} = \mathbf{x}^{(0)}$  and note that the claims then hold for  $k = -1$  as we compute the initial leverage scores,  $\boldsymbol{\tau}^{(0)}$ , exactly and since the polytope is symmetric we have  $\delta_{e(\boldsymbol{\tau}^{(0)}, \mathbf{x}^{(0)})}(\mathbf{0}) = 0$ . We now suppose they hold for all  $t < k$  and show that they hold for  $t = k$ .

We first bound  $\delta$ . For notational simplicity, let  $\eta_t \stackrel{\text{def}}{=} \sqrt{c_e + \min(\mu(\mathbf{x}^{(t)}), c_d)}$ . By the induction hypothesis we know that  $\delta_{P^{(t)}, e(\boldsymbol{\tau}^{(t)}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) \leq \frac{1}{400}\eta_t$ . Now, when we update  $\boldsymbol{\tau}^{(t)}$  to  $\boldsymbol{\tau}^{(t+\frac{1}{3})}$ , we set  $e_{i(t)}$  to 0. Consequently, Lemma 8.3.8 and the induction hypothesis  $\|e(\boldsymbol{\tau}^{(t)}, \mathbf{x}^{(t)})\|_\infty \leq \frac{1}{400}c_e$  show that

$$\begin{aligned} \delta_{P^{(t)}, e(\boldsymbol{\tau}^{(t+\frac{1}{3})}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) &\leq \delta_{P^{(t)}, e(\boldsymbol{\tau}^{(t)}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) + \frac{1}{\sqrt{c_e + \mu(\mathbf{x}^{(t)})}} e_{i(t)}(\boldsymbol{\tau}^{(t)}, \mathbf{x}^{(t)}) \\ &\leq \frac{1}{400}\eta_t + \frac{\sqrt{c_e}}{400} \leq \frac{\eta_t}{200} \end{aligned} \quad (8.18)$$

Next, we estimate the  $\delta$  changes when we remove or add a constraint.

For the case of removal, we note that it happens only if  $\mu(\mathbf{x}^{(t)}) \leq \min_i w_i \leq c_d \leq \frac{1}{10^6}$ . Also, the row we remove has leverage score at most  $1.1\mu(\mathbf{x}^{(t)})$  because we pick the row with minimum  $w$ . Hence, Lemma 8.3.13 show that

$$\begin{aligned} \delta_{P^{(t+1)}, e(\boldsymbol{\tau}^{(t+\frac{2}{3})}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) &\leq \frac{1}{\sqrt{1 - 2\mu(\mathbf{x}^{(t)})}} \delta_{P^{(t)}, e(\boldsymbol{\tau}^{(t+\frac{1}{3})}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) + 3(c_e + \mu(\mathbf{x}^{(t)})) \\ &\leq \frac{1}{\sqrt{1 - 2 \cdot 10^{-6}}} \left( \frac{\eta_t}{200} \right) + 3(c_e + \mu(\mathbf{x}^{(t)})) \leq \frac{\eta_t}{100} \end{aligned}$$

where we used the fact  $\mu(\mathbf{x}^{(t)}) \leq c_d$  and hence  $c_e + \mu(\mathbf{x}^{(t)}) \leq \sqrt{2c_d}\eta_t \leq \frac{\sqrt{2}}{1000}\eta_t$ .

For the case of addition, we note that it happens only if  $2\mu(\mathbf{x}^{(t)}) \geq \min_i w_i \geq c_d$ . Furthermore, in this case the hyperplane we add is chosen precisely so that  $\psi_a = c_a$ . Furthermore, since  $c_e \leq c_d \leq c_a \leq \frac{1}{100}$  by Lemma 8.3.12 we have that

$$\begin{aligned} \delta_{P^{(t+1)}, e(\boldsymbol{\tau}^{(t+\frac{2}{3})}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) &\leq \sqrt{1 + c_a} \delta_{P^{(t)}, e(\boldsymbol{\tau}^{(t+\frac{1}{3})}, \mathbf{x}^{(t)})} + (c_e \sqrt{1 + c_a} + c_a) \sqrt{\frac{c_a}{\mu(\mathbf{x}^{(t)})}} + c_a \\ &\leq \frac{\eta_t}{190} + 4c_a \sqrt{\frac{c_a}{c_d}}. \end{aligned}$$

Furthermore, since  $c_a \sqrt{c_a} \leq \frac{c_d}{1000}$ ,  $\mu(\mathbf{x}^{(t)}) \geq c_d/2$ , and  $c_d \leq 10^{-6}$  we know that  $4c_a \sqrt{c_a/c_d} \leq \frac{1}{250}\eta_t$  and consequently in both cases we have  $\delta_{P^{(t+1)}, e(\boldsymbol{\tau}^{(t+\frac{2}{3})}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) \leq \frac{1}{100}\eta_t$ .

Now, note that Lemmas 8.3.12 and 8.3.13 show that  $e$  does not change during the addition or removal of an constraint. Hence, we have  $\|e(\boldsymbol{\tau}^{(t+\frac{2}{3})}, \mathbf{x}^{(t)})\|_\infty \leq \|e(\boldsymbol{\tau}^{(t+\frac{1}{3})}, \mathbf{x}^{(t)})\|_\infty$ . Furthermore, we know the step " $\boldsymbol{\tau}_{i^{(k)}}^{(k+\frac{1}{3})} = \psi_{P^{(k)}}(\mathbf{x}^{(k)})_{i^{(k)}}$ " only decreases  $\|e\|_\infty$  and hence we have  $\|e(\boldsymbol{\tau}^{(t+\frac{2}{3})}, \mathbf{x}^{(t)})\|_\infty \leq \|e(\boldsymbol{\tau}^{(t)}, \mathbf{x}^{(t)})\|_\infty \leq \frac{c_e}{400}$ . Thus, we have all the conditions needed for Lemma 8.3.9 and consequently

$$\delta_{P^{(t+1)}, e(\boldsymbol{\tau}^{(t+1)}, \mathbf{x}^{(t+1)})}(\mathbf{x}^{(t+1)}) \leq 2 \left(1 - \frac{1}{64}\right)^{200} \delta_{P^{(t+1)}, e(\boldsymbol{\tau}^{(t+\frac{2}{3})}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) \leq \frac{1}{1000}\eta_t \quad .$$

Lemma 8.3.9 also shows that that  $\left\| \frac{s(\mathbf{x}^{(t+1)}) - s(\mathbf{x}^{(t)})}{s(\mathbf{x}^{(t)})} \right\|_2 \leq \frac{1}{10}$  and hence  $\psi_i(\mathbf{x}^{(t)}) \leq 2\psi_i(\mathbf{x}^{(t+1)})$  for all  $i$ .

Therefore,  $\eta_t \leq 2\eta_{t+1}$  and thus

$$\delta_{P^{(t+1)}, \mathbf{e}(\boldsymbol{\tau}^{(t+1)}, \mathbf{x}^{(t+1)})}(\mathbf{x}^{(t+1)}) \leq \frac{\sqrt{c_e + \min(c_d, \mu(\mathbf{x}^{(t+1)}))}}{400}.$$

completing the induction case for  $\delta$ .

Now, we bound  $\|\mathbf{e}\|_\infty$ . Lemma 8.3.12 and 8.3.13 show that  $\mathbf{e}$  does not change during the addition or removal of an constraint. Hence,  $\mathbf{e}$  is affected by only the update step “ $\tau_{i^{(k)}}^{(k+\frac{1}{2})} = \psi_{P^{(k)}}(\mathbf{x}^{(k)})_{i^{(k)}}$ ” and the centering step. Using the induction hypothesis  $\delta_{P^{(t)}, \mathbf{e}(\boldsymbol{\tau}^{(t+\frac{2}{3})}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) \leq \frac{1}{100}\eta_t$  and Lemma 8.3.9 shows that  $\mathbf{E}\mathbf{e}(\boldsymbol{\tau}^{(t+1)}, \mathbf{x}^{(t+1)}) = \mathbf{e}(\boldsymbol{\tau}^{(t+\frac{2}{3})}, \mathbf{x}^{(t)})$  and  $\|\mathbf{e}(\boldsymbol{\tau}^{(t+1)}, \mathbf{x}^{(t+1)}) - \mathbf{e}(\boldsymbol{\tau}^{(t+\frac{2}{3})}, \mathbf{x}^{(t)})\|_2 \leq \frac{1}{10}c_\Delta$ . The goal for the update step is to decrease  $\mathbf{e}$  by updating  $\boldsymbol{\tau}$ . In Section 8.4.2, we give a self-contained analysis of the effect of this step as a game. In each round, the vector  $\mathbf{e}$  is corrupted by some mean 0 and bounded variance noise and the problem is how to update  $\mathbf{e}$  such that  $\|\mathbf{e}\|_\infty$  is bounded. Theorem 8.4.3 shows that we can do this by setting the  $e_i = 0$  for the almost maximum coordinate in each iteration. This is exactly what the update step is doing. Since our algorithm would run only  $O(n \log(nR/\varepsilon))$  many iterations, Theorem 8.4.3 shows that this strategy guarantees that after the update step, we have

$$\left\| \mathbf{e}(\boldsymbol{\tau}^{(t+\frac{1}{3})}, \mathbf{x}^{(t)}) \right\|_\infty = O(c_\Delta \log(n \log(R/\varepsilon))).$$

Now, by our choice of  $c_\Delta$ , we have  $\|\mathbf{e}(\boldsymbol{\tau}^{(t+\frac{1}{3})}, \mathbf{x}^{(t)})\|_\infty \leq \frac{1}{1000}c_e$ . Lemma 8.3.12 and 8.3.13 show that  $\mathbf{e}$  does not change during the addition or removal of an constraint. Hence, we have  $\|\mathbf{e}(\boldsymbol{\tau}^{(t+\frac{2}{3})}, \mathbf{x}^{(t)})\|_\infty \leq \frac{1}{1000}c_e$ . Now, we note that again Lemma 8.3.9 shows  $\|\mathbf{e}(\boldsymbol{\tau}^{(t+1)}, \mathbf{x}^{(t+1)}) - \mathbf{e}(\boldsymbol{\tau}^{(t+\frac{2}{3})}, \mathbf{x}^{(t)})\|_2 \leq \frac{1}{10}c_\Delta \leq \frac{1}{1000}c_e$ , and we have  $\|\mathbf{e}(\boldsymbol{\tau}^{(t+1)}, \mathbf{x}^{(t+1)})\|_\infty \leq \frac{c_e}{400}$ . This finishes the induction case for  $\|\mathbf{e}\|_\infty$  and proves this lemma.  $\square$

Next, we show the number of constraints is always linear to  $n$ .

**Lemma 8.3.17.** *Throughout our cutting plane method, there are at most  $1 + \frac{2n}{c_d}$  constraints.*

*Proof.* We only add a constraint if  $\min_i w_i \geq c_d$ . Since  $2\psi_i \geq w_i$ , we have  $\psi_i \geq \frac{c_d}{2}$  for all  $i$ . Letting  $m$  denote the number of constraints after we add that row, we have  $n \geq \sum_i \psi_i \geq (m-1)(c_d/2)$ .  $\square$

Using  $K \neq \emptyset$  and  $K \subset B_\infty(R)$ , here we show that the points are bounded.

**Lemma 8.3.18.** *During our Cutting Plane Method, for all  $k$ , we have  $\|\mathbf{x}^{(k)}\|_2 \leq 6\sqrt{n/\lambda} + 2\sqrt{n}R$ .*

*Proof.* By Lemma 8.3.16 and Lemma 8.3.14 we know that  $\|\mathbf{x}^{(k)}\|_2^2 \leq 4\lambda^{-1}(mc_e + n) + 2\|\mathbf{y}\|_2^2$  for any  $\mathbf{y} \in P^{(k)}$ . Since our method never cuts out any point in  $K$  and since  $K$  is nonempty, there is some  $\mathbf{y} \in K \subset P^{(k)}$ . Since  $K \subset B_\infty(R)$ , we have  $\|\mathbf{y}\|_2^2 \leq nR^2$ . Furthermore, by Lemma 8.3.17 we have that  $mc_e \leq c_e + 2n \leq 3n$  yielding the result.  $\square$

**Lemma 8.3.19.**  *$s_i(\mathbf{x}^{(k)}) \leq 12\sqrt{n/\lambda} + 4\sqrt{n}R + \sqrt{\frac{1}{c_a\lambda}}$  for all  $i$  and  $k$  in the our cutting plane method.*

*Proof.* Let  $\mathbf{x}^{(j)}$  be the current point at the time that the constraint corresponding to  $s_i$ , denoted  $\{\mathbf{x} : \mathbf{a}_i^T \mathbf{x} \geq \mathbf{a}_i^T \mathbf{x}^{(j)} - s_i(\mathbf{x}^{(j)})\}$ , was added. Clearly

$$s_i(\mathbf{x}^{(k)}) = \mathbf{a}_i^T \mathbf{x}^{(k)} - \mathbf{a}_i^T \mathbf{x}^{(j)} + s_i(\mathbf{x}^{(j)}) \leq \|\mathbf{a}_i\| \cdot \|\mathbf{x}^{(k)}\| + \left| \mathbf{a}_i^T \mathbf{x}^{(j)} - s_i(\mathbf{x}^{(j)}) \right| .$$

On the one hand, if the constraint for  $s_i$  comes from the initial symmetric polytope  $P^{(0)} = B_\infty(R)$ , we know  $|\mathbf{a}^T \mathbf{x}^{(j)} - s_i(\mathbf{x}^{(j)})| \leq R$ . On the other hand, if the constraint was added later then we know that

$$s_i(\mathbf{x}^{(j)}) = c_a^{-1/2} \sqrt{\mathbf{a}^T (\mathbf{A}^T \mathbf{S}_{\mathbf{x}^{(j)}}^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a}} \leq (c_a \lambda)^{-1/2}$$

and  $|\mathbf{a}^T \mathbf{x}^{(j)} - s_i(\mathbf{x}^{(j)})| \leq \|\mathbf{a}_i\| \cdot \|\mathbf{x}^{(j)}\| + |s_i(\mathbf{x}^{(j)})|$ . Since  $\|\mathbf{a}_i\|_2 = 1$  by design and  $\|\mathbf{x}^{(j)}\|_2$  and  $\|\mathbf{x}^{(k)}\|_2$  are upper bounded by  $6\sqrt{n/\lambda} + 2\sqrt{n}R$  by Lemma 8.3.18, in either case the result follows.  $\square$

Now, we have everything we need to prove that the potential function is increasing in expectation.

**Lemma 8.3.20.** *Under the assumptions of Lemma 8.3.16 if  $\lambda = \frac{1}{c_a R^2}$ ,  $c_e = \frac{c_d}{4 \ln(6nR/\varepsilon)}$ , and  $24c_d \leq c_a \leq \frac{1}{1000}$  then for all  $k$  we have*

$$\mathbb{E} p_{e(\tau^{(k+1)}, \mathbf{x}^{(k+1)})}(\mathbf{x}^{(k+1)}) \geq p_{e(\tau^{(k)}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) - c_d + \ln(1 + \beta)$$

where  $\beta = c_a$  for the case of adding a constraint and  $\beta = -c_d$  for the case of removal.

*Proof.* Note that there are three places which affect the function value, namely the update step for  $\tau^{(k+\frac{1}{3})}$ , the addition/removal of constraints, and the centering step. We bound the effect of each separately.

First, for the update step, we have

$$p_{e(\tau^{(k+\frac{1}{3})}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) = (c + e_{i(k)}) \log(s_{i(k)}(\mathbf{x}^{(k)})/R) + p_{e(\tau^{(k)}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}).$$

Lemma 8.3.19, the termination condition,  $\lambda = \frac{1}{c_a R^2}$  and  $c_a < \frac{1}{1000}$  ensure that

$$\varepsilon \leq s_{i(k)}(\mathbf{x}^{(k)}) \leq 12\sqrt{n/\lambda} + 4\sqrt{n}R + \sqrt{\frac{1}{c_a \lambda}} \leq 6\sqrt{n}R \quad (8.19)$$

and Lemma 8.3.16 shows that  $|e_{i(k)}| \leq c_e$ . Hence, we have

$$p_{e(\tau^{(k+\frac{1}{3})}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) \geq p_{e(\tau^{(k)}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) + 2c_e \log(\varepsilon/R).$$

For the addition step, Lemma 8.3.12 shows that

$$\begin{aligned} p_{e(\tau^{(k+\frac{2}{3})}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) &= p_{e(\tau^{(k+\frac{1}{3})}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) - c_e \ln(s(\mathbf{x})_{m+1}/R) + \ln(1 + c_a) \\ &\geq p_{e(\tau^{(k+\frac{1}{3})}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) - c_e \log(6\sqrt{n}) + \ln(1 + c_a) \end{aligned}$$

and for the removal step, Lemma 8.3.13 and  $|e_i| \leq c_e$  shows that

$$\begin{aligned} p_{e(\tau^{(k+\frac{2}{3})}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) &\geq p_{e(\tau^{(k+\frac{1}{3})}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) + [c_e + e_P(\tau, \mathbf{x})_m] \ln(s(\mathbf{x})_m/R) + \ln(1 - c_d) \\ &\geq p_{e(\tau^{(k+\frac{1}{3})}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) - 2c_e \log(R/\varepsilon) + \ln(1 - c_d) \end{aligned}$$

After the addition or removal of a constraint, Lemma 8.3.16 shows that

$$\delta_{P^{(k)}, e(\tau^{(k+\frac{2}{3})}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) \leq \frac{1}{100} \sqrt{c_e + \min(\mu(\mathbf{x}^{(k)}), c_d)}$$

and therefore Lemma 8.3.9 and  $c_e \leq c_d$  show that

$$\begin{aligned} \mathbb{E}p_{e(\tau^{(k+1)}, \mathbf{x}^{(k+1)})}(\mathbf{x}^{(k+1)}) &\geq p_{e(\tau^{(k+\frac{2}{3})}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) - 8 \left( \frac{\sqrt{c_e + \min(\mu(\mathbf{x}^{(k)}), c_d)}}{100} \right)^2 \\ &\geq p_{e(\tau^{(k+\frac{2}{3})}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) - \frac{c_d}{625}. \end{aligned}$$

Combining them with  $c_e = \frac{c_d}{4 \ln(6nR/\varepsilon)}$ , we have

$$\begin{aligned} \mathbb{E}p_{e(\tau^{(k+1)}, \mathbf{x}^{(k+1)})}(\mathbf{x}^{(k+1)}) &\geq p_{e(\tau^{(k)}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) - 2c_e \log(6nR/\varepsilon) - \frac{c_d}{625} + \ln(1 + \beta) \\ &\geq p_{e(\tau^{(k)}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) - c_d + \ln(1 + \beta) \end{aligned}$$

where  $\beta = c_a$  for the case of addition and  $\beta = -c_d$  for the case of removal.  $\square$

**Theorem 8.3.21.** For  $c_a = \frac{1}{10^{10}}$ ,  $c_d = \frac{1}{10^{12}}$ ,  $c_e = \frac{c_d}{4 \ln(6nR/\varepsilon)}$ ,  $c_\Delta = \frac{C c_e}{\log(n \log(R/\varepsilon))}$  and  $\lambda = \frac{1}{c_a R^2}$  for some small enough universal constant  $C$ , then we have

$$\mathbb{E}p_{e(\tau^{(k+1)}, \mathbf{x}^{(k+1)})}(\mathbf{x}^{(k+1)}) \geq p_{e(\tau^{(k)}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) - \frac{1}{10^{11}} + \frac{9\beta}{10^{11}}$$

where  $\beta = 1$  for the case of addition and  $\beta = 0$  for the case of removal.

*Proof.* It is easy to see that these parameters satisfy the requirements of Lemma 8.3.20.  $\square$

### ■ 8.3.5 Guarantees of the Algorithm

In this section we put everything together to prove Theorem 8.3.26, the main result of this section, providing the guarantees of our cutting plane method.

For the remainder of this section we assume that  $c_a = \frac{1}{10^{10}}$ ,  $c_d = \frac{1}{10^{12}}$ ,  $c_e = \frac{c_d}{4 \ln(6nR/\varepsilon)}$ ,  $c_\Delta = \frac{C c_e}{\log(n \ln(R/\varepsilon))}$  and  $\lambda = \frac{1}{c_a R^2}$ . Consequently, throughout the algorithm we have

$$\|\mathbf{x}\|_2 \leq 6\sqrt{n/\lambda} + 2\sqrt{n}R \leq 3\sqrt{n}R. \quad (8.20)$$

**Lemma 8.3.22.** If  $s_i(\mathbf{x}^{(k)}) < \varepsilon$  for some  $i$  and  $k$  during our Cutting Plane Method then

$$\max_{\mathbf{y} \in P^{(k)} \cap B_\infty(R)} \langle \mathbf{a}_i, \mathbf{y} \rangle - \min_{\mathbf{y} \in P^{(k)} \cap B_\infty(R)} \langle \mathbf{a}_i, \mathbf{y} \rangle \leq \frac{8n\varepsilon}{c_a c_e}.$$

*Proof.* Let  $\mathbf{y} \in P^{(k)} \cap B_\infty(R)$  be arbitrary. Since  $\mathbf{y} \in B_\infty(R)$  clearly  $\|\mathbf{y}\|_2^2 \leq nR^2$ . Furthermore, by Lemma 8.3.17 and the choice of parameters  $mc_e + n \leq 3n$ . Consequently, by Lemma 8.3.14 and the fact that  $\lambda = \frac{1}{c_a R^2}$  and  $c_a < 1$  we have

$$\|\mathbf{x} - \mathbf{y}\|_{H(\mathbf{x})} \leq \frac{12mc_e + 6n + 2\lambda\|\mathbf{y}\|_2^2}{\sqrt{c_e + \mu(\mathbf{x})}} \leq \frac{30n + 2\frac{n}{c_a}}{\sqrt{c_e + \mu(\mathbf{x})}} \leq \frac{4n}{c_a \sqrt{c_e}}$$

and therefore

$$\left\| \mathbf{S}_{\mathbf{x}^{(k)}}^{-1}(s(\mathbf{x}^{(k)}) - s(\mathbf{y})) \right\|_\infty \leq \frac{1}{\sqrt{c_e}} \left\| \mathbf{S}_{\mathbf{x}^{(k)}}^{-1}(s(\mathbf{x}^{(k)}) - s(\mathbf{y})) \right\|_{c_e \mathbf{I} + \Psi} \leq \frac{4n}{c_a c_e}.$$

Consequently, we have  $(1 - \frac{4n}{c_a c_e})s_i(\mathbf{x}^{(k)}) \leq s_i(\mathbf{y}) \leq (1 + \frac{4n}{c_a c_e})s_i(\mathbf{x}^{(k)})$  for all  $\mathbf{y} \in P^{(k)} \cap B_\infty(R)$ .  $\square$

Now let us show how to compute a proof (or certificate) that the feasible region has small width on the direction  $\mathbf{a}_i$ .

**Lemma 8.3.23.** *Suppose that during some iteration  $k$  for  $i = \arg \min_j s_j(\mathbf{x}^{(k)})$  we have  $s_i(\mathbf{x}^{(k)}) \leq \varepsilon$ . Let  $(\mathbf{x}_*, \boldsymbol{\tau}_*) = \text{Centering}(\mathbf{x}^{(k)}, \boldsymbol{\tau}^{(k)}, 64 \log(2R/\varepsilon), c_\Delta)$  where  $\boldsymbol{\tau}^{(k)}$  is the  $\boldsymbol{\tau}$  at that point in the algorithm and let*

$$\mathbf{a}^* = \sum_{j \neq i} t_j \mathbf{a}_j \text{ where } t_j = \left( \frac{s(\mathbf{x}_*)_i}{s(\mathbf{x}_*)_j} \right) \left( \frac{c_e + e_j(\mathbf{x}_*, \boldsymbol{\tau}_*) + \psi_j(\mathbf{x}_*)}{c_e + e_i(\mathbf{x}_*, \boldsymbol{\tau}_*) + \psi_i(\mathbf{x}_*)} \right).$$

Then, we have that  $\|\mathbf{a}_i + \mathbf{a}^*\|_2 \leq \frac{8\sqrt{n}\varepsilon}{c_a c_e R}$  and  $t_j \geq 0$  for all  $j$ . Furthermore, we have

$$\left( \sum_{j \neq i}^{O(n)} t_j \mathbf{a}_j \right)^T \mathbf{x}_* - \sum_{j \neq i}^{O(n)} t_j b_j \leq \frac{4n\varepsilon}{c_e}.$$

*Proof.* By Lemma 8.3.9 and Lemma 8.3.16 we know that  $e(\mathbf{x}_*, \boldsymbol{\tau}_*) \leq \frac{1}{2}c_e$  and  $\delta_{e(\mathbf{x}_*, \boldsymbol{\tau}_*)} \leq \frac{\varepsilon}{R} \sqrt{c_e + \mu(\mathbf{x}_*)}$ . Since  $e(\mathbf{x}_*, \boldsymbol{\tau}_*) \leq \frac{1}{2}c_e$ , we have  $t_j \geq 0$  for all  $j$ . Furthermore, by Lemma 8.3.15 and (8.20), we then have that with high probability in  $n$

$$\begin{aligned} \|\mathbf{a}_i + \mathbf{a}^*\|_2 &\leq \frac{2\varepsilon}{(c_e + \mu(\mathbf{x}_*))} \left[ \lambda \|\mathbf{x}_*\|_2 + \delta_e(\mathbf{x}_*) \sqrt{\frac{mc_e + n}{\varepsilon^2} + \lambda} \right] \\ &\leq \frac{2\varepsilon}{c_e} \left[ \frac{1}{c_a R^2} (3\sqrt{n}R) + \frac{2\varepsilon}{R} \sqrt{\frac{3n}{\varepsilon^2} + \frac{1}{c_a R^2}} \right] \\ &\leq \frac{2\varepsilon}{c_e} \left[ \frac{3\sqrt{n}}{c_a R} + \frac{2\sqrt{3n}}{R} + \frac{2}{\sqrt{c_a} R} \right] \leq \frac{8\sqrt{n}\varepsilon}{c_a c_e R}. \end{aligned}$$

By Lemma 8.3.16 we know that  $e(\mathbf{x}_*, \boldsymbol{\tau}_*) \leq \frac{1}{2}c_e$  and hence

$$\begin{aligned} \left( \sum_{j \neq i}^{O(n)} t_j \mathbf{a}_j \right)^T \mathbf{x}_* - \sum_{j \neq i}^{O(n)} t_j b_j &= \sum_{j \neq i}^{O(n)} t_j s(\mathbf{x}_*)_j = s_i(\mathbf{x}_*) \sum_{j \neq i}^{O(n)} \left( \frac{c_e + e_j(\mathbf{x}_*, \boldsymbol{\tau}_*) + \psi_j(\mathbf{x}_*)}{c_e + e_i(\mathbf{x}_*, \boldsymbol{\tau}_*) + \psi_i(\mathbf{x}_*)} \right) \\ &\leq s_i(\mathbf{x}_*) \sum_{j \neq i}^{O(n)} \left( \frac{\frac{3}{2}c_e + \psi_j(\mathbf{x}_*)}{\frac{1}{2}c_e + \psi_i(\mathbf{x}_*)} \right) \leq s_i(\mathbf{x}_*) \frac{3mc_e + 2n}{c_e} \\ &\leq \frac{3n}{c_e} s_i(\mathbf{x}_*) \leq \frac{4n\varepsilon}{c_e} \end{aligned}$$

where the last line follows the fact that  $s_i(\mathbf{x}_*) \leq 1.1s_i(\mathbf{x}^{(k)}) \leq 1.1\varepsilon$  (Lemma 8.3.9).  $\square$

**Lemma 8.3.24.** *During our Cutting Plane Method, if  $p_e(\mathbf{x}^{(k)}) \geq 2n \log(\frac{nR}{c_a \varepsilon}) + \frac{6n}{c_a}$ , then we have  $s_i(\mathbf{x}^{(k)}) \leq \varepsilon$  for some  $i$ .*

*Proof.* Recall that

$$p_e(\mathbf{x}^{(k)}) = - \sum_{i \in [m]} (c_e + e_i) \log \left( s_i(\mathbf{x}^{(k)})_i / R \right) + \frac{1}{2} \log \det \left( R^2 \left( \mathbf{A}^T \mathbf{S}_{\mathbf{x}^{(k)}}^{-2} \mathbf{A} + \lambda \mathbf{I} \right) \right) + \frac{\lambda}{2} \|\mathbf{x}^{(k)}\|_2^2.$$

Using  $\|\mathbf{x}^{(k)}\| \leq 3\sqrt{n}R$  (8.20) and  $\lambda = \frac{1}{c_a R^2}$ , we have

$$p_e(\mathbf{x}^{(k)}) \leq - \sum_{i \in [m]} (c_e + e_i) \log(s(\mathbf{x}^{(k)})_i / R) + \frac{1}{2} \log \det \left( R^2 \left( \mathbf{A}^T \mathbf{S}_{\mathbf{x}^{(k)}}^{-2} \mathbf{A} + \lambda \mathbf{I} \right) \right) + \frac{5n}{c_a}.$$

Next, we note that  $\|e_i\|_\infty \leq c_e \leq \frac{1}{4 \ln(6nR/\varepsilon)}$  and  $s_i(\mathbf{x}^{(k)}) \leq 12\sqrt{n/\lambda} + 4\sqrt{n}R + \sqrt{\frac{1}{c_a \lambda}} \leq 6\sqrt{n}R$  (Lemma 8.3.19). Hence, we have

$$p_e(\mathbf{x}^{(k)}) \leq \frac{1}{2} \log \det \left( R^2 \left( \mathbf{A}^T \mathbf{S}_{\mathbf{x}^{(k)}}^{-2} \mathbf{A} + \lambda \mathbf{I} \right) \right) + \frac{6n}{c_a} - \frac{n}{2} \log \left( \frac{\min_i s_i}{R} \right).$$

We prove  $s_i(\mathbf{x}^{(k)}) \leq \varepsilon$  by contradiction. Since  $p_e(\mathbf{x}^{(k)}) \geq 2n \log(\frac{nR}{c_a \varepsilon}) + \frac{6n}{c_a}$  and  $s_i(\mathbf{x}^{(k)}) > \varepsilon$  for all  $i$ , we have that  $\frac{1}{2} \log \det \left( R^2 \left( \mathbf{A}^T \mathbf{S}_{\mathbf{x}^{(k)}}^{-2} \mathbf{A} + \lambda \mathbf{I} \right) \right) \geq n \log(\frac{nR}{c_a \varepsilon})$ . Using  $\varepsilon < R$ , we have that

$$\sum_i \log \lambda_i \left( R^2 \left( \mathbf{A}^T \mathbf{S}_{\mathbf{x}^{(k)}}^{-2} \mathbf{A} + \lambda \mathbf{I} \right) \right) \geq n \log \left( \frac{n^2 R^2}{c_a^2 \varepsilon^2} \right) \geq n \log \left( \frac{nR^2}{2c_a^2 \varepsilon^2} + R^2 \lambda \right).$$

Therefore, we have  $\log \lambda_{\max} \left( R^2 \left( \mathbf{A}^T \mathbf{S}_{\mathbf{x}^{(k)}}^{-2} \mathbf{A} + \lambda \mathbf{I} \right) \right) \geq \log \left( \frac{nR^2}{2c_a^2 \varepsilon^2} + R^2 \lambda \right)$ . Hence, we have some unit vector  $\mathbf{v}$  such that  $\mathbf{v} \mathbf{A}^T \mathbf{S}_{\mathbf{x}^{(k)}}^{-2} \mathbf{A} \mathbf{v} + \lambda \mathbf{v}^T \mathbf{v} \geq \frac{n}{2c_a^2 \varepsilon^2} + \lambda$ . Thus,

$$\sum_i \frac{(\mathbf{A} \mathbf{v})_i^2}{s(\mathbf{x}^{(k)})_i^2} \geq \frac{n}{2c_a^2 \varepsilon^2}.$$

Lemma 8.3.17 shows that the number of constraints is bounded by  $1 + 2n/c_d \leq n/(2c_a^2)$  and hence there is some  $i$  such that  $\frac{(\mathbf{A} \mathbf{v})_i^2}{s(\mathbf{x}^{(k)})_i^2} \geq \frac{1}{\varepsilon^2}$ . Since  $\mathbf{a}_i$  and  $\mathbf{v}$  are unit vectors, we have  $1 \geq \langle \mathbf{a}_i, \mathbf{v} \rangle^2 \geq s(\mathbf{x}^{(k)})_i^2 / \varepsilon^2$  and hence  $s(\mathbf{x}^{(k)})_i \leq \varepsilon$  (contradiction).  $\square$

**Lemma 8.3.25.** *With constant probability, the algorithm ends in  $10^{24}n \log(\frac{nR}{\varepsilon})$  iterations.*<sup>1</sup>

*Proof.* Theorem 8.3.21 shows that for all  $k$

$$\mathbb{E} p_{e(\tau^{(k+1)}, \mathbf{x}^{(k+1)})}(\mathbf{x}^{(k+1)}) \geq p_{e(\tau^{(k)}, \mathbf{x}^{(k)})}(\mathbf{x}^{(k)}) - \frac{1}{10^{11}} + \frac{9\beta}{10^{11}} \quad (8.21)$$

where  $\beta = 1$  for the case of adding a constraint and  $\beta = 0$  for the case of removing a constraint. Now, for all  $t$  consider the random variable

$$\mathbf{X}_t = p_{e(\tau^{(t)}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) - \frac{4.5m^{(t)}}{10^{11}} - \frac{3.5t}{10^{11}}$$

where  $m^{(t)}$  is the number of constraints in iteration  $t$  of the algorithm. Then, since  $m^{(t+1)} = m^{(t)} - 1 + 2\beta$ , (8.21) shows that

$$\begin{aligned} \mathbf{E} \mathbf{X}_{t+1} &\geq p_{e(\tau^{(t)}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) - \frac{1}{10^{11}} + \frac{9\beta}{10^{11}} - \frac{4.5m^{(t+1)}}{10^{11}} - \frac{3.5(t+1)}{10^{11}} \\ &= \mathbf{X}_t - \frac{1}{10^{11}} + \frac{9\beta}{10^{11}} - \frac{4.5(-1+2\beta)}{10^{11}} - \frac{3.5}{10^{11}} = \mathbf{X}_t. \end{aligned}$$

<sup>1</sup>We have made no effort on improving this constant and we believe it can be improved to less than 300 using techniques in [11, 14].

Hence, it is a sub-martingale. Let  $T$  be the iteration the algorithm outputs  $\mathbf{x}^{(k)}$  or  $P^{(k)}$ . Optional stopping theorem shows that

$$\mathbb{E}\mathbf{X}_{\min(T,t)} \geq \mathbb{E}\mathbf{X}_0. \quad (8.22)$$

Since the polytope is  $B_\infty(R)$ , we have

$$\begin{aligned} p_0(\mathbf{0}) &= - \sum_{i \in [m^{(0)}]} c_e \log(s_i(\mathbf{0})/R) + \frac{1}{2} \log \det (R^2 (\mathbf{A}^T \mathbf{S}_0^{-2} \mathbf{A} + \lambda \mathbf{I})) + \frac{\lambda}{2} \|\mathbf{0}\|_2^2 \\ &\geq \frac{1}{2} \log \det \left( R^2 \left( \frac{2}{R^2} \mathbf{I} + \frac{1}{c_a R^2} \mathbf{I} \right) \right) \\ &\geq -15n \end{aligned}$$

where we used  $c_a = 1/10^{10}$  on the last line. Hence, we have

$$\begin{aligned} \mathbf{X}_0 &\geq -15n - \frac{4.5m^{(0)}}{10^{11}} \\ &\geq -20n. \end{aligned}$$

Therefore, (8.22) shows that for all  $t$  we have

$$\begin{aligned} -20n &\leq \mathbb{E}\mathbf{X}_{\min(T,t)} \\ &= p \mathbb{E} [\mathbf{X}_{\min(T,t)} | T < t] + (1-p) \mathbb{E} [\mathbf{X}_{\min(T,t)} | T \geq t] \end{aligned} \quad (8.23)$$

where  $p \stackrel{\text{def}}{=} \mathbb{P}(T < t)$ .

Note that

$$\begin{aligned} \mathbb{E} [\mathbf{X}_{\min(T,t)} | T \geq t] &\leq \mathbb{E} \left[ p_{e(\tau^{(t)}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) | T \geq t \right] - \frac{4.5m^{(t)}}{10^{11}} - \frac{3.5t}{10^{11}} \\ &\leq \mathbb{E} \left[ p_{e(\tau^{(t)}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) | T \geq t \right] - \frac{3.5t}{10^{11}}. \end{aligned}$$

Furthermore, by Lemma 8.3.24 we know that when  $p_{e(\tau^{(t)}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) \geq 2n \log(\frac{nR}{c_a \varepsilon}) + \frac{6n}{c_a}$ , there is a slack that is too small and the algorithm terminates. Hence, we have

$$\mathbb{E} [\mathbf{X}_{\min(T,t)} | T \geq t] \leq 2n \log\left(\frac{nR}{c_a \varepsilon}\right) + \frac{6n}{c_a} - \frac{3.5t}{10^{11}}.$$

The proof of Lemma 8.3.16 shows that the function value does not change by more than 1 in one iteration by changing  $\mathbf{x}$  and can change by at most  $mc_e \log(\frac{3nR}{\varepsilon})$  by changing  $\tau$ . Since by Lemma 8.3.17 we know that  $m \leq 1 + \frac{2n}{c_a}$  and  $c_e = \frac{c_d}{4 \ln(6nR/\varepsilon)}$ , we have that  $p_e(\mathbf{x}) \leq 2n \log(\frac{nR}{c_a \varepsilon}) + \frac{7n}{c_a}$  throughout the execution of the algorithm. Therefore, we have

$$\mathbb{E} [\mathbf{X}_{\min(T,t)} | T \leq t] \leq \mathbb{E}_{T < t} p_{e(\tau^{(t)}, \mathbf{x}^{(t)})}(\mathbf{x}^{(t)}) \leq 2n \log\left(\frac{nR}{c_a \varepsilon}\right) + \frac{7n}{c_a}.$$

Therefore, (8.23) shows that

$$-20n \leq 2n \log\left(\frac{nR}{c_a \varepsilon}\right) + \frac{7n}{c_a} - (1-p) \frac{3.5t}{10^{11}}.$$



Hence, we have

$$\begin{aligned}
(1-p)\frac{3.5t}{10^{11}} &\leq 2n \log\left(\frac{nR}{c_a \varepsilon}\right) + \frac{7n}{c_a} + 20n \\
&\leq 2n \log\left(\frac{nR}{c_a \varepsilon}\right) + \frac{7n}{c_a} + 20n \\
&= 2n \log\left(\frac{Rn}{c_a \varepsilon}\right) + 8 \cdot 10^{10}n.
\end{aligned}$$

Thus, we have

$$\mathbb{P}(T < t) = p \geq 1 - \frac{1}{t} \left( 10^{11}n \log\left(\frac{nR}{\varepsilon}\right) + 10^{22}n \right).$$

□

Now, we gather all the result as follows:

**Theorem 8.3.26** (Our Cutting Plane Method). *Let  $K \subseteq \mathbb{R}^n$  be a non-empty set contained in a box of radius  $R$ , i.e.  $K \subseteq B_\infty(R)$ . For any  $\varepsilon \in (0, R)$  in expected time  $O(n\text{SO}_{\Omega(\varepsilon/\sqrt{n})}(K) \log(nR/\varepsilon) + n^3 \log^{O(1)}(nR/\varepsilon))$  our cutting plane method either outputs  $\mathbf{x} \in K$  or finds a polytope  $P = \{\mathbf{x} : \mathbf{A}\mathbf{x} \geq \mathbf{b}\} \supseteq K$  such that*

1.  $P$  has  $O(n)$  many constraints (i.e.  $\mathbf{A} \in \mathbb{R}^{O(n) \times n}$  and  $\mathbf{b} \in \mathbb{R}^{O(n)}$ ).
2. Each constraint of  $P$  is either an initial constraint from  $B_\infty(R)$  or of the form  $\langle \mathbf{a}, \mathbf{x} \rangle \geq b - \delta$  where  $\langle \mathbf{a}, \mathbf{x} \rangle \geq b$  is a normalized hyperplane (i.e.  $\|\mathbf{a}\|_2 = 1$ ) returned by the separation oracle and  $\delta = \Omega\left(\frac{\varepsilon}{\sqrt{n}}\right)$ .
3. The polytope  $P$  has small width with respect to some direction  $\mathbf{a}_1$  given by one of the constraints, i.e.

$$\max_{\mathbf{y} \in P \cap B_\infty(R)} \langle \mathbf{a}_1, \mathbf{y} \rangle - \min_{\mathbf{y} \in P \cap B_\infty(R)} \langle \mathbf{a}_1, \mathbf{y} \rangle \leq O(n\varepsilon \ln(nR/\varepsilon))$$

4. Furthermore, the algorithm produces a proof of the fact above involving convex combination of the constraints, namely, non-negatives  $t_2, \dots, t_{O(n)}$  and  $\mathbf{x} \in P$  such that

- (a)  $\|\mathbf{x}\|_2 = O(\sqrt{n}R)$ ,
- (b)  $\left\| \mathbf{a}_1 + \sum_{i=2}^{O(n)} t_i \mathbf{a}_i \right\|_2 = O\left(\frac{\varepsilon}{R} \sqrt{n} \log(nR/\varepsilon)\right)$ ,
- (c)  $\mathbf{a}_1^T \mathbf{x} - b_1 \leq \varepsilon$ ,
- (d)  $\left( \sum_{i=2}^{O(n)} t_i \mathbf{a}_i \right)^T \mathbf{x} - \sum_{i=2}^{O(n)} t_i b_i \leq O(n\varepsilon \log(nR/\varepsilon))$  .

*Proof.* Our algorithm either finds  $\mathbf{x} \in K$  or we have  $s_i(\mathbf{x}^{(k)}) < \varepsilon$ . When  $s_i(\mathbf{x}^{(k)}) < \varepsilon$ , we apply Lemma 8.3.23 to construct the polytope  $P$  and the linear combination  $\sum_{i=2}^{O(n)} t_i \mathbf{a}_i$ .

Notice that each iteration of our algorithm needs to solve constant number of linear systems and implements the sampling step to find  $\Delta^{(k)} \in \mathbb{R}^n$  s.t.  $\mathbf{E}[\Delta^{(k)}] = \boldsymbol{\psi}(\mathbf{x}^{(k)}) - \boldsymbol{\psi}(\mathbf{x}^{(k-1)})$ . Theorem 8.4.2 shows how to do the sampling in  $\tilde{O}(1)$  many linear systems. Hence, in total, each iterations needs to solve  $\tilde{O}(1)$  many linear systems plus nearly linear work. To output the proof for (4), we use Lemma 8.3.23.

Note that the linear systems the whole algorithm need to solve is of the form

$$(\mathbf{A}^T \mathbf{S}_x^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{x} = \mathbf{y}.$$

where the matrix  $\mathbf{A}^T \mathbf{S}_x^{-2} \mathbf{A} + \lambda \mathbf{I}$  can be written as  $\overline{\mathbf{A}}^T \mathbf{D} \overline{\mathbf{A}}$  for the matrix  $\overline{\mathbf{A}} = [\mathbf{A} \ \mathbf{I}]$  and diagonal matrix

$$\mathbf{D} = \begin{bmatrix} \mathbf{S}^{-2} & \mathbf{0} \\ \mathbf{0} & \lambda \mathbf{I} \end{bmatrix}.$$

Note that Lemma 8.3.9 shows that  $\|(\mathbf{S}^{(k)})^{-1} (\mathbf{s}^{(k+1)} - \mathbf{s}^{(k)})\|_2 \leq \frac{1}{10}$  for the  $k^{\text{th}}$  and  $(k+1)^{\text{th}}$  linear systems we solved in the algorithm. Hence, we have  $\|(\mathbf{D}^{(k)})^{-1} (\mathbf{d}^{(k+1)} - \mathbf{d}^{(k)})\|_2 \leq \frac{1}{10}$ . In [164], they showed how to solve such sequence of systems in  $\tilde{O}(n^2)$  amortized cost. Moreover, since our algorithm always changes the constraints by  $\delta$  amount where  $\delta = \Omega(\frac{\varepsilon}{\sqrt{n}})$  an inexact separation oracle  $\text{SO}_{\Omega(\varepsilon/\sqrt{n})}$  suffices. (see Def 2.3.5). Consequently, the total work  $O(n \text{SO}_{\Omega(\varepsilon/\sqrt{n})}(K) \log(nR/\varepsilon) + n^3 \log^{O(1)}(nR/\varepsilon))$ . Note that as the running time holds with only constant probability, we can restart the algorithm whenever the running time is too large.

To prove (2), we note that from the algorithm description, we know the constraints are either from  $B_\infty(R)$  or of the form  $\mathbf{a}^T \mathbf{x} \geq \mathbf{a}^T \mathbf{x}^{(k)} - \delta$  where

$$\delta = \sqrt{\frac{\mathbf{a}^T (\mathbf{A}^T \mathbf{S}_{x^{(k)}}^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a}}{c_a}}.$$

From the proof of Lemma 8.3.24, we know that if  $\lambda_{\max}(\mathbf{A}^T \mathbf{S}_x^{-2} \mathbf{A} + \lambda \mathbf{I}) \geq \frac{n}{c_a^2 \varepsilon^2}$ , then there is  $s_i < \varepsilon$ . Hence, we have  $\lambda_{\min}((\mathbf{A}^T \mathbf{S}_x^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1}) \geq \frac{c_a^2 \varepsilon^2}{n}$ . Since  $\mathbf{a}$  is a unit vector, we have

$$\sqrt{\frac{\mathbf{a}^T (\mathbf{A}^T \mathbf{S}_{x^{(k)}}^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a}}{c_a}} \geq \sqrt{\frac{c_a \varepsilon^2}{n}}.$$

□

## ■ 8.4 Technical Tools

In this section we provide stand-alone technical tools we use in our cutting plane method in Section 8.3. In Section 8.4.1 we show how to efficiently compute accurate estimates of changes in leverage scores using access to a linear system solver. In Section 8.4.2 we study what we call the “Stochastic Chasing 0 Game” and show how to maintain that a vector is small in  $\ell_\infty$  norm by making small coordinate updates while the vector changes randomly in  $\ell_2$ .

### ■ 8.4.1 Estimating Changes in Leverage Scores

In previous sections, we needed to compute leverage scores accurately and efficiently for use in our cutting plane method. Note that the leverage score definition we used was

$$\psi(\mathbf{w})_i = \mathbf{1}_i^T \sqrt{\overline{\mathbf{W}}} \mathbf{A} (\mathbf{A}^T \mathbf{W} \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \sqrt{\overline{\mathbf{W}}} \mathbf{1}_i$$

for some  $\lambda > 0$  which is different from the standard definition

$$\sigma(\mathbf{w})_i = \mathbf{1}_i^T \sqrt{\overline{\mathbf{W}}} \mathbf{A} (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \sqrt{\overline{\mathbf{W}}} \mathbf{1}_i.$$

However, note that the matrix  $\mathbf{A}^T \mathbf{W} \mathbf{A} + \lambda \mathbf{I}$  can be written as  $\overline{\mathbf{A}}^T \mathbf{D} \overline{\mathbf{A}}$  for the matrix  $\overline{\mathbf{A}}^T = [\mathbf{A}^T \ \mathbf{I}]$  and diagonal matrix

$$\mathbf{D} = \begin{bmatrix} \mathbf{W} & \mathbf{0} \\ \mathbf{0} & \lambda \mathbf{I} \end{bmatrix}.$$

and therefore computing  $\psi$  is essentially same as computing typical leverage scores. Consequently, we use the standard definition  $\sigma$  to simplify notation.

We explained in Section 2.3.1 that leverage scores can be approximated efficiently using dimension reduction. Unfortunately, the error incurred by this approximation is too large to use inside the cutting point method. In this section, we show how to efficiently approximate the *change* of leverage score more accurately.

In particular, we show how to approximate  $\sigma(\mathbf{w}) - \sigma(\mathbf{v})$  for any given  $\mathbf{w}, \mathbf{v}$  with  $\|\log(\mathbf{w}) - \log(\mathbf{v})\|_2 \ll 1$ . Our algorithm breaks  $\sigma(\mathbf{w})_i - \sigma(\mathbf{v})_i$  into the sum of the norm of small vectors and then uses the Johnson-Lindenstrauss dimension reduction to approximate the norm of each vector separately. Our algorithm makes use of the following version of Johnson-Lindenstrauss.

**Lemma 8.4.1** ([2]). *Let  $0 \leq \varepsilon \leq \frac{1}{2}$  and let  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  be arbitrary  $m$  points. For  $k = O(\varepsilon^{-2} \log(m))$  let  $\mathbf{Q}$  be a  $k \times n$  random matrix with each entry sampled from  $\{-\frac{1}{\sqrt{k}}, \frac{1}{\sqrt{k}}\}$  uniformly and independently. Then,  $\mathbb{E} \|\mathbf{Q} \mathbf{x}_i\|^2 = \|\mathbf{x}_i\|^2$  for all  $i \in [m]$  and with high probability in  $m$  we have that for all  $i \in [m]$*

$$(1 - \varepsilon) \|\mathbf{x}_i\|^2 \leq \|\mathbf{Q} \mathbf{x}_i\|^2 \leq (1 + \varepsilon) \|\mathbf{x}_i\|^2 \quad .$$

---

**Algorithm 23:**  $\widehat{h} = \text{LeverageChange}(\mathbf{A}, \mathbf{v}, \mathbf{w}, \varepsilon)$

---

**Input:**  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{v}, \mathbf{w} \in \mathbb{R}_{>0}^m$ ,  $\varepsilon \in (0, 0.5)$ .

**Given:**  $\|\mathbf{V}^{-1}(\mathbf{v} - \mathbf{w})\|_2 \leq \frac{1}{10}$ ;  $\mathbf{A}^T \mathbf{V} \mathbf{A}$  and  $\mathbf{A}^T \mathbf{W} \mathbf{A}$  are invertible.

Sample  $\mathbf{Q}_d \in \mathbb{R}^{O(\varepsilon^{-2} \log(m)) \times n}$  as in Lemma 8.4.1.

Let  $\hat{d}_i = \|\mathbf{Q}_d \sqrt{\mathbf{W} \mathbf{A}} (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{1}_i\|_2^2$  for all  $i \in [n]$ .

Let  $t = O(\log(\varepsilon^{-1}))$ .

Sample  $\mathbf{Q}_f \in \mathbb{R}^{O(\varepsilon^{-2} \log(mt)) \times n}$  as in Lemma 8.4.1.

Pick positive integer  $u$  randomly such that  $\Pr[u = i] = (\frac{1}{2})^i$ .

**for**  $j \in \{1, 2, \dots, t\} \cup \{t + u\}$  **do**

**if**  $j$  is even **then**

        Let  $\hat{f}_i^{(j)} = \|\mathbf{Q}_f \sqrt{\mathbf{V} \mathbf{A}} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} (\mathbf{A}^T (\mathbf{V} - \mathbf{W}) \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1})^{\frac{j}{2}} \mathbf{A}^T \mathbf{1}_i\|_2^2$ .

**else**

        Let  $\Delta^+ \stackrel{\text{def}}{=} (\mathbf{V} - \mathbf{W})^+$ , i.e. the matrix  $\mathbf{V} - \mathbf{W}$  with negative entries set to 0.

        Let  $\Delta^- \stackrel{\text{def}}{=} (\mathbf{W} - \mathbf{V})^+$ , i.e. the matrix  $\mathbf{W} - \mathbf{V}$  with negative entries set to 0.

        Let  $\hat{\alpha}_i^{(j)} = \|\mathbf{Q}_f \sqrt{\Delta^+} \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} (\mathbf{A}^T (\mathbf{V} - \mathbf{W}) \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1})^{\frac{j-1}{2}} \mathbf{A}^T \mathbf{1}_i\|_2^2$ .

        Let  $\hat{\beta}_i^{(j)} = \|\mathbf{Q}_f \sqrt{\Delta^-} \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} (\mathbf{A}^T (\mathbf{V} - \mathbf{W}) \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1})^{\frac{j-1}{2}} \mathbf{A}^T \mathbf{1}_i\|_2^2$ .

        Let  $\hat{f}_i^{(j)} = \hat{\alpha}_i^{(j)} - \hat{\beta}_i^{(j)}$ .

**end**

**end**

Let  $\hat{f}_i = 2^u \hat{f}_i^{(t+u)} + \sum_{j=1}^t \hat{f}_i^{(j)}$ .

**Output:**  $\hat{h}_i = (w_i - v_i) \hat{d}_i + v_i \hat{f}_i$ . for all  $i \in [m]$

---

**Theorem 8.4.2.** *Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{v}, \mathbf{w} \in \mathbb{R}_{>0}^m$  be such that  $\alpha \stackrel{\text{def}}{=} \|\mathbf{V}^{-1}(\mathbf{v} - \mathbf{w})\|_2 \leq \frac{1}{10}$  and both  $\mathbf{A}^T \mathbf{V} \mathbf{A}$  and  $\mathbf{A}^T \mathbf{W} \mathbf{A}$  are invertible. For any  $\varepsilon \in (0, 0.5)$ , Algorithm 23 generates a random variable  $\widehat{h}$*

such that  $\mathbf{E}\hat{\mathbf{h}} = \sigma(\mathbf{w}) - \sigma(\mathbf{v})$  and with high probability in  $m$ , we have  $\|\hat{\mathbf{h}} - (\sigma(\mathbf{w}) - \sigma(\mathbf{v}))\|_2 \leq O(\alpha\varepsilon)$ . Furthermore, the expected running time is  $\tilde{O}((\text{nnz}(\mathbf{A}) + \text{LO})/\varepsilon^2)$  where  $\text{LO}$  is the amount of time needed to apply  $(\mathbf{A}^T\mathbf{V}\mathbf{A})^{-1}$  and  $(\mathbf{A}^T\mathbf{W}\mathbf{A})^{-1}$  to a vector.

*Proof.* First we bound the running time. To compute  $\hat{d}_i, \hat{f}_i^{(j)}, \hat{\alpha}_i^{(j)}, \hat{\beta}_i^{(j)}$ , we simply perform matrix multiplications from the left and then consider the dot products with each of the rows of  $\mathbf{A}$ . Naively this would take time  $\tilde{O}((t+u)^2 \log(mt)(\text{nnz}(\mathbf{A}) + \text{LO}))$ . However, we can reuse the computation in computing high powers of  $j$  to only take time  $\tilde{O}((t+u) \log(mt)(\text{nnz}(\mathbf{A}) + \text{LO}))$ . Now since  $\mathbf{E}[u]$  is constant we see that the total running time is as desired. It only remains to prove the desired properties of  $\hat{\mathbf{h}}$ .

First we note that we can re-write leverage score differences using

$$\sigma(\mathbf{w})_i - \sigma(\mathbf{v})_i = (w_i - v_i) \left[ \mathbf{A} (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \right]_{ii} + v_i \left[ \mathbf{A} \left( (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} - (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} \right) \mathbf{A}^T \right]_{ii} .$$

Consequently, for all  $i \in [m]$ , if we let

$$\begin{aligned} d_i &\stackrel{\text{def}}{=} \mathbf{1}_i^T \mathbf{A} (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{1}_i, \\ f_i &\stackrel{\text{def}}{=} \mathbf{1}_i^T \mathbf{A} \left[ (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} - (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} \right] \mathbf{A}^T \mathbf{1}_i. \end{aligned}$$

then

$$\sigma(\mathbf{w})_i - \sigma(\mathbf{v})_i = (w_i - v_i)d_i + (v_i)f_i . \quad (8.24)$$

We show that  $\hat{d}_i$  approximates  $d$  and  $\hat{f}_i$  approximate  $f$  well enough to satisfy the statements in the Theorem.

First we bound the quality of  $\hat{d}_i$ . Note that  $d_i = \|\sqrt{\mathbf{W}}\mathbf{A} (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{1}_i\|_2^2$ . Consequently, Lemma 8.4.1 shows that  $\mathbf{E}[\hat{d}_i] = d_i$  and that with high probability in  $m$  we have  $(1-\varepsilon)d_i \leq \hat{d}_i \leq (1+\varepsilon)d_i$  for all  $i \in [m]$ . Therefore, with high probability in  $m$ , we have

$$\begin{aligned} \|(\mathbf{w} - \mathbf{v})\hat{\mathbf{d}} - (\mathbf{w} - \mathbf{v})\mathbf{d}\|_2^2 &= \sum_{i \in [m]} (w_i - v_i)^2 (\hat{d}_i - d_i)^2 \leq \varepsilon^2 \sum_{i \in [m]} (w_i - v_i)^2 d_i^2 \\ &= \varepsilon^2 \sum_{i \in [m]} (w_i - v_i)^2 \left( \frac{\sigma(\mathbf{w})_i}{w_i} \right)^2 \leq 2\varepsilon^2 \sum_{i \in [m]} \left( \frac{w_i - v_i}{v_i} \right)^2 . \end{aligned}$$

Next we show how to estimate  $f$ . Let  $\mathbf{X} \stackrel{\text{def}}{=} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} \mathbf{A}^T (\mathbf{V} - \mathbf{W}) \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2}$ . By the assumption on  $\alpha$  we know  $-\frac{1}{2}\mathbf{V} \prec \mathbf{V} - \mathbf{W} \prec \frac{1}{2}\mathbf{V}$  and therefore  $-\frac{1}{2}\mathbf{I} \prec \mathbf{X} \prec \frac{1}{2}\mathbf{I}$ . Consequently we have that

$$\begin{aligned} (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} &= (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} (\mathbf{I} - \mathbf{X})^{-1} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} \\ &= \sum_{j=0}^{\infty} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} \mathbf{X}^j (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} . \end{aligned}$$

and therefore

$$\begin{aligned} f_i &= \mathbf{1}_i^T \mathbf{A} \left( \sum_{j=0}^{\infty} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} \mathbf{X}^j (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} - (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} \right) \mathbf{A}^T \mathbf{1}_i \\ &= \sum_{j=1}^{\infty} f_i^{(j)} \quad \text{where} \quad f_i^{(j)} \stackrel{\text{def}}{=} \mathbf{1}_i^T \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} \mathbf{X}^j (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} \mathbf{A}^T \mathbf{1}_i . \end{aligned}$$

Furthermore, using the definition of  $\mathbf{X}$  we have that for even  $j$

$$\begin{aligned} f_i^{(j)} &= \left\| \mathbf{X}^{\frac{j}{2}} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} \mathbf{A}^T \mathbf{1}_i \right\|_2^2 \\ &= \left\| (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} \left( \mathbf{A}^T (\mathbf{V} - \mathbf{W}) \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} \right)^{\frac{j}{2}} \mathbf{A}^T \mathbf{1}_i \right\|_2^2 \\ &= \left\| \sqrt{\mathbf{V}} \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} \left( \mathbf{A}^T (\mathbf{V} - \mathbf{W}) \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} \right)^{\frac{j}{2}} \mathbf{A}^T \mathbf{1}_i \right\|_2^2 \end{aligned}$$

For odd  $j$ , using our definition of  $\mathbf{\Delta}^+$  and  $\mathbf{\Delta}^-$  we have that

$$\begin{aligned} f_i^{(j)} &= \mathbf{1}_i^T \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} \mathbf{X}^j (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} \mathbf{A}^T \mathbf{1}_i \\ &= \mathbf{1}_i^T \mathbf{A} \left( (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} \mathbf{A}^T (\mathbf{V} - \mathbf{W}) \mathbf{A} \right)^{\frac{j-1}{2}} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} \mathbf{A}^T (\mathbf{V} - \mathbf{W}) \\ &\quad \times \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} \left( \mathbf{A}^T (\mathbf{V} - \mathbf{W}) \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} \right)^{\frac{j-1}{2}} \mathbf{A}^T \mathbf{1}_i \\ &= \alpha_i^{(j)} - \beta_i^{(j)} \end{aligned}$$

where

$$\begin{aligned} \alpha_i^{(j)} &\stackrel{\text{def}}{=} \left\| \sqrt{\mathbf{\Delta}^+} \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} \left( \mathbf{A}^T (\mathbf{W} - \mathbf{V}) \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} \right)^{\frac{j-1}{2}} \mathbf{A}^T \mathbf{1}_i \right\|_2^2, \\ \beta_i^{(j)} &\stackrel{\text{def}}{=} \left\| \sqrt{\mathbf{\Delta}^-} \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} \left( \mathbf{A}^T (\mathbf{W} - \mathbf{V}) \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1} \right)^{\frac{j-1}{2}} \mathbf{A}^T \mathbf{1}_i \right\|_2^2. \end{aligned}$$

Consequently, by Lemma 8.4.1 and the construction, we see that

$$\mathbf{E} \hat{f}_i = \sum_{j=1}^{\infty} f_i^{(j)} = f_i$$

and therefore  $\mathbf{E} \hat{h} = \sigma(\mathbf{w}) - \sigma(\mathbf{v})$  as desired. All that remains is to bound the variance of  $\hat{f}_i$ .

To bound the variance of  $\hat{f}$ , let  $|\mathbf{X}| = (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} \mathbf{A}^T |\mathbf{W} - \mathbf{V}| \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2}$ . Note that  $-\frac{1}{4} \mathbf{I} \preceq -|\mathbf{X}| \preceq \mathbf{X} \preceq |\mathbf{X}| \preceq \frac{1}{4} \mathbf{I}$  and consequently for all  $j$

$$\begin{aligned} g_i^{(j)} &\stackrel{\text{def}}{=} \mathbf{1}_i^T \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} |\mathbf{X}|^j (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} \mathbf{A}^T \mathbf{1}_i \\ &\leq \frac{1}{4^{j-1}} \mathbf{1}_i^T \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} |\mathbf{X}| (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} \mathbf{A}^T \mathbf{1}_i \\ &\stackrel{\text{def}}{=} \frac{1}{v_i 4^{j-1}} \mathbf{1}_i^T \mathbf{P}_v \mathbf{\Delta} \mathbf{P}_v \mathbf{1}_i \end{aligned}$$

where  $\mathbf{P}_v = \sqrt{\mathbf{V}} \mathbf{A} (\mathbf{A}^T \mathbf{V} \mathbf{A})^{-1/2} \mathbf{A}^T \sqrt{\mathbf{V}}$  and  $\mathbf{\Delta}$  is a diagonal matrix with  $\Delta_{ii} = \left| \frac{w_i - v_i}{v_i} \right|$ . Using that  $\mathbf{0} \preceq \mathbf{P}_v \preceq \mathbf{I}$ , we have that for all  $j$

$$\begin{aligned}
(4^{j-1})^2 \sum_{i=1}^m \left( v_i g_i^{(j)} \right)^2 &= \sum_{i=1}^m \left( \mathbf{1}_i^T \mathbf{P}_v \Delta \mathbf{P}_v \mathbf{1}_i \right)^2 = \text{Tr}(\mathbf{P}_v \Delta \mathbf{P}_v \mathbf{P}_v \Delta \mathbf{P}_v) \\
&\leq \text{Tr}(\mathbf{P}_v \Delta \Delta \mathbf{P}_v) = \text{Tr}(\Delta \mathbf{P}_v \mathbf{P}_v \Delta) \\
&\leq \text{Tr}(\Delta^2) = \sum_{i=1}^m \left( \frac{w_i - v_i}{v_i} \right)^2 \leq \alpha^2
\end{aligned}$$

and thus  $\|\mathbf{V} \mathbf{g}^{(j)}\|_2 \leq \frac{4\alpha}{4^j}$ . Furthermore, since  $\Delta^+ \preceq |\mathbf{W} - \mathbf{V}|$  and  $\Delta^- \preceq |\mathbf{W} - \mathbf{V}|$  we have that  $|\alpha_i^{(j)}| \leq g_i^{(j)}$  and  $|\beta_i^{(j)}| \leq g_i^{(j)}$ . Consequently, by Lemma 8.4.1 again, we have

$$\begin{aligned}
\|\mathbf{V} \hat{\mathbf{f}}^{(j)} - \mathbf{V} \mathbf{f}^{(j)}\|_2^2 &= \sum_i v_i^2 \left( \hat{f}_i^{(j)} - f_i^{(j)} \right)^2 \\
&\leq 2 \sum_i v_i^2 \left( \hat{\alpha}_i^{(j)} - \alpha_i^{(j)} \right)^2 + 2 \sum_i v_i^2 \left( \hat{\beta}_i^{(j)} - \beta_i^{(j)} \right)^2 \\
&\leq 2\varepsilon^2 \sum_i v_i^2 \left( \left( \alpha_i^{(j)} \right)^2 + \left( \beta_i^{(j)} \right)^2 \right) \\
&\leq 4\varepsilon^2 \sum_i \left( v_i g_i^{(j)} \right)^2 \leq \frac{4\alpha^2 \varepsilon^2}{(4^{j-1})^2}.
\end{aligned}$$

Putting this all together we have that

$$\begin{aligned}
\|\mathbf{V} \hat{\mathbf{f}} - \mathbf{V} \mathbf{f}\|_2 &\leq \|2^u \mathbf{V} \hat{\mathbf{f}}^{(t+u)} + \sum_{j=1}^t \mathbf{V} \hat{\mathbf{f}}^{(j)} - \sum_{j=1}^{\infty} \mathbf{V} \mathbf{f}^{(j)}\|_2 \\
&\leq 2^u \|\mathbf{V} \hat{\mathbf{f}}^{(t+u)}\|_2 + \sum_{j=1}^t \|\mathbf{V} \hat{\mathbf{f}}^{(j)} - \mathbf{V} \mathbf{f}^{(j)}\|_2 + \sum_{j=t+1}^{\infty} \|\mathbf{V} \mathbf{f}^{(j)}\|_2 \\
&\leq 2^u \frac{4\alpha}{4^{t+u-1}} + \sum_{j=1}^t \frac{2\alpha\varepsilon}{4^{j-1}} + \sum_{j=t+1}^{\infty} \frac{2\alpha}{4^{j-1}} \\
&= O\left(\alpha\varepsilon + \frac{\alpha}{4^t}\right).
\end{aligned}$$

Consequently, since  $t = O(\log(\varepsilon^{-1}))$  we have the desired result.  $\square$

## ■ 8.4.2 The Stochastic Chasing 0 Game

To avoid computing leverage scores exactly, in Section 8.4.1 we showed how to estimate the difference of leverage scores and use these to update the leverage scores. However, if we only applied this technique, the error of leverage scores would accumulate in the algorithm and we need to fix it. Naturally, one may wish to use dimension reduction to compute a multiplicative approximation to the leverage scores and update our computed value if the error is too large. However, this strategy would fail if there are too many rows with inaccurate leverage scores in the same iteration. In this case, we would change the central point too much that we are not able to recover. In this section, we present this update problem in a general form that we call *Stochastic Chasing 0 game* and provide an effective strategy for playing this game.

The *Stochastic chasing 0 game* is as follows. There is a player, a stochastic adversary, and a point

$\mathbf{x} \in \mathbb{R}^m$ . The goal of the player is to keep the point close to  $\mathbf{0} \in \mathbb{R}^m$  in  $\ell_\infty$  norm and the goal of the stochastic adversary is to move  $\mathbf{x}$  away from  $\mathbf{0}$ . The game proceeds for an infinite number of iterations where in each iteration the stochastic adversary moves the current point  $\mathbf{x}^{(k)} \in \mathbb{R}^m$  to some new point  $\mathbf{x}^{(k)} + \Delta^{(k)} \in \mathbb{R}^m$  and the player needs to respond. The stochastic adversary cannot move the  $\Delta^{(k)}$  arbitrarily, instead he is only allowed to choose a probability distribution  $\mathcal{D}^{(k)}$  and sample  $\Delta^{(k)}$  from it. Furthermore, it is required that  $\mathbb{E}_{\mathcal{D}^{(k)}} \Delta = \mathbf{0}$  and  $\|\Delta\|_2^2 \leq c$  for some fixed  $c$  and all  $\Delta \in \mathcal{D}^{(k)}$ . The player does not know  $\mathbf{x}^{(k)}$  or the distribution  $\mathcal{D}^{(k)}$  or the move  $\Delta^{(k)}$  of the stochastic adversary. All the player knows is some  $\mathbf{y}^{(k)} \in \mathbb{R}^n$  that is close to  $\mathbf{x}^{(k)}$  in  $\ell_\infty$  norm. With this information, the player is allowed to choose one coordinate  $i$  and set  $x_i^{(k+1)}$  to be zero and for other  $j$ , we have  $x_j^{(k+1)} = x_j^{(k)} + \Delta_j^{(k)}$ .

The question we would like to address is, what strategy the player should choose to keep  $\mathbf{x}^{(k)}$  close to  $\mathbf{0}$  in  $\ell_\infty$  norm? We show that there is a trivial strategy that performs well: simply pick the largest coordinate and set it to 0.

---

**Algorithm 24:** Stochastic chasing  $\mathbf{0}$  game

---

**Constant:**  $c > 0, R > 0$ .

Let  $\mathbf{x}^{(1)} = \mathbf{0} \in \mathbb{R}^m$ .

**for**  $k = 1$  to  $\infty$  **do**

**Stochastic Adversary:** Pick  $\mathcal{D}^{(k)}$  such that  $\mathbb{E}_{\mathcal{D}^{(k)}} \Delta = \mathbf{0}$  and  $\|\Delta\|_2 \leq c$  all  $\Delta \in \mathcal{D}^{(k)}$ .

**Stochastic Adversary:** Pick  $\mathbf{y}^{(k)} \in \mathbb{R}^m$  such that  $\|\mathbf{y}^{(k)} - \mathbf{x}^{(k)}\|_\infty \leq R$ .

**Player:** Pick a coordinate  $i^{(k)}$  using only  $\mathbf{y}^{(k)}$ .

    Sample  $\Delta^{(k)}$  from  $\mathcal{D}^{(k)}$ .

    Set  $x_{i^{(k)}}^{(k+1)} = 0$  and  $x_j^{(k+1)} = x_j^{(k)} + \Delta_j^{(k)}$  for all  $j \neq i^{(k)}$ .

**end**

---

**Theorem 8.4.3.** Using the strategy  $i^{(k)} = \arg \max_i |y_i^{(k)}|$ , with probability at least  $1 - p$ , we have

$$\|\mathbf{x}^{(k)}\|_\infty \leq 2(c + R) \log(4mk^2/p)$$

for all  $k$  in the Stochastic Chasing  $\mathbf{0}$  Game.

*Proof.* Consider the potential function  $\Phi(\mathbf{x}) = \sum_i e^{\alpha x_i} + \sum_i e^{-\alpha x_i}$  where  $\alpha$  is to be determined. Now for all  $x$  we know that  $e^x \leq 1 + x + \frac{x^2}{2} e^{|x|}$  and therefore for all  $|\delta| \leq c$ ,  $x$  and  $\alpha$ , we have

$$e^{\alpha x + \alpha \delta} \leq e^{\alpha x} + \alpha \delta e^{\alpha x} + \frac{1}{2} \alpha^2 \delta^2 e^{\alpha x + |\alpha| c} .$$

Consequently,

$$\begin{aligned} \mathbb{E}_{\Delta \in \mathcal{D}^{(k)}} \Phi(\mathbf{x}^{(k)} + \Delta) &\leq \Phi(\mathbf{x}^{(k)}) + \alpha \mathbb{E}_{\Delta \in \mathcal{D}^{(k)}} \left( \sum_{i \in [m]} e^{\alpha x_i^{(k)}} \Delta_i - \sum_{i \in [m]} e^{-\alpha x_i^{(k)}} \Delta_i \right) \\ &\quad + \frac{\alpha^2}{2} e^{\alpha \|\Delta\|_\infty} \mathbb{E}_{\Delta \in \mathcal{D}^{(k)}} \left( \sum_{i \in [m]} e^{\alpha x_i^{(k)}} \Delta_i^2 + \sum_{i \in [m]} e^{-\alpha x_i^{(k)}} \Delta_i^2 \right) . \end{aligned}$$

Since  $\mathbb{E}_{\mathcal{D}^{(k)}} \Delta = \mathbf{0}$  and  $\|\Delta\|_2 \leq c$ , we have  $\mathbb{E}_{\Delta \in \mathcal{D}^{(k)}} \left( \sum_i e^{\alpha x_i^{(k)}} \Delta_i - \sum_i e^{-\alpha x_i^{(k)}} \Delta_i \right) = 0$  and

$$\begin{aligned} \mathbb{E}_{\Delta \in \mathcal{D}^{(k)}} \left( \sum_i e^{\alpha x_i^{(k)}} \Delta_i^2 + \sum_i e^{-\alpha x_i^{(k)}} \Delta_i^2 \right) &\leq \mathbb{E}_{\Delta \in \mathcal{D}^{(k)}} \left( \sum_i \Delta_i^2 \right) \left( \max_i e^{\alpha x_i^{(k)}} + \max_i e^{-\alpha x_i^{(k)}} \right) \\ &\leq 2c^2 \max_i e^{\alpha |x_i^{(k)}|} . \end{aligned}$$

Letting  $\eta^{(k)} = \max_i e^{\alpha |x_i^{(k)}|}$ , we then have

$$\mathbb{E}_{\Delta \in \mathcal{D}^{(k)}} \Phi(\mathbf{x}^{(k)} + \Delta) \leq \Phi(\mathbf{x}^{(k)}) + \alpha^2 e^{\alpha c} c^2 \eta^{(k)} .$$

Since  $i^{(k)} = \arg \max_i |y_i^{(k)}|$  and  $\|\mathbf{y}^{(k)} - \mathbf{x}^{(k)}\|_\infty \leq R$ , the player setting  $x_{i^{(k)}}^{(k+1)} = 0$  decreases  $\Phi$  by at least  $e^{-\alpha(R+c)} \eta^{(k)}$ . Hence, we have

$$\mathbb{E}_{\Delta \in \mathcal{D}^{(k)}} \Phi(\mathbf{x}^{(k+1)}) \leq \Phi(\mathbf{x}^{(k)}) + \alpha^2 e^{\alpha c} c^2 \eta^{(k)} - e^{-\alpha(R+c)} \eta^{(k)} .$$

Picking  $\alpha = \frac{1}{2(c+R)}$ , we have  $e^{2\alpha(c+R)} (\alpha(c+R))^2 \leq 1$  and hence  $\alpha^2 e^{\alpha c} c^2 \leq e^{-\alpha(R+c)}$ . Therefore, we have that

$$\mathbb{E}_{\Delta \in \mathcal{D}^{(k)}} \Phi(\mathbf{x}^{(k+1)}) \leq \mathbf{E} \Phi(\mathbf{x}^{(k)}) \leq \dots \leq \Phi(\mathbf{x}^{(1)}) = 2m .$$

Consequently, by Markov's inequality we have that  $\Pr[\Phi(\mathbf{x}^{(k)}) \geq \lambda_k] \leq \frac{2m}{\lambda_k}$  for any  $\lambda_k$ . Furthermore, since clearly  $\Phi(\mathbf{x}) \geq e^{\alpha \|\mathbf{x}\|_\infty}$  we have that  $\Pr[\|\mathbf{x}^{(k)}\|_\infty \geq \log(\lambda_k)/\alpha] \leq \frac{2m}{\lambda_k}$  for all  $k$ . Choosing  $\lambda_k = \frac{4mk^2}{p}$  and taking a union bound over all  $k$ , we have that

$$\|\mathbf{x}^{(k)}\|_\infty \leq 2(c+R) \log(4mk^2/p)$$

for all  $k$  with probability at least

$$1 - \sum_{i=1}^{\infty} \frac{2m}{\lambda_k} = 1 - \sum_{k=1}^{\infty} \frac{p}{2k^2} \geq 1 - p .$$

□

## ■ 8.5 Glossary

Here we summarize problem specific notation we use throughout this section. For many quantities we included the typical order of magnitude as they appear during our algorithms.

- Our algorithm maintains a polytope  $\{\mathbf{A}\mathbf{x} \geq \mathbf{b}\}$  which contains the set of solution  $K$ .
- Slacks  $\mathbf{s}(\mathbf{x}) = \mathbf{A}\mathbf{x} - \mathbf{b}$ ,  $\mathbf{S}(\mathbf{x})$  is the diagonal matrix corresponds to  $\mathbf{s}(\mathbf{x})$ . Rescaled constraint matrix  $\mathbf{A}_s = \mathbf{S}^{-1}\mathbf{A}$ .
- Leverage Score:  $\psi(\mathbf{x}) = \mathbf{diag} \left( \mathbf{A}_x (\mathbf{A}_x^T \mathbf{A}_x + \lambda \mathbf{I})^{-1} \mathbf{A}_x^T \right)$ ,  $\mu(\mathbf{x}) = \min \psi_i(\mathbf{x})$ ,  $\lambda = \frac{1}{c_a R^2}$  and  $c_a = \frac{1}{10^{10}}$ .
- Each iteration, if there is a constraint with leverage score less than  $c_d = \frac{1}{10^{12}}$ , then we delete that constraint. Otherwise, we add a constraint and put the leverage score  $c_a$ .
- Each iteration, our algorithm computes the change of leverage score with accuracy  $c_\Delta = \frac{C c_e}{\log(n \log(R/\varepsilon))}$  where  $c_e = \frac{c_d}{4 \ln(6nR/\varepsilon)}$  and  $R$  is the diameter of the box containing  $K$ .



- Potential function  $p_e(\mathbf{x}) = -\sum_{i \in [m]} (c_e + e_i) \log(s_i(\mathbf{x})/R) + \frac{1}{2} \log \det(R^2 (\mathbf{A}^T \mathbf{S}_x^{-2} \mathbf{A} + \lambda \mathbf{I})) + \frac{\lambda}{2} \|\mathbf{x}\|_2^2$ .
- The algorithm starts with  $p_e(\mathbf{x}) \sim -\Theta(1)n$  and ends when  $p_e(\mathbf{x}) \sim \Theta(1)n \log(nR/\varepsilon)$ .
- Centrality  $\delta_e(\mathbf{x}) = \|\nabla p_e(\mathbf{x})\|_{\mathbf{H}(\mathbf{x})^{-1}}$  where  $\mathbf{H}(\mathbf{x}) = \mathbf{A}_x^T (c_e \mathbf{I} + \mathbf{\Psi}(\mathbf{x})) \mathbf{A}_x + \lambda \mathbf{I}$ .  $\delta_e(\mathbf{x}) \approx \Theta(\sqrt{c_e + \mu(\mathbf{x})})$ .
- Approximate Hessian  $\mathbf{Q}(\mathbf{x}, \mathbf{w}) = \mathbf{A}_x^T (c_e \mathbf{I} + \mathbf{W}) \mathbf{A}_x + \lambda \mathbf{I}$ .



# Effective Use of Cutting Plane Method

## ■ 9.1 Introduction

Cutting plane methods have long been employed to obtain polynomial time algorithms for solving optimization problems. However, for many problems cutting plane methods are often regarded as inefficient both in theory and in practice. Here, in this chapter, we provide several techniques for applying cutting plane methods efficiently. Moreover, we illustrate the efficacy and versatility of these techniques by applying them to achieve improved running times for solving multiple problems including semidefinite programming, matroid intersection, and submodular flow.

We hope these results revive interest in ellipsoid and cutting plane methods. We believe these results demonstrate how cutting plan methods are often useful not just for showing that a problem is solvable in polynomial time, but in many yield substantial running time improvements. We stress that while some results in this chapter are problem-specific, the techniques introduced here are quite general and are applicable to a wide range of problems.

In the remainder of this introduction we survey the key techniques we use to apply our cutting plane method (Section 9.1.1) and the key results we obtain on improving the running time for solving various optimization problems (Section 9.1.2). We conclude in Section 9.1.3 by providing an overview of where to find additional technical result in this chapter.

### ■ 9.1.1 Techniques

Although cutting plane methods are typically introduced as algorithms for finding a point in a convex set (as we did with the feasibility problem in Chapter 8), this is often not the easiest way to apply the methods. Moreover, improperly applying results on the feasibility problem to solve convex optimization problems can lead to vastly sub-optimal running times. Our central goal, here, in this chapter is to provide tools that allow cutting plane methods to be efficiently applied to solve complex optimization problems. Some of these tools are new and some are extensions of previously known techniques. Here we briefly survey the techniques we cover in Section 9.3 and Section 9.4.

### Technique 0: From Feasibility to Optimization

In Section 9.3.1, we explain how to use our cutting plane method to solve convex optimization problems using an approximate subgradient oracle. Our result is based on a result of Nemirovski [202] in which he showed how to use a cutting plane method to solve convex optimization problems without smoothness assumptions on the function and with minimal assumptions on the size of the function's domain. We adopt his proof to accommodate for an approximate separation oracle, an extension which is essential for our applications. We use this result as the starting point for two new techniques we discuss below.

### Technique 1: Dimension Reduction through Duality

In Section 9.3.2, we discuss how cutting plane methods can be applied to obtain both primal and dual solutions to convex optimization problems. Moreover, we show how this can be achieved while only applying the cutting plane method in the space, primal or dual, which has a fewer number of variables. Thus we show how to use duality to improve the convergence of cutting plane methods while still solving the original problem.

To illustrate this idea consider the following very simple linear program (LP)

$$\min_{x_i \geq 0, \sum x_i = 1} \sum_{i=1}^n w_i x_i$$

where  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{w} \in \mathbb{R}^n$ . Although this LP has  $n$  variables, it should be easy to solve purely on the grounds that it only has one equality constraint and thus dual linear program is simply

$$\max_{y \leq w_i \forall i} y,$$

i.e. a LP with only one variable. Consequently, we can apply our cutting plane method to solve it efficiently.

However, while this simple example demonstrates how we can use duality to decrease dimensions, it is not always obvious how to recover the optimal primal solution  $x$  variable given the optimal dual solution  $y$ . Indeed, for many problems their dual is significantly simpler than itself (primal), so some work is required to show that working in the space suffices to require a primal solution.

One such recent example of this approach proving successful is the linear program result explained in Chapter 6. In this result, the authors show how to take advantage of this observation and get a faster LP solver and maximum flow algorithm. It is interesting to study how far this technique can extend, that is, in what settings can one recover the solution to a more difficult dual problem from the solution to its easier primal problem?

There is in fact another precedent for such an approach. Grötschel, Lovász and Schrijver[109] showed how to obtain the primal solution for linear program by using a cutting plane method to solve the linear program exactly. This is based on the observation that cutting plane methods are able to find the active constraints of the optimal solution and hence one can take dual of the linear program to get the dual solution. This idea was further extended in [151] which also observed that cutting plane methods are incrementally building up a LP relaxation of the optimization problem. Hence, one can find a dual solution by taking the dual of that relaxation.

In Section 9.3.2, we provide a fairly general technique to recover a dual optimal solution from an approximately optimal primal solution. Unfortunately, the performance of this technique seems quite problem-dependent. We therefore only analyze this technique for semidefinite programming (SDP), a classic and popular convex optimization problem. As a result, we obtain a faster SDP solver in both the primal and dual formulations of the problem.

### Technique 2: Using Optimization Oracles Directly

In the seminal works of Grötschel, Lovász, Schrijver and independently Karp and Papadimitriou [110, 138], they showed the equivalence between optimization oracles and separation oracles, and gave a general method to construct a separation oracle for a convex set given an optimization oracle for that set, that is an oracle for minimizing linear functionals over the set. This seminal result led to the first weakly polynomial time algorithm for many algorithms such as submodular function minimization. Since then, this idea has been used extensively in various settings [127, 44, 45, 65].

Unfortunately, while this equivalence of separation and optimization is a beautiful and powerful tool for polynomial time solvability of problems, in many case it may lead to inefficient algorithms. In order to use this reduction to get a separation oracle, the optimization oracle may need to be called multiple times – essentially the number of times needed to run a cutting plane method and hence may be detrimental to obtaining small asymptotic running times. Therefore, it is an interesting question of whether there is a way of using an optimization oracle more directly.

In Section 9.4 we provide a partial answer to this question for the case of a broad class of problems, that we call the *intersection problem*. For these problems we demonstrate how to achieve running time improvements by using optimization oracles directly. The problem we consider is as follows. We wish to solve the problem for some cost vector  $\mathbf{c} \in \mathbb{R}^n$  and convex set  $K$ . We assume that the convex set  $K$  can be decomposed as  $K = K_1 \cap K_2$  such that  $\max_{\mathbf{x} \in K_1} \langle \mathbf{c}, \mathbf{x} \rangle$  and  $\max_{\mathbf{x} \in K_2} \langle \mathbf{c}, \mathbf{x} \rangle$  can each be solved efficiently. Our goal is to obtain a running time for this problem comparable to that of minimizing  $K$  given only a separation oracle for it.

We show that by considering a carefully regularized variant, we obtain a problem such that optimization oracles for  $K_1$  and  $K_2$  immediately yield a separation oracle for this regularized problem. By analyzing the regularizer and bounding the domains of the problem we are able to show that this allows us to efficiently compute highly accurate solutions to the intersection problem by applying our cutting plane method once. In other words, we do not need to use a complicated iterative scheme or directly invoke the equivalence between separation and optimization and thereby save  $O(\text{poly}(n))$  factors in our running times.

We note that this intersection problem can be viewed as a generalization of the matroid intersection problem and in Section 9.4.2, we show our reduction gives a faster algorithm in certain parameter regimes. As another example, in Section 9.4.3 we show our reduction gives a substantial polynomial improvement for the submodular flow problem. Furthermore, in Section 9.4.4 we show how our techniques allow us to minimize a linear function over the intersection of a convex set and an affine subspace in a number of iterations that depends only on the co-dimension of the affine space.

## ■ 9.1.2 Applications

Our main goal in this chapter is to provide general techniques for efficiently using cutting plane methods for various problems. Hence, in this chapter, we use minimally problem-specific techniques to achieve the best possible running time. However, we also demonstrate the efficacy of our approach by showing how techniques improve upon the previous best known running times for solve several classic problems in combinatorial and continuous optimization. Here we provide a brief overview of these applications, previous work on these problems, and our results.

In order to avoid deviating from our main discussion, our coverage of previous methods and techniques is brief. Given the large body of prior works on SDP, matroid intersection and submodular flow, it would be impossible to have an in-depth discussion on all of them. Therefore, this section focuses on running time comparisons and explanations of relevant previous techniques.

## Semidefinite Programming

In Section 9.3.2 we consider the classic semidefinite programming (SDP) problem:

$$\max_{\mathbf{X} \succeq \mathbf{0}} \mathbf{C} \bullet \mathbf{X} \text{ s.t. } \mathbf{A}_i \bullet \mathbf{X} = b_i \text{ (primal)} \quad \min_{\mathbf{y}} \mathbf{b}^T \mathbf{y} \text{ s.t. } \sum_{i=1}^n y_i \mathbf{A}_i \succeq \mathbf{C} \text{ (dual)}$$

where  $\mathbf{X}$ ,  $\mathbf{C}$ ,  $\mathbf{A}_i$  are  $m \times m$  symmetric matrices,  $\mathbf{b}, \mathbf{y} \in \mathbb{R}^n$ , and  $\mathbf{A} \bullet \mathbf{B} \stackrel{\text{def}}{=} \text{Tr}(\mathbf{A}^T \mathbf{B})$ . For many problems,  $n \ll m^2$  and hence the dual problem has fewer variables than the primal. There are many

Years	Authors	Running times
1992	Nesterov, Nemirovsky [210]	$\tilde{O}(\sqrt{n}(nm^\omega + n^{\omega-1}m^2))$
2000	Anstreicher [13]	$\tilde{O}((mn)^{1/4}(nm^\omega + n^{\omega-1}m^2))$
2003	Krishnan, Mitchell [152]	$\tilde{O}(m(n^\omega + m^\omega + S))$ (dual SDP)
-	This Chapter	$\tilde{O}(m(n^2 + m^\omega + S))$

**Table 9.1:** Previous algorithms for solving a  $n \times n$  SDP with  $m$  constraints and  $S$  non-zeros entries

results and applications of SDP; see [259, 248, 196] for a survey on this topic. Since our focus is on polynomial time algorithms, we do not discuss pseudo-polynomial algorithms such as the spectral bundle method [115], multiplicative weight update methods [17, 19, 126, 8], etc.

Currently, there are two competing approaches for solving SDP problems, namely interior point methods (IPM) and cutting plane methods. Typically, IPMs require fewer iterations than the cutting plane methods, however each iteration of these methods is more complicated and possibly more computationally expensive. For SDP problems, interior point methods require the computations of the Hessian of the function  $-\log \det(\mathbf{C} - \sum_{i=1}^n y_i \mathbf{A}_i)$  whereas cutting plane methods usually only need to compute minimum eigenvectors of the slack matrix  $\mathbf{C} - \sum_{i=1}^n y_i \mathbf{A}_i$ .

In [13], Anstreicher provided the current fastest IPM for solving the dual SDP problem using a method based on the volumetric barrier function. This method takes  $O((mn)^{1/4})$  iterations and each iteration is as cheap as usual IPMs. For general matrices  $\mathbf{C}, \mathbf{X}, \mathbf{A}_i$ , each iteration takes  $O(nm^\omega + n^{\omega-1}m^2)$  time where  $\omega$  is the fast matrix multiplication exponent. If the constraint matrices  $\mathbf{A}_i$  are rank one matrices, the iteration cost can be improved to  $O(m^\omega + nm^2 + n^2m)$  [153]. If the matrices are sparse, then [97, 197] show how to use matrix completion inside the IPM. However, the running time depends on the extended sparsity patterns which can be much larger than the total number of non-zeros.

In [152], Krishnan and Mitchell observed that the separation oracle for dual SDP takes only  $O(m^\omega + S)$  time, where  $S = \sum_{i=1}^n \text{nnz}(\mathbf{A}_i)$  be the total number of non-zeros in the constant matrix. Hence, the cutting plane method by [253] gives a faster algorithm for SDP for many regimes. For  $\omega = 2.38$ , the cutting plane method is faster when  $\mathbf{A}_i$  is not rank 1 and the problem is not too dense, i.e.  $\sum_{i=1}^n \text{nnz}(\mathbf{A}_i) < n^{0.63}m^{2.25}$ . While there are previous methods for using cutting plane methods to obtain primal solutions [151], to the best of our knowledge, there are no worst case running time analysis for these techniques.

In Section 9.3.2, show how to alleviate this issue. We provide an improved algorithm for finding the dual solution and prove carefully how to obtain a comparable primal solution as well. See Figure 9.1 for a summary of the algorithms for SDP and their running times.

## Matroid Intersection

In Section 9.4.2 we show how our optimization oracle technique can be used to improve upon the previous best known running times for matroid intersection. Matroid intersection is one of the most fundamental problems in combinatorial optimization. The first algorithm for matroid intersection is due to the seminal paper by Edmonds [81]. In Figures 9.2 and 9.3 we provide a summary of the previous algorithms for unweighted and weighted matroid intersection as well as the new running times we obtain in this chapter. While there is no total ordering on the running times of these algorithms due to the different dependence on various parameters, we would like to point out that our algorithms outperform the previous ones in regimes where  $r$  is close to  $n$  and/or the oracle query costs are relatively expensive. In particular, in terms of oracle query complexity our algorithms are the first to

Years	Authors	Running times
1968	Edmonds [81]	not stated
1971	Aigner, Dowling [6]	$O(nr^2\mathcal{T}_{\text{ind}})$
1974	Tomizawa, Iri [249]	not stated
1975	Lawler [157]	$O(nr^2\mathcal{T}_{\text{ind}})$
1979	Edmonds [80]	not stated
1986	Cunningham [60]	$O(nr^{1.5}\mathcal{T}_{\text{ind}})$
-	This Chapter	$O(n^2 \log n \mathcal{T}_{\text{ind}} + n^3 \log^{O(1)} n)$ $O(nr \log^2 n \mathcal{T}_{\text{rank}} + n^3 \log^{O(1)} n)$

**Table 9.2:** Previous algorithms for (unweighted) matroid intersection. Here  $n$  is the size of the ground set,  $r = \max\{r_1, r_2\}$  is the maximum rank of the two matroids,  $\mathcal{T}_{\text{ind}}$  is the time needed to check if a set is independent (independence oracle), and  $\mathcal{T}_{\text{rank}}$  is the time needed to compute the rank of a given set (rank oracle).

Years	Authors	Running times
1968	Edmonds [81]	not stated
1974	Tomizawa, Iri [249]	not stated
1975	Lawler [157]	$O(nr^2\mathcal{T}_{\text{ind}} + nr^3)$
1979	Edmonds [80]	not stated
1981	Frank [90]	$O(n^2r(\mathcal{T}_{\text{circuit}} + n))$
1983	Orlin, Ahuja [219]	not stated
1986	Brezovec, Cornuéjols, Glover [39]	$O(nr(\mathcal{T}_{\text{circuit}} + r + \log n))$
1995	Fujishige, Zhang [96]	$O(n^2r^{0.5} \log rM \cdot \mathcal{T}_{\text{ind}})$
1995	Shigeno, Iwata [234]	$O((n + \mathcal{T}_{\text{circuit}})nr^{0.5} \log rM)$
-	This Chapter	$O((n^2 \log n \mathcal{T}_{\text{ind}} + n^3 \log^{O(1)} n) \log nM)$ $O((nr \log^2 n \mathcal{T}_{\text{rank}} + n^3 \log^{O(1)} n) \log nM)$

**Table 9.3:** Previous algorithms for weighted matroid intersection. In additions to the notations used in the unweighted table,  $\mathcal{T}_{\text{circuit}}$  is the time needed to find a fundamental circuit and  $M$  is the bit complexity of the weights.

achieve the quadratic bounds of  $\tilde{O}(n^2)$  and  $\tilde{O}(nr)$  for independence and rank oracles. We hope our work will revive the interest in the problem of which progress has been mostly stagnated for the past 20-30 years.

### Minimum-Cost Submodular Flow

In Section 9.4.3 we show how our optimization oracle technique can be used to improve upon the previous best known running times for (Minimum-cost) Submodular Flow. Submodular flow is a very general problem in combinatorial optimization which generalizes many problems such as minimum cost flow, the graph orientation, polymatroid intersection, directed cut covering [94]. In Figure 9.4 we provide an overview of the previous algorithms for submodular flow as well as the new running times we obtain in this chapter.

Many of the running times are in terms of a parameter  $h$ , which is the time required for computing an “exchange capacity”. To the best of our knowledge, the most efficient way of computing an exchange capacity is to solve an instance of submodular minimization which previously took time  $\tilde{O}(n^4\text{EO} + n^5)$  (and now takes  $\tilde{O}(n^2\text{EO} + n^3)$  time using our result in Chapter 10). Readers may wish to substitute  $h = \tilde{O}(n^2\text{EO} + n^3)$  when reading the table.

Years	Authors	Running times
1978	Fujishige [92]	not stated
1981	Grötschel, Lovász, Schrijver[110]	weakly polynomial
1982	Zimmermann [276]	not stated
1984	Barahona, Cunningham [30]	not stated
1985	Cunningham, Frank [62]	$\rightarrow O(n^4 h \log C)$
1987	Fujishige [93]	not stated
1987	Frank, Tardos [91]	strongly polynomial
1988	Cui, Fujishige [265]	not stated
1989	Fujishige, Röck, Zimmermann[95]	$\rightarrow O(n^6 h \log n)$
1991	Chung, Tcha [52]	not stated
1992	Zimmermann [277]	not stated
1993	McCormick, Ervolina [190]	$O(n^7 h^* \log nCU)$
1994	Wallacher, Zimmermann [266]	$O(n^8 h \log nCU)$
1997	Iwata [117]	$O(n^7 h \log U)$
1998	Iwata, McCormick, Shigeno [123]	$O(n^4 h \min \{\log nC, n^2 \log n\})$
1999	Iwata, McCormick, Shigeno [124]	$O(n^6 h \min \{\log nU, n^2 \log n\})$
1999	Fleischer, Iwata, McCormick[89]	$O(n^4 h \min \{\log U, n^2 \log n\})$
1999	Iwata, McCormick, Shigeno [122]	$O(n^4 h \min \{\log C, n^2 \log n\})$
2000	Fleischer, Iwata [88]	$O(mn^5 \log nU \cdot \text{EO})$
-	This Chapter	$O(n^2 \log nCU \cdot \text{EO} + n^3 \log^{O(1)} nCU)$

**Figure 9-1:** Previous algorithms for Submodular Flow with  $n$  vertices, maximum cost  $C$  and maximum capacity  $U$ . The factor  $h$  is the time for an exchange capacity oracle,  $h^*$  is the time for a “more complicated exchange capacity oracle” and EO is the time for evaluation oracle of the submodular function. The arrow,  $\rightarrow$ , indicates that it used currently best maximum submodular flow algorithm as subroutine which was non-existent at the time of the publication.

The previous fastest weakly polynomial algorithms for submodular flow are by [122, 88, 89], which take time  $\tilde{O}(n^6 \text{EO} + n^7)$  and  $O(mn^5 \log nU \cdot \text{EO})$ , assuming  $h = \tilde{O}(n^2 \text{EO} + n^3)$ . Our algorithm for submodular flow has a running time of  $\tilde{O}(n^2 \text{EO} + n^3)$ , which is significantly faster by roughly a factor of  $\tilde{O}(n^4)$ .

For strongly polynomial algorithms, our results do not yield a speedup but we remark that our faster strongly polynomial algorithm for submodular minimization in Chapter 10 improves the previous algorithms by a factor of  $\tilde{O}(n^2)$  as a corollary (because  $h$  requires solving an instance of submodular minimization).

### ■ 9.1.3 Overview

After providing covering some preliminaries on convex analysis in Section 9.2 we split the remainder of Chapter 9 into Section 9.3 and Section 9.4. In Section 9.3 we cover our algorithm for convex optimization using an approximate subgradient oracle (Section 9.3.1) as well as our technique on using duality to decrease dimensions and improve the running time of semidefinite programming (Section 9.3.2). In Section 9.4 we provide our technique for using minimization oracles to minimize functions over the intersection of convex sets and provide several applications including matroid intersection (Section 9.4.2), submodular flow (Section 9.4.3), and minimizing a linear function over the intersection of an affine subspace and a convex set (Section 9.4.4).



## ■ 9.2 Preliminaries

In this section we review basic facts about convex functions that we use throughout Chapter 9. We also introduce two oracles that we use throughout Chapter 9, i.e. subgradient and optimization oracles, and provide some basic reductions between them. Note that we have slightly extended some definitions and facts to accommodate for the noisy separation oracles used in this chapter.

First we recall the definition of strong convexity

**Definition 9.2.1** (Strong Convexity). A real valued function  $f$  on a convex set  $\Omega$  is  $\alpha$ -strongly convex if for any  $\mathbf{x}, \mathbf{y} \in \Omega$  and  $t \in [0, 1]$ , we have

$$f(t\mathbf{x} + (1-t)\mathbf{y}) + \frac{1}{2}\alpha t(1-t)\|\mathbf{x} - \mathbf{y}\|^2 \leq tf(\mathbf{x}) + (1-t)f(\mathbf{y}).$$

Next we define an approximate subgradient.

**Definition 9.2.2** (Subgradient). For any convex function  $f$  on a convex set  $\Omega$ , the  $\delta$ -subgradients of  $f$  at  $\mathbf{x}$  are defined to be

$$\partial_\delta f(\mathbf{x}) \stackrel{\text{def}}{=} \{\mathbf{g} \in \Omega : f(\mathbf{y}) + \delta \geq f(\mathbf{x}) + \langle \mathbf{g}, \mathbf{y} - \mathbf{x} \rangle \text{ for all } \mathbf{y} \in \Omega\}.$$

Here we provide some basic facts regarding convexity and subgradients. These statements are natural extensions of well known facts regarding convex functions and their proof can be found in any standard textbook on convex optimization.

**Fact 9.2.3.** For any convex set  $\Omega$  and  $\mathbf{x}$  be a point in the interior of  $\Omega$ , we have the following:

1. If  $f$  is convex on  $\Omega$ , then  $\partial_0 f(\mathbf{x}) \neq \emptyset$  and  $\partial_s f(\mathbf{x}) \subseteq \partial_t f(\mathbf{x})$  for all  $0 \leq s \leq t$ . Otherwise, we have  $\|\mathbf{g}\|_2 > \frac{1}{2}\sqrt{\frac{\delta}{D}}$ . For any  $f(\mathbf{y}) \leq f(\mathbf{x})$ , we have  $\delta \geq \langle \mathbf{g}, \mathbf{y} - \mathbf{x} \rangle$  and hence
2. If  $f$  is a differential convex function on  $\Omega$ , then  $\nabla f(\mathbf{x}) \in \partial_0 f(\mathbf{x})$ .
3. If  $f_1$  and  $f_2$  are convex function on  $\Omega$ ,  $\mathbf{g}_1 \in \partial_{\delta_1} f_1(\mathbf{x})$  and  $\mathbf{g}_2 \in \partial_{\delta_2} f_2(\mathbf{x})$ , then  $\alpha\mathbf{g}_1 + \beta\mathbf{g}_2 \in \partial_{\alpha\delta_1 + \beta\delta_2} (\mathbf{g}_1 + \mathbf{g}_2)(\mathbf{x})$ .
4. If  $f$  is  $\alpha$ -strongly convex on  $\Omega$  with minimizer  $\mathbf{x}^*$ , then for any  $\mathbf{y}$  with  $f(\mathbf{y}) \leq f(\mathbf{x}^*) + \varepsilon$ , we have  $\frac{1}{2}\alpha\|\mathbf{x}^* - \mathbf{y}\|^2 \leq \varepsilon$ .

Next we provide a reduction from subgradients to separation oracles. We will use this reduction several times in Chapter 9 to simplify our construction of separation oracles.

**Lemma 9.2.4.** Let  $f$  be a convex function. Suppose we have  $\mathbf{x}$  and  $\mathbf{g} \in \partial_\delta f(\mathbf{x})$  with  $\|\mathbf{x}\|_2 \leq 1 \leq D$  and  $\delta \leq 1$ . If  $\|\mathbf{g}\|_2 \leq \frac{1}{2}\sqrt{\frac{\delta}{D}}$ , then  $f(\mathbf{x}) \leq \min_{\|\mathbf{y}\|_2 \leq D} f(\mathbf{y}) + 2\sqrt{\delta D}$  and if  $\|\mathbf{g}\|_2 \leq \frac{1}{2}\sqrt{\frac{\delta}{D}}$  then

$$\{\|\mathbf{y}\|_2 \leq D : f(\mathbf{y}) \leq f(\mathbf{x})\} \subset \{\mathbf{y} : \mathbf{d}^T \mathbf{y} \leq \mathbf{d}^T \mathbf{x} + 2\sqrt{\delta D}\}$$

with  $\mathbf{d} = \mathbf{g}/\|\mathbf{g}\|_2$ . Hence, this gives a  $(2\sqrt{\delta D}, 2\sqrt{\delta D})$ -separation oracle on the set  $\{\|\mathbf{x}\|_2 \leq D\}$ .

*Proof.* Let  $\mathbf{y}$  such that  $\|\mathbf{y}\|_2 \leq D$ . By the definition of  $\delta$ -subgradient, we have

$$f(\mathbf{y}) + \delta \geq f(\mathbf{x}) + \langle \mathbf{g}, \mathbf{y} - \mathbf{x} \rangle.$$

If  $\|\mathbf{g}\| \leq \frac{1}{2}\sqrt{\frac{\delta}{D}}$ , then, we have  $|\langle \mathbf{g}, \mathbf{y} - \mathbf{x} \rangle| \leq \sqrt{\delta D}$  because  $\|\mathbf{x}\| \leq D$  and  $\|\mathbf{y}\|_2 \leq D$ . Therefore,

$$\min_{\|\mathbf{y}\|_2 \leq D} f(\mathbf{y}) + 2\sqrt{\delta D} \geq f(\mathbf{x}).$$

Otherwise, we have  $\|\mathbf{g}\|_2 > \frac{1}{2}\sqrt{\frac{\delta}{D}}$ . For any  $f(\mathbf{y}) \leq f(\mathbf{x})$ , we have  $\delta \geq \langle \mathbf{g}, \mathbf{y} - \mathbf{x} \rangle$  and hence

$$2\sqrt{\delta D} \geq \left\langle \frac{\mathbf{g}}{\|\mathbf{g}\|}, \mathbf{y} - \mathbf{x} \right\rangle.$$

□

At several times in Chapter 9 we will wish to construct subgradient oracles or separation oracles given only the ability to approximately maximize a linear function over a convex set. In the remainder of this section we formally define such a *optimization oracle* and prove this equivalence.

**Definition 9.2.5** (Optimization Oracle). Given a convex set  $K$  and  $\delta > 0$  a  $\delta$ -optimization oracle for  $K$  is a function on  $\mathbb{R}^n$  such that for any input  $\mathbf{c} \in \mathbb{R}^n$ , it outputs  $\mathbf{y}$  such that

$$\max_{\mathbf{x} \in K} \langle \mathbf{c}, \mathbf{x} \rangle \leq \langle \mathbf{c}, \mathbf{y} \rangle + \delta.$$

We denote by  $\text{OO}_\delta(K)$  the time complexity of this oracle.

**Lemma 9.2.6.** *Given a convex set  $K$ , any  $\varepsilon$ -optimization oracle for  $K$  is a  $\varepsilon$ -subgradient oracle for  $f(\mathbf{c}) \stackrel{\text{def}}{=} \max_{\mathbf{x} \in K} \langle \mathbf{c}, \mathbf{x} \rangle$ .*

*Proof.* Let  $\mathbf{x}_c$  be the output of  $\varepsilon$ -optimization oracle on the cost vector  $\mathbf{c}$ . We have

$$\max_{\mathbf{x} \in K} \langle \mathbf{c}, \mathbf{x} \rangle \leq \langle \mathbf{c}, \mathbf{x}_c \rangle + \varepsilon.$$

Hence, for all  $\mathbf{d}$ , we have and therefore

$$\langle \mathbf{x}_c, \mathbf{d} - \mathbf{c} \rangle + f(\mathbf{c}) \leq f(\mathbf{d}) + \varepsilon.$$

Hence,  $\mathbf{x}_c \in \partial_\delta f(\mathbf{c})$ . □

Combining these lemmas shows that having an  $\varepsilon$ -optimization oracle for a convex set  $K$  contained in a ball of radius  $D$  yields a  $O(\sqrt{D\varepsilon}, \sqrt{D\varepsilon})$  separation oracle for  $\max_{\mathbf{x} \in K} \langle \mathbf{c}, \mathbf{x} \rangle$ . We use these ideas to construction separation oracles throughout Chapter 9.

## ■ 9.3 Convex Optimization

In this section we show how to apply our cutting plane method to efficiently solve problems in convex optimization. First, in Section 9.3.1 we show how to use our result to minimize a convex function given an approximate subgradient oracle. Then, in Section 9.3.2 we illustrate how this result can be used to obtain both primal and dual solutions for a standard convex optimization problems. In particular, we show how our result can be used to obtain improved running times for semidefinite programming across a range of parameters.

### ■ 9.3.1 From Feasibility to Optimization

In this section we consider the following standard optimization problem. We are given a convex function  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  and we want to find a point  $\mathbf{x}$  that approximately solves the minimization

problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

given only a subgradient oracle for  $f$ .

Here we show how to apply the cutting plane method from Chapter 8 turning the small width guarantee of the output of that algorithm into a tool to find an approximate minimizer of  $f$ . Our result is applicable to any convex optimization problem armed with a separation or subgradient oracle. This result will serve as the foundation for many of our applications in this chapter.

Our approach is an adaptation of Nemiroski's method [202] which applies the cutting plane method to solve convex optimization problems, with only minimal assumption on the cutting plane method. The proof here is a generalization that accommodates for the noisy separation oracle. In the remainder of this subsection we provide a key definition we will use in our algorithm (Definition 9.3.1), provide our main result (Theorem 9.3.2), and conclude with a brief discussion of this result.

**Definition 9.3.1.** For any compact set  $K$ , we define the *minimum width* by  $\text{MinWidth}(K) \stackrel{\text{def}}{=} \min_{\|\mathbf{a}\|_2=1} \max_{\mathbf{x}, \mathbf{y} \in K} \langle \mathbf{a}, \mathbf{x} - \mathbf{y} \rangle$ .

**Theorem 9.3.2.** Let  $f$  be a convex function on  $\mathbb{R}^n$  and  $\Omega$  be a convex set that contains a minimizer of  $f$ . Suppose we have a  $(\eta, \delta)$ -separation oracle for  $f$  and  $\Omega$  is contained inside  $B_\infty(R)$ . Using  $B_\infty(R)$  as the initial polytope for our Cutting Plane Method, for any  $0 < \alpha < 1$ , we can compute  $\mathbf{x} \in \mathbb{R}^n$  such that

$$f(\mathbf{x}) - \min_{\mathbf{y} \in \Omega} f(\mathbf{y}) \leq \eta + \alpha \left( \max_{\mathbf{y} \in \Omega} f(\mathbf{y}) - \min_{\mathbf{y} \in \Omega} f(\mathbf{y}) \right) . \quad (9.1)$$

with an expected running time of

$$O \left( n \text{SO}_{\eta, \delta}(f) \log \left( \frac{n\kappa}{\alpha} \right) + n^3 \log^{O(1)} \left( \frac{n\kappa}{\alpha} \right) \right),$$

where  $\delta = \Theta \left( \frac{\alpha \text{MinWidth}(\Omega)}{n^{3/2} \ln(\kappa)} \right)$  and  $\kappa = \frac{R}{\text{MinWidth}(\Omega)}$ . Furthermore, we only need the oracle defined on the set  $B_\infty(R)$ .

*Proof.* Let  $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \Omega} f(\mathbf{x})$ . Since  $B_\infty(R) \supset \Omega$  contains a minimizer of  $f$ , by the definition of  $(\eta, \delta)$ -separation oracles, our Cutting Plane Method (Theorem 8.3.26) either returns a point  $\mathbf{x}$  that is almost optimal or returns a polytope  $P$  of small width. In the former case we have a point  $\mathbf{x}$  such that  $f(\mathbf{x}) \leq \min_{\mathbf{y}} f(\mathbf{y}) + \eta$ . Hence, the error is clearly at most  $\eta + \alpha (\max_{\mathbf{z} \in \Omega} f(\mathbf{z}) - \min_{\mathbf{x} \in \Omega} f(\mathbf{x}))$  as desired. Consequently, we assume the latter case.

Theorem 8.3.26 shows  $\text{MinWidth}(P) < Cn\epsilon \ln(R/\epsilon)$  for some universal constant  $C$ . Picking

$$\epsilon = C' \frac{\alpha \text{MinWidth}(\Omega)}{n \ln \left( \frac{n\kappa}{\alpha} \right)} \quad (9.2)$$

for small enough constant  $C'$ , we have  $\text{MinWidth}(P^{(i)}) < \alpha \text{MinWidth}(\Omega)$ . Let  $\Omega^\alpha = \mathbf{x}^* + \alpha(\Omega - \mathbf{x}^*)$ , namely,  $\Omega^\alpha = \{\mathbf{x}^* + \alpha(\mathbf{z} - \mathbf{x}^*) : \mathbf{z} \in \Omega\}$ . Then, we have

$$\text{MinWidth}(\Omega^\alpha) = \alpha \text{MinWidth}(\Omega) > \text{MinWidth}(P).$$

Therefore,  $\Omega^\alpha$  is not a subset of  $P^{(i)}$  and hence there is some point  $\mathbf{y} \in \Omega^\alpha \setminus P$ . Since  $\Omega^\alpha \subseteq \Omega \subseteq B_\infty(R)$ , we know that  $\mathbf{y}$  does not violate any of the constraints of  $P^{(0)}$  and therefore must violate one of the constraints added by querying the separation oracle. Therefore, for some  $j \leq i$ , we have

$$\langle \mathbf{c}^{(j-1)}, \mathbf{y} \rangle > \langle \mathbf{c}^{(j-1)}, \mathbf{x}^{(j-1)} \rangle + c_s \epsilon / \sqrt{n} .$$

By the definition of  $(\eta, c_s \varepsilon / \sqrt{n})$ -separation oracle (Definition 2.3.6), we have  $f(\mathbf{y}) > f(\mathbf{x}^{(j-1)})$ . Since  $\mathbf{y} \in \Omega^\alpha$ , we have  $\mathbf{y} = (1 - \alpha)\mathbf{x}^* + \alpha\mathbf{z}$  for some  $\mathbf{z} \in \Omega$ . Thus, the convexity of  $f$  implies that

$$f(\mathbf{y}) \leq (1 - \alpha)f(\mathbf{x}^*) + \alpha f(\mathbf{z}).$$

Therefore, we have

$$\min_{1 \leq k \leq i} f(\mathbf{x}^{(k)}) - \min_{\mathbf{x} \in \Omega} f(\mathbf{x}) < f(\mathbf{y}) - f(\mathbf{x}^*) \leq \alpha \left( \max_{\mathbf{z} \in \Omega} f(\mathbf{z}) - \min_{\mathbf{x} \in \Omega} f(\mathbf{x}) \right).$$

Hence, we can simply output the best  $\mathbf{x}$  among all  $\mathbf{x}^{(j)}$  and in either case  $\mathbf{x}$  satisfies (9.1).

Note that we need to call  $(\eta, \delta)$ -separation oracle with  $\delta = \Omega(\varepsilon / \sqrt{n})$  to ensure we do not cut out  $\mathbf{x}^*$ . Theorem 8.3.26 shows that the algorithm takes  $O(n\text{SO}_{\eta, \delta}(f) \log(nR/\varepsilon) + n^3 \log^{O(1)}(nR/\varepsilon))$  expected time, as promised. Furthermore, the oracle needs only be defined on  $B_\infty(R)$  as our cutting plane method guarantees  $\mathbf{x}^{(k)} \in B_\infty(R)$  for all  $k$  (although if needed, an obvious separating hyperplane can be returned for a query point outside  $B_\infty(R)$ ).  $\square$

Observe that this algorithm requires no information about  $\Omega$  (other than that  $\Omega \subseteq B_\infty(R)$ ) and does not guarantee that the output is in  $\Omega$ . Hence, even though  $\Omega$  can be complicated to describe, the algorithm still gives a guarantee related to the gap  $\max_{\mathbf{x} \in \Omega} f(\mathbf{x}) - \min_{\mathbf{x} \in \Omega} f(\mathbf{x})$ . For specific applications, it is therefore advantageous to pick a  $\Omega$  as large as possible while the bound on function value is as small as possible.

Before indulging into specific applications, we remark on the dependence on  $\kappa$ . Using John's ellipsoid, it can be shown that any convex set  $\Omega$  can be transformed linearly such that (1)  $B_\infty(1)$  contains  $\Omega$  and, (2)  $\text{MinWidth}(\Omega) = \Omega(n^{-3/2})$ . In other words,  $\kappa$  can be effectively chosen as  $O(n^{3/2})$ . Therefore if we are able to find such a linear transformation, the running time is simply  $O\left(n\text{SO}(f) \log(n/\alpha) + n^3 \log^{O(1)}(n/\alpha)\right)$ . Often this can be done easily using the structure of the particular problem and the running time does not depend on the size of domain at all.

### ■ 9.3.2 Duality and Semidefinite Programming

In this section we illustrate how our result in Section 9.3.1 can be used to obtain both primal and dual solutions for standard problems in convex optimization. In particular we show how to obtain improved running times for semidefinite programming.

To explain our approach, consider the following minimax problem

$$\min_{\mathbf{y} \in Y} \max_{\mathbf{x} \in X} \langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{c}, \mathbf{x} \rangle + \langle \mathbf{d}, \mathbf{y} \rangle \quad (9.3)$$

where  $\mathbf{x} \in \mathbb{R}^m$  and  $\mathbf{y} \in \mathbb{R}^n$ . When  $m \gg n$ , solving this problem by directly using Chapter 8 could lead to an inefficient algorithm with running time at least  $m^3$ . In many situations, for any fixed  $\mathbf{y}$ , the problem  $\max_{\mathbf{x} \in X} \langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle$  is very easy and hence one can use it as a separation oracle and apply Chapter 8 and this would give a running time almost independent of  $m$ . However, this would only give us the  $\mathbf{y}$  variable and it is not clear how to recover  $\mathbf{x}$  variable from it.

In this section we show how to alleviate this issue and give semidefinite programming (SDP) as a concrete example of how to apply this general technique. We do not write down the general version as the running time of the technique seems to be problem specific and faster SDP is already an interesting application.

For the remainder of this section we focus on the semidefinite programming (SDP) problem:

$$\max_{\mathbf{X} \succeq \mathbf{0}} \mathbf{C} \bullet \mathbf{X} \text{ s.t. } \mathbf{A}_i \bullet \mathbf{X} = b_i \quad (9.4)$$

and its dual

$$\min_{\mathbf{y}} \mathbf{b}^T \mathbf{y} \text{ s.t. } \sum_{i=1}^n y_i \mathbf{A}_i \succeq \mathbf{C} \quad (9.5)$$

where  $\mathbf{X}$ ,  $\mathbf{C}$ ,  $\mathbf{A}_i$  are  $m \times m$  symmetric matrices and  $\mathbf{b}, \mathbf{y} \in \mathbb{R}^n$ . Our approach is partially inspired by one of the key ideas of [115, 152]. These results write down the dual SDP in the form

$$\min_{\mathbf{y}} \mathbf{b}^T \mathbf{y} - K \min(\lambda_{\min}(\sum_{i=1}^n y_i \mathbf{A}_i - \mathbf{C}), 0) \quad (9.6)$$

for some large number  $K$  and use non-smooth optimization techniques to solve the dual SDP problem. Here, we follow the same approach but instead write it as a max-min problem  $\min_{\mathbf{y}} f_K(\mathbf{y})$  where

$$f_K(\mathbf{y}) = \max_{\text{Tr} \mathbf{X} \leq K, \mathbf{X} \succeq \mathbf{0}} \left( \mathbf{b}^T \mathbf{y} + \left\langle \mathbf{X}, \mathbf{C} - \sum_{i=1}^n y_i \mathbf{A}_i \right\rangle \right). \quad (9.7)$$

Thus the SDP problem in fact assumes the form (9.3) and many ideas in this section can be generalized to the minimax problem (9.3).

To get a dual solution, we notice that the cutting plane method maintains a subset of the primal feasible solution  $\text{conv}(\mathbf{X}_i)$  such that

$$\min_{\mathbf{y}} \mathbf{b}^T \mathbf{y} + \max_{\text{Tr} \mathbf{X} \leq K, \mathbf{X} \succeq \mathbf{0}} \left\langle \mathbf{X}, \mathbf{C} - \sum_{i=1}^n y_i \mathbf{A}_i \right\rangle \sim \min_{\mathbf{y}} \mathbf{b}^T \mathbf{y} + \max_{\mathbf{X} \in \text{conv}(\mathbf{X}_i)} \left\langle \mathbf{X}, \mathbf{C} - \sum_{i=1}^n y_i \mathbf{A}_i \right\rangle.$$

Applying minimax theorem, this shows that there exists an approximation solution  $\mathbf{X}$  in  $\text{conv}(\mathbf{X}_i)$  for the primal problem. Hence, we can restrict the primal SDP on the polytope  $\text{conv}(\mathbf{X}_i)$ , this reduces the primal SDP into a linear program which can be solved very efficiently. This idea of getting primal/dual solution from the cutting plane method is quite general and is the main purpose of this example. As a by-product, we have a faster SDP solver in both primal and dual! We remark that this idea has been used as a heuristic to obtain [151] for getting the primal SDP solution and our contribution here is mainly the asymptotic time analysis.

We first show how to construct the separation oracle for SDP. For that we need to compute smallest eigenvector of a matrix. Below, for completeness we provide a folklore result showing we can do this using fast matrix multiplication.

**Lemma 9.3.3.** *Given a  $n \times n$  symmetric matrix  $\mathbf{Y}$  such that  $-R\mathbf{I} \preceq \mathbf{Y} \preceq R\mathbf{I}$ , for any  $\varepsilon > 0$ , with high probability in  $n$  in time  $O(n^{\omega+o(1)} \log^{O(1)}(R/\varepsilon))$  we can find a unit vector  $\mathbf{u}$  such that  $\mathbf{u}^T \mathbf{Y} \mathbf{u} \geq \lambda_{\max}(\mathbf{Y}) - \varepsilon$ .*

*Proof.* Let  $\mathbf{B} \stackrel{\text{def}}{=} \frac{1}{R} \mathbf{Y} + \mathbf{I}$ . Note that  $\mathbf{B} \succeq \mathbf{0}$ . Now, we consider the repeated squaring  $\mathbf{B}_0 = \mathbf{B}$  and  $\mathbf{B}_{k+1} = \frac{\mathbf{B}_k^2}{\text{Tr} \mathbf{B}_k^2}$ . Let  $0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  be the eigenvalues of  $\mathbf{B}$  and  $\mathbf{v}_i$  be the corresponding eigenvectors. Then, it is easy to see the the eigenvalues of  $\mathbf{B}_k$  are  $\frac{\lambda_i^{2^k}}{\sum_{i=1}^n \lambda_i^{2^k}}$ .

Let  $\mathbf{q}$  be a random unit vector and  $\mathbf{r} \stackrel{\text{def}}{=} \mathbf{B}_k \mathbf{q}$ . Now  $\mathbf{q} = \sum \alpha_i \mathbf{v}_i$  for some  $\alpha_i$  such that  $\sum \alpha_i^2 = 1$ . Letting

$$\mathbf{p} = \frac{\sum_{\lambda_i > (1-\delta)\lambda_n} \alpha_i \lambda_i^{2^k} \mathbf{v}_i}{\sum_{i=1}^n \lambda_i^{2^k}}$$

we have

$$\|\mathbf{r} - \mathbf{p}\|_2 = \left\| \frac{\sum_{\lambda_i \leq (1-\delta)\lambda_n} \alpha_i \lambda_i^{2k} \mathbf{v}_i}{\sum_{i=1}^n \lambda_i^{2k}} \right\|_2 \leq \frac{\sum_{\lambda_i \leq (1-\delta)\lambda_n} \lambda_i^{2k}}{\sum_{i=1}^n \lambda_i^{2k}} \leq (1-\delta)^{2k} n.$$

Letting  $k = \log_2 \left( \frac{\log(n^{3/2}/\delta)}{\delta} \right)$ , we have  $\|\mathbf{r} - \mathbf{p}\|_2 \leq \delta/\sqrt{n}$ . Since  $\mathbf{0} \preceq \mathbf{B} \preceq 2\mathbf{I}$ , we have

$$\begin{aligned} \sqrt{\mathbf{r}^T \mathbf{B} \mathbf{r}} &\geq \sqrt{\mathbf{p}^T \mathbf{B} \mathbf{p}} - \sqrt{(\mathbf{r} - \mathbf{p})^T \mathbf{B} (\mathbf{r} - \mathbf{p})} \\ &\geq \sqrt{\mathbf{p}^T \mathbf{B} \mathbf{p}} - 2\delta/\sqrt{n}. \end{aligned}$$

Note that  $\mathbf{p}$  involves only eigenvectors between  $(1-\delta)\lambda_n$  to  $\lambda_n$ . Hence, we have

$$\sqrt{\mathbf{r}^T \mathbf{B} \mathbf{r}} \geq \sqrt{(1-\delta)\lambda_n} \|\mathbf{p}\|_2 - 2\delta/\sqrt{n}.$$

With constant probability, we have  $\alpha_n = \Omega(1/\sqrt{n})$ . Hence, we have  $\|\mathbf{p}\|_2 = \Omega(1/\sqrt{n})$ . Using  $\mathbf{B} \preceq 2\mathbf{I}$  and  $\|\mathbf{p}\|_2 \geq \|\mathbf{r}\|_2 - \delta/\sqrt{n}$  we have that so long as  $\delta$  is a small enough universal constant

$$\begin{aligned} \frac{\sqrt{\mathbf{r}^T \mathbf{B} \mathbf{r}}}{\|\mathbf{r}\|_2} &\geq \frac{\sqrt{(1-\delta)\lambda_n} \|\mathbf{p}\|_2 - 2\delta/\sqrt{n}}{\|\mathbf{p}\|_2 + \delta/\sqrt{n}} \\ &= (1 - O(\delta)) \sqrt{\lambda_n} - O(\delta) \\ &= \sqrt{\lambda_n} - O(\delta\sqrt{R}). \end{aligned}$$

Therefore, we have  $\frac{\mathbf{r}^T \mathbf{Y} \mathbf{r}}{\|\mathbf{r}\|_2^2} \geq \lambda_{\max}(\mathbf{Y}) - O(R\delta)$ . Hence, we can find vector  $\mathbf{r}$  by computing  $k$  matrix multiplications. [66] showed that fast matrix multiplication is stable under Frobenius norm, i.e., for any  $\eta > 0$ , using  $O(\log(n/b))$  bits, we can find  $\mathbf{C}$  such that  $\|\mathbf{C} - \mathbf{A}\mathbf{B}\|_F \leq \frac{1}{b} \|\mathbf{A}\| \|\mathbf{B}\|$  in time  $O(n^{\omega+\eta})$  where  $\omega$  is the matrix multiplicative constant. Hence, this algorithm takes only  $O(n^{\omega+o(1)} \log^{O(1)}(\delta^{-1}))$  time. The result follows from renormalizing the vector  $\mathbf{r}$ , repeating the algorithm  $O(\log n)$  times to boost the probability and taking  $\delta = \Omega(\varepsilon/R)$ .  $\square$

The following lemma shows how to compute a separation for  $f_K$  defined in (9.7).

**Lemma 9.3.4.** *Suppose that  $\|\mathbf{A}_i\|_F \leq M$  and  $\|\mathbf{C}\|_F \leq M$ . For any  $0 < \varepsilon < 1$  and  $\mathbf{y}$  with  $\|\mathbf{y}\|_2 = O(L)$ , with high probability in  $m$ , we can compute a  $(\varepsilon, \varepsilon)$ -separation of  $f_K$  on  $\{\|\mathbf{x}\|_2 \leq L\}$  at  $\mathbf{y}$  in time  $O(S + m^{\omega+o(1)} \log^{O(1)}(nKML/\varepsilon))$  where where  $S$  is the sparsity of the problem defined as  $\text{nnz}(\mathbf{C}) + \sum_{i=1}^n \text{nnz}(\mathbf{A}_i)$ .*

*Proof.* Note that  $-O(nML)\mathbf{I} \preceq \mathbf{C} - \sum_{i=1}^n y_i \mathbf{A}_i \preceq O(nML)\mathbf{I}$ . Using Lemma 9.3.3, we can find a vector  $\mathbf{v}$  with  $\|\mathbf{v}\|_2 = K$  in time  $O(m^{\omega+o(1)} \log^{O(1)}(nKML/\delta))$  such that

$$\mathbf{v}^T \left( \mathbf{C} - \sum_{i=1}^n y_i \mathbf{A}_i \right) \mathbf{v} \geq \max_{\text{Tr} \mathbf{X} \leq K, \mathbf{X} \succeq \mathbf{0}} \left\langle \mathbf{X}, \mathbf{C} - \sum_{i=1}^n y_i \mathbf{A}_i \right\rangle - \delta. \quad (9.8)$$

In other words, we have a  $\delta$ -optimization oracle for the function  $f_K$ . Lemma 9.2.6 shows this yields a  $\delta$ -subgradient oracle and Lemma 9.2.4 then shows this yields a  $(O(\sqrt{\delta L}), O(\sqrt{\delta L}))$ -separation oracle on the set  $\{\|\mathbf{x}\|_2 \leq L\}$ . By picking  $\delta = \varepsilon^2/L$ , we have the promised oracle.  $\square$

With the separation oracle in hand, we are ready to give the algorithm for SDP:

**Theorem 9.3.5.** *Given a primal-dual semidefinite programming problem in the form (9.4) and (9.5), suppose that for some  $M \geq 1$  we have*

1.  $\|b\|_2 \leq M$ ,  $\|C\|_F \leq M$  and  $\|A_i\|_F \leq M$  for all  $i$ .
2. The primal feasible set lies inside the region  $\text{Tr} \mathbf{X} \leq M$ .
3. The dual feasible set lies inside the region  $\|\mathbf{y}\|_\infty \leq M$ .

Let  $\text{OPT}$  be the optimum solution of (9.4) and (9.5). Then, with high probability, we can find  $\mathbf{X}$  and  $\mathbf{y}$  such that

1.  $\mathbf{X} \succeq \mathbf{0}$ ,  $\text{Tr} \mathbf{X} = O(M)$ ,  $\sum_i |b_i - \langle \mathbf{X}, \mathbf{A}_i \rangle| \leq \varepsilon$  for all  $i$  and  $C \bullet \mathbf{X} \geq \text{OPT} - \varepsilon$ .
2.  $\|\mathbf{y}\|_\infty = O(M)$ ,  $\sum_{i=1}^n y_i \mathbf{A}_i \succeq C - \varepsilon \mathbf{I}$  and  $\mathbf{b}^T \mathbf{y} \leq \text{OPT} + \varepsilon$ .

in expected time  $O\left(\left(nS + n^3 + nm^{\omega+o(1)}\right) \log^{O(1)}\left(\frac{nM}{\varepsilon}\right)\right)$  where  $S$  is the sparsity of the problem defined as  $\text{nnz}(C) + \sum_{i=1}^n \text{nnz}(A_i)$  and  $\omega$  is the fast matrix multiplication constant.

*Proof.* Let  $K \geq M$  be some parameter to be determined. Since the primal feasible set is lies inside the region  $\text{Tr} \mathbf{X} \leq M \leq K$ , we have

$$\begin{aligned} \min_{\sum_{i=1}^n y_i \mathbf{A}_i \succeq C} \mathbf{b}^T \mathbf{y} &= \max_{\mathbf{X} \succeq \mathbf{0}, \text{Tr} \mathbf{X} \leq K, \mathbf{A}_i \bullet \mathbf{X} = b_i} C \bullet \mathbf{X} \\ &= \max_{\mathbf{X} \succeq \mathbf{0}, \text{Tr} \mathbf{X} \leq K} \min_{\mathbf{y}} C \bullet \mathbf{X} - \sum_i y_i (\mathbf{A}_i \bullet \mathbf{X} - b_i) \\ &= \min_{\mathbf{y}} \max_{\mathbf{X} \succeq \mathbf{0}, \text{Tr} \mathbf{X} \leq K} \left( \mathbf{b}^T \mathbf{y} + (C - \sum_i y_i \mathbf{A}_i) \bullet \mathbf{X} \right) \\ &= \min_{\mathbf{y}} f_K(\mathbf{y}). \end{aligned}$$

Lemma 9.3.4 shows that it takes  $\text{SO}_{\delta, \delta}(f_K) = O(S + m^{\omega+o(1)} \log(nKML/\delta))$  time to compute a  $(\delta, \delta)$ -separation oracle of  $f_K$  for any point  $\mathbf{y}$  with  $\|\mathbf{y}\|_\infty = O(L)$  where  $L$  is some parameter with  $L \geq M$ . Taking the radius  $R = L$ , Theorem 9.3.2 shows that it takes  $O\left(n \text{SO}_{\delta, \delta}(f_K) \log\left(\frac{n}{\alpha}\right) + n^3 \log^{O(1)}\left(\frac{n}{\alpha}\right)\right)$  expected time with  $\delta = \Theta(\alpha n^{-3/2} L)$  to find  $\mathbf{y}$  such that

$$f_K(\mathbf{y}) - \min_{\|\mathbf{y}\|_\infty \leq L} f_K(\mathbf{y}) \leq \delta + \alpha \left( \max_{\|\mathbf{y}\|_\infty \leq L} f_K(\mathbf{y}) - \min_{\|\mathbf{y}\|_\infty \leq L} f_K(\mathbf{y}) \right) \leq \delta + 2\alpha(nML + 2nKML).$$

Picking  $\alpha = \frac{\varepsilon}{7nMKL}$ , we have  $f_K(\mathbf{y}) \leq \min_{\mathbf{y}} f_K(\mathbf{y}) + \varepsilon$ . Therefore,

$$\mathbf{b}^T \mathbf{y} + K \max(\lambda_{\max}(C - \sum_{i=1}^n y_i \mathbf{A}_i), 0) \leq \text{OPT} + \varepsilon.$$

Let  $\beta = \max(\lambda_{\max}(C - \sum_{i=1}^n y_i \mathbf{A}_i), 0)$ . Then, we have that  $\sum_{i=1}^n y_i \mathbf{A}_i \succeq C - \beta \mathbf{I}$  and

$$\begin{aligned} \mathbf{b}^T \mathbf{y} &\geq \min_{\sum_{i=1}^n y_i \mathbf{A}_i \succeq C - \beta \mathbf{I}} \mathbf{b}^T \mathbf{y} \\ &= \max_{\mathbf{X} \succeq \mathbf{0}, \mathbf{A}_i \bullet \mathbf{X} = b_i} (C - \beta \mathbf{I}) \bullet \mathbf{X} \\ &\geq \text{OPT} - \beta M \end{aligned}$$

because  $\text{Tr} \mathbf{X} \leq M$ . Hence, we have

$$\text{OPT} - \beta M + \beta K \leq \mathbf{b}^T \mathbf{y} + K \max(\lambda_{\max}(C - \sum_{i=1}^n y_i \mathbf{A}_i), 0) \leq \text{OPT} + \varepsilon$$

Putting  $K = M + 1$ , we have  $\beta \leq \varepsilon$ . Thus,

$$\sum_{i=1}^n y_i \mathbf{A}_i \succeq \mathbf{C} - \varepsilon \mathbf{I}.$$

This gives the result for the dual with the running time  $O\left(\left(nS + n^3 + nm^{\omega+o(1)}\right) \log^{O(1)}\left(\frac{nML}{\varepsilon}\right)\right)$ .

Our Cutting Plane Method accesses the sub-problem

$$\max_{\mathbf{X} \succeq \mathbf{0}, \text{Tr} \mathbf{X} \leq K} \left( \mathbf{C} - \sum_i y_i \mathbf{A}_i \right) \bullet \mathbf{X}$$

only through the separation oracle. Let  $\mathbf{z}$  be the output of our Cutting Plane Method and  $\{\mathbf{v}_i \mathbf{v}_i^T\}_{i=1}^{O(n)}$  be the matrices used to construct the separation for the  $O(n)$  hyperplanes the algorithm maintains at the end. Let  $\mathbf{u}$  be the maximum eigenvector of  $\mathbf{C} - \sum_{i=1}^n z_i \mathbf{A}_i$ . Now, we consider a realization of  $f_K$

$$\tilde{f}_K(\mathbf{y}) = \mathbf{b}^T \mathbf{y} + \max_{\mathbf{X} \in \text{conv}(K\mathbf{u}\mathbf{u}^T, \mathbf{v}_i \mathbf{v}_i^T)} \left\langle \mathbf{X}, \mathbf{C} - \sum_{i=1}^n y_i \mathbf{A}_i \right\rangle.$$

Since applying our Cutting Plane Method to either  $f_K$  or  $\tilde{f}_K$  gives the same result, the correctness of the our Cutting Plane Method shows that

$$\tilde{f}_K(\mathbf{z}) \leq \min_{\|\mathbf{y}\|_\infty \leq L} \tilde{f}_K(\mathbf{y}) + \varepsilon.$$

Note that the function  $\tilde{f}_K$  is defined such that  $\tilde{f}_K(\mathbf{z}) = f_K(\mathbf{z})$ . Hence, we have

$$\min_{\|\mathbf{y}\|_\infty \leq L} f_K(\mathbf{y}) \leq f_K(\mathbf{z}) \leq \tilde{f}_K(\mathbf{z}) \leq \min_{\|\mathbf{y}\|_\infty \leq L} \tilde{f}_K(\mathbf{y}) + \varepsilon.$$

Also, note that  $\tilde{f}_K(\mathbf{x}) \leq f_K(\mathbf{x})$  for all  $\mathbf{x}$ . Hence, we have

$$\min_{\|\mathbf{y}\|_\infty \leq L} f_K(\mathbf{y}) - \varepsilon \leq \min_{\|\mathbf{y}\|_\infty \leq L} \tilde{f}_K(\mathbf{y}) \leq \min_{\|\mathbf{y}\|_\infty \leq L} f_K(\mathbf{y}).$$

Now, we consider the primal version of  $\tilde{f}$ , namely

$$g(\mathbf{X}) \stackrel{\text{def}}{=} \min_{\|\mathbf{y}\|_\infty \leq L} \mathbf{b}^T \mathbf{y} + \left\langle \mathbf{X}, \mathbf{C} - \sum_{i=1}^n y_i \mathbf{A}_i \right\rangle.$$

Sion's minimax theorem [239] shows that

$$\text{OPT} \geq \max_{\mathbf{X} \in \text{conv}(K\mathbf{u}\mathbf{u}^T, \mathbf{v}_i \mathbf{v}_i^T)} g(\mathbf{X}) = \min_{\|\mathbf{y}\|_\infty \leq L} \tilde{f}(\mathbf{y}) \geq \text{OPT} - \varepsilon.$$

Therefore, to get the primal solution, we only need to find  $\mathbf{u}$  by Lemma 9.3.3 and solve the maximization problem on  $g$ . Note that

$$\begin{aligned} g(\mathbf{X}) &= \min_{\|\mathbf{y}\|_\infty \leq L} \sum_{i=1}^n y_i (b_i - \langle \mathbf{X}, \mathbf{A}_i \rangle) + \langle \mathbf{X}, \mathbf{C} \rangle \\ &= -L \sum_i |b_i - \langle \mathbf{X}, \mathbf{A}_i \rangle| + \langle \mathbf{X}, \mathbf{C} \rangle. \end{aligned}$$

For notation simplicity, we write  $K\mathbf{u}\mathbf{u}^T = \mathbf{v}_0 \mathbf{v}_0^T$ . Then,  $\mathbf{X} = \sum_{j=0}^{O(n)} \alpha_j \mathbf{v}_j \mathbf{v}_j^T$  for some  $\sum \alpha_j = 1$  and



$\alpha_j \geq 0$ . Substituting this into the function  $g$ , we have

$$g(\boldsymbol{\alpha}) = -L \sum_j \left| b_i - \sum_j \alpha_j \mathbf{v}_j^T \mathbf{A}_i \mathbf{v}_j \right| + \sum_j \alpha_j \mathbf{v}_j^T \mathbf{C} \mathbf{v}_j.$$

Hence, this can be easily written as a linear program with  $O(n)$  variables and  $O(n)$  constraints in time  $O(nS)$ . Now, we can apply interior point method to find  $\boldsymbol{\alpha}$  such that

$$g(\boldsymbol{\alpha}) \geq \max_{\mathbf{X} \in \text{conv}(K \mathbf{u} \mathbf{u}^T, \mathbf{v}_i \mathbf{v}_i^T)} g(\mathbf{X}) - \varepsilon \geq \text{OPT} - 2\varepsilon.$$

Let the corresponding approximate solution be  $\widetilde{\mathbf{X}} = \sum \alpha_j \mathbf{v}_j \mathbf{v}_j^T$ . Then, we have

$$\langle \widetilde{\mathbf{X}}, \mathbf{C} \rangle - L \sum_i |b_i - \langle \widetilde{\mathbf{X}}, \mathbf{A}_i \rangle| \geq \text{OPT} - 2\varepsilon.$$

Now, we let  $\tilde{b}_i = \langle \widetilde{\mathbf{X}}, \mathbf{A}_i \rangle$ . Then, we note that

$$\begin{aligned} \langle \widetilde{\mathbf{X}}, \mathbf{C} \rangle &\leq \max_{\mathbf{X} \succeq \mathbf{0}, \mathbf{A}_i \bullet \mathbf{X} = \tilde{b}_i} \mathbf{C} \bullet \mathbf{X} \\ &= \min_{\sum_{i=1}^n y_i \mathbf{A}_i \succeq \mathbf{C}} \tilde{b}_i^T \mathbf{y} \\ &\leq \text{OPT} + M \sum_i |b_i - \langle \widetilde{\mathbf{X}}, \mathbf{A}_i \rangle| \end{aligned}$$

because  $\|\mathbf{y}\|_\infty \leq M$ . Hence, we have

$$\text{OPT} + (M - L) \sum_i |b_i - \langle \widetilde{\mathbf{X}}, \mathbf{A}_i \rangle| \geq \langle \widetilde{\mathbf{X}}, \mathbf{C} \rangle - L \sum_i |b_i - \langle \widetilde{\mathbf{X}}, \mathbf{A}_i \rangle| \geq \text{OPT} - 2\varepsilon.$$

Now, we put  $L = M + 2$ , we have

$$\sum_i |b_i - \langle \widetilde{\mathbf{X}}, \mathbf{A}_i \rangle| \leq \varepsilon.$$

This gives the result for the primal. Note that it only takes  $O(n^{5/2} \log^{O(1)}(nM/\varepsilon))$  to solve a linear program with  $O(n)$  variables and  $O(n)$  constraints because we have an explicit interior point deep inside the feasible set, i.e.  $\alpha_i = \frac{1}{m}$  for some parameter  $m$  [164].<sup>1</sup> Hence, the running time is dominated by the cost of cutting plane method which is  $O\left((nS + n^3 + nm^{\omega+o(1)}) \log^{O(1)}\left(\frac{nM}{\varepsilon}\right)\right)$  by putting  $L = M + 2$ .  $\square$

We leave it as an open problem if it is possible to improve this result by reusing the computation in the separation oracle and achieve a running time of  $O\left((nS + n^3 + nm^2) \log^{O(1)}\left(\frac{nM}{\varepsilon}\right)\right)$ .

## ■ 9.4 Intersection of Convex Sets

In this section we introduce a general technique to optimize a linear function over the intersection of two convex sets, whenever the linear optimization problem on each of them can be done efficiently. At the very high level, this is accomplished by applying cutting plane to a suitably regularized version of the problem. In Section 9.4.1 we present the technique and in the remaining sections we provide

<sup>1</sup>Without this, the running time of interior point method depends on the bit complexity of the linear programs.

several applications including, matroid intersection (Section 9.4.2), submodular flow (Section 9.4.3), and minimizing over the intersection of an affine subspace and a convex set (Section 9.4.4).

### ■ 9.4.1 The Technique

Throughout this section we consider variants of the following general optimization problem

$$\max_{\mathbf{x} \in K_1 \cap K_2} \langle \mathbf{c}, \mathbf{x} \rangle \quad (9.9)$$

where  $\mathbf{x}, \mathbf{c} \in \mathbb{R}^n$ ,  $K_1$  and  $K_2$  are convex subsets of  $\mathbb{R}^n$ . We assume that

$$\max_{\mathbf{x} \in K_1} \|\mathbf{x}\|_2 < M, \quad \max_{\mathbf{x} \in K_2} \|\mathbf{x}\|_2 < M, \quad \|\mathbf{c}\|_2 \leq M \quad (9.10)$$

for some constant  $M \geq 1$  and we assume that

$$K_1 \cap K_2 \neq \emptyset. \quad (9.11)$$

Instead of a separation oracle, we assume that  $K_1$  and  $K_2$  each have optimization oracles (see Section 9.2).

To solve this problem we first introduce a relaxation for the problem (9.9) that we can optimize efficiently. Because we have only the optimization oracles for  $K_1$  and  $K_2$ , we simply have variables  $\mathbf{x}$  and  $\mathbf{y}$  for each of them in the objective. Since the output should (approximately) be in the intersection of  $K_1$  and  $K_2$ , a regularization term  $-\frac{\lambda}{2}\|\mathbf{x} - \mathbf{y}\|_2^2$  is added to force  $\mathbf{x} \approx \mathbf{y}$  where  $\lambda$  is a large number to be determined later. Furthermore, we add terms to make the problem strong concave.

**Lemma 9.4.1.** *Assume (9.10) and (9.11). For  $\lambda \geq 1$ , let*

$$f_\lambda(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \frac{1}{2} \langle \mathbf{c}, \mathbf{x} \rangle + \frac{1}{2} \langle \mathbf{c}, \mathbf{y} \rangle - \frac{\lambda}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 - \frac{1}{2\lambda} \|\mathbf{x}\|_2^2 - \frac{1}{2\lambda} \|\mathbf{y}\|_2^2. \quad (9.12)$$

*There is an unique maximizer  $(\mathbf{x}_\lambda, \mathbf{y}_\lambda)$  for the problem  $\max_{\mathbf{x} \in K_1, \mathbf{y} \in K_2} f_\lambda(\mathbf{x}, \mathbf{y})$ . The maximizer  $(\mathbf{x}_\lambda, \mathbf{y}_\lambda)$  is a good approximation of the solution of (9.9), i.e.  $\|\mathbf{x}_\lambda - \mathbf{y}_\lambda\|_2^2 \leq \frac{6M^2}{\lambda}$  and*

$$\max_{\mathbf{x} \in K_1 \cap K_2} \langle \mathbf{c}, \mathbf{x} \rangle \leq f_\lambda(\mathbf{x}_\lambda, \mathbf{y}_\lambda) + \frac{M^2}{\lambda}. \quad (9.13)$$

*Proof.* Let  $\mathbf{x}^*$  be a maximizer of  $\max_{\mathbf{x} \in K_1 \cap K_2} \langle \mathbf{c}, \mathbf{x} \rangle$ . By assumption (9.10),  $\|\mathbf{x}^*\|_2 \leq M$ , and therefore

$$f_\lambda(\mathbf{x}^*, \mathbf{x}^*) = \langle \mathbf{c}, \mathbf{x}^* \rangle - \frac{\|\mathbf{x}^*\|_2^2}{\lambda} \geq \max_{\mathbf{x} \in K_1 \cap K_2} \langle \mathbf{c}, \mathbf{x} \rangle - \frac{M^2}{\lambda}. \quad (9.14)$$

This shows (9.13). Since  $f_\lambda$  is strongly concave in  $\mathbf{x}$  and  $\mathbf{y}$ , there is a unique maximizer  $(\mathbf{x}_\lambda, \mathbf{y}_\lambda)$ . Let  $\text{OPT}_\lambda = f_\lambda(\mathbf{x}_\lambda, \mathbf{y}_\lambda)$ . Then, we have

$$\begin{aligned} \text{OPT}_\lambda &\leq \frac{1}{2} \|\mathbf{c}\|_2 \|\mathbf{x}_\lambda\|_2 + \frac{1}{2} \|\mathbf{c}\|_2 \|\mathbf{y}_\lambda\|_2 - \frac{\lambda}{2} \|\mathbf{x}_\lambda - \mathbf{y}_\lambda\|_2^2 \\ &\leq \frac{M^2}{2} + \frac{M^2}{2} - \frac{\lambda}{2} \|\mathbf{x}_\lambda - \mathbf{y}_\lambda\|_2^2. \end{aligned}$$

On the other hand, using  $\lambda \geq 1$ , (9.14) shows that

$$\text{OPT}_\lambda \geq f_\lambda(\mathbf{x}^*, \mathbf{x}^*) \geq \max_{\mathbf{x} \in K_1 \cap K_2} \langle \mathbf{c}, \mathbf{x} \rangle - \frac{M^2}{\lambda} \geq -2M^2.$$

Hence, we have

$$\|\mathbf{x}_\lambda - \mathbf{y}_\lambda\|_2^2 \leq \frac{2(M^2 - \text{OPT}_\lambda)}{\lambda} \leq \frac{6M^2}{\lambda}. \quad (9.15)$$

□

Now we write  $\max f_\lambda(\mathbf{x}, \mathbf{y})$  as a max-min problem. The reason for doing this is that the dual approximate solution is much easier to obtain and there is a way to read off a primal approximate solution from a dual approximate solution. This is analogous to the idea in [161] which showed how to convert a cut solution to a flow solution by adding regularization terms into the problem.

**Lemma 9.4.2.** *Assume (9.10) and (9.11). Let  $\lambda \geq 2$ . For any  $\mathbf{x} \in K_1$  and  $\mathbf{y} \in K_2$ , the function  $f_\lambda$  can be represented as*

$$f_\lambda(\mathbf{x}, \mathbf{y}) = \min_{(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) \in \Omega} g_\lambda(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) \quad (9.16)$$

where  $\Omega = \{(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) : \|\boldsymbol{\theta}_1\|_2 \leq 2M, \|\boldsymbol{\theta}_2\|_2 \leq M, \|\boldsymbol{\theta}_3\|_2 \leq M\}$  and

$$g_\lambda(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) = \left\langle \frac{\mathbf{c}}{2} + \lambda\boldsymbol{\theta}_1 + \frac{\boldsymbol{\theta}_2}{\lambda}, \mathbf{x} \right\rangle + \left\langle \frac{\mathbf{c}}{2} - \lambda\boldsymbol{\theta}_1 + \frac{\boldsymbol{\theta}_3}{\lambda}, \mathbf{y} \right\rangle + \frac{\lambda}{2}\|\boldsymbol{\theta}_1\|_2^2 + \frac{1}{2\lambda}\|\boldsymbol{\theta}_2\|_2^2 + \frac{1}{2\lambda}\|\boldsymbol{\theta}_3\|_2^2. \quad (9.17)$$

Let  $h_\lambda(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) = \max_{\mathbf{x} \in K_1, \mathbf{y} \in K_2} g_\lambda(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3)$ . For any  $(\boldsymbol{\theta}'_1, \boldsymbol{\theta}'_2, \boldsymbol{\theta}'_3)$  such that  $h_\lambda(\boldsymbol{\theta}'_1, \boldsymbol{\theta}'_2, \boldsymbol{\theta}'_3) \leq \min_{(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) \in \Omega} h_\lambda(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) + \varepsilon$ , we know  $\mathbf{z} = \frac{1}{2}(\boldsymbol{\theta}'_2 + \boldsymbol{\theta}'_3)$  satisfies

$$\max_{\mathbf{x} \in K_1 \cap K_2} \langle \mathbf{c}, \mathbf{x} \rangle \leq \langle \mathbf{c}, \mathbf{z} \rangle + \frac{20M^2}{\lambda} + 20\lambda^3\varepsilon.$$

and  $\|\mathbf{z} - \mathbf{x}_\lambda\|_2 + \|\mathbf{z} - \mathbf{y}_\lambda\|_2 \leq 4\sqrt{2\lambda\varepsilon} + \sqrt{\frac{6M^2}{\lambda}}$  where  $(\mathbf{x}_\lambda, \mathbf{y}_\lambda)$  is the unique maximizer for the problem  $\max_{\mathbf{x} \in K_1, \mathbf{y} \in K_2} f_\lambda(\mathbf{x}, \mathbf{y})$ .

*Proof.* Note that for any  $\|\boldsymbol{\zeta}\|_2 \leq \alpha$ , we have

$$-\frac{1}{2}\|\boldsymbol{\zeta}\|_2^2 = \min_{\|\boldsymbol{\theta}\|_2 \leq \alpha} \langle \boldsymbol{\theta}, \boldsymbol{\zeta} \rangle + \frac{1}{2}\|\boldsymbol{\theta}\|_2^2$$

Using this and (9.10), we have (9.16) for all  $\mathbf{x} \in K_1$  and  $\mathbf{y} \in K_2$  as desired. Since  $\Omega$  is closed and bounded set and the function  $g_\lambda$  is concave in  $(\mathbf{x}, \mathbf{y})$  and convex in  $(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3)$ , Sion's minimax theorem [239] shows that

$$\max_{\mathbf{x} \in K_1, \mathbf{y} \in K_2} f_\lambda(\mathbf{x}, \mathbf{y}) = \min_{(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) \in \Omega} h_\lambda(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) \quad (9.18)$$

Since  $f_\lambda$  is strongly concave, there is a unique maximizer  $(\mathbf{x}_\lambda, \mathbf{y}_\lambda)$  of  $f_\lambda$ . Since  $h_\lambda$  is strongly convex, there is a unique minimizer  $(\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \boldsymbol{\theta}_3^*)$ . By the definition of  $f_\lambda$  and  $h_\lambda$ , we have

$$h_\lambda(\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \boldsymbol{\theta}_3^*) \geq g_\lambda(\mathbf{x}_\lambda, \mathbf{y}_\lambda, \boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \boldsymbol{\theta}_3^*) \geq f_\lambda(\mathbf{x}_\lambda, \mathbf{y}_\lambda) \quad .$$

Using (9.18), the equality above holds and hence  $(\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \boldsymbol{\theta}_3^*)$  is the minimizer of  $g_\lambda(\mathbf{x}_\lambda, \mathbf{y}_\lambda, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3)$  over  $(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3)$ . Since the domain  $\Omega$  is large enough that  $(\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \boldsymbol{\theta}_3^*)$  is an interior point in  $\Omega$ , the optimality condition of  $g_\lambda$  shows that we have  $\boldsymbol{\theta}_2^* = \mathbf{x}_\lambda$  and  $\boldsymbol{\theta}_3^* = \mathbf{y}_\lambda$ .

Since  $h_\lambda$  is  $\frac{1}{\lambda}$  strongly convex, we have  $\|\boldsymbol{\theta}'_1 - \boldsymbol{\theta}_1^*\|_2^2 + \|\boldsymbol{\theta}'_2 - \boldsymbol{\theta}_2^*\|_2^2 + \|\boldsymbol{\theta}'_3 - \boldsymbol{\theta}_3^*\|_2^2 \leq 2\lambda\varepsilon$  (Fact 9.2.3). Since  $\boldsymbol{\theta}_2^* = \mathbf{x}_\lambda$  and  $\boldsymbol{\theta}_3^* = \mathbf{y}_\lambda$ , we have

$$\|\boldsymbol{\theta}'_2 - \mathbf{x}_\lambda\|_2^2 + \|\boldsymbol{\theta}'_3 - \mathbf{y}_\lambda\|_2^2 \leq 2\lambda\varepsilon. \quad (9.19)$$

Therefore, we have  $\|\mathbf{x}_\lambda - \mathbf{y}_\lambda\|_2 \geq \|\boldsymbol{\theta}'_2 - \boldsymbol{\theta}'_3\|_2 - 2\sqrt{2\lambda\varepsilon}$ ,  $\|\mathbf{x}_\lambda\|_2 \geq \|\boldsymbol{\theta}'_2\|_2 - \sqrt{2\lambda\varepsilon}$  and  $\|\mathbf{y}_\lambda\|_2 \geq \|\boldsymbol{\theta}'_3\|_2 - \sqrt{2\lambda\varepsilon}$ . Using these,  $\|\mathbf{x}_\lambda\|_2 \leq M$  and  $\|\mathbf{y}_\lambda\|_2 \leq M$ , we have

$$\begin{aligned}
f_\lambda(\boldsymbol{\theta}'_2, \boldsymbol{\theta}'_3) &= \frac{1}{2} \langle \mathbf{c}, \boldsymbol{\theta}'_2 \rangle + \frac{1}{2} \langle \mathbf{c}, \boldsymbol{\theta}'_3 \rangle - \frac{\lambda}{2} \|\boldsymbol{\theta}'_2 - \boldsymbol{\theta}'_3\|_2^2 - \frac{1}{2\lambda} \|\boldsymbol{\theta}'_2\|_2^2 - \frac{1}{2\lambda} \|\boldsymbol{\theta}'_3\|_2^2 \\
&\geq \frac{1}{2} \langle \mathbf{c}, \mathbf{x}_\lambda \rangle + \frac{1}{2} \langle \mathbf{c}, \mathbf{y}_\lambda \rangle - M\sqrt{2\lambda\varepsilon} \\
&\quad - \frac{\lambda}{2} \left( \|\mathbf{x}_\lambda - \mathbf{y}_\lambda\|_2 + 2\sqrt{2\lambda\varepsilon} \right)^2 \\
&\quad - \frac{1}{2\lambda} \left( \|\mathbf{x}_\lambda\|_2 + \sqrt{2\lambda\varepsilon} \right)^2 - \frac{1}{2\lambda} \left( \|\mathbf{y}_\lambda\|_2 + \sqrt{2\lambda\varepsilon} \right)^2 \\
&= \frac{1}{2} \langle \mathbf{c}, \mathbf{x}_\lambda \rangle + \frac{1}{2} \langle \mathbf{c}, \mathbf{y}_\lambda \rangle - \frac{\lambda}{2} \|\mathbf{x}_\lambda - \mathbf{y}_\lambda\|_2^2 - \frac{1}{2\lambda} \|\mathbf{x}_\lambda\|_2^2 - \frac{1}{2\lambda} \|\mathbf{y}_\lambda\|_2^2 \\
&\quad - M\sqrt{2\lambda\varepsilon} - 2\lambda\sqrt{2\lambda\varepsilon} \|\mathbf{x}_\lambda - \mathbf{y}_\lambda\|_2 - 4\lambda^2\varepsilon \\
&\quad - \frac{1}{\lambda} \|\mathbf{x}_\lambda\|_2\sqrt{2\lambda\varepsilon} - \varepsilon - \frac{1}{\lambda} \|\mathbf{y}_\lambda\|_2\sqrt{2\lambda\varepsilon} - \varepsilon.
\end{aligned}$$

Using  $\|\mathbf{x}_\lambda - \mathbf{y}_\lambda\|_2 \leq \sqrt{\frac{6M^2}{\lambda}}$  (Lemma 9.4.1),  $\|\mathbf{x}_\lambda\|_2 < M$  and  $\|\mathbf{y}_\lambda\|_2 < M$ , we have

$$\begin{aligned}
f_\lambda(\boldsymbol{\theta}'_2, \boldsymbol{\theta}'_3) &\geq f_\lambda(\mathbf{x}_\lambda, \mathbf{y}_\lambda) \\
&\quad - M\sqrt{2\lambda\varepsilon} - 2\lambda\sqrt{2\lambda\varepsilon} \|\mathbf{x}_\lambda - \mathbf{y}_\lambda\|_2 - 4\lambda^2\varepsilon \\
&\quad - \frac{1}{\lambda} \|\mathbf{x}_\lambda\|_2\sqrt{2\lambda\varepsilon} - \varepsilon - \frac{1}{\lambda} \|\mathbf{y}_\lambda\|_2\sqrt{2\lambda\varepsilon} - \varepsilon. \\
&\geq f_\lambda(\mathbf{x}_\lambda, \mathbf{y}_\lambda) \\
&\quad - M\sqrt{2\lambda\varepsilon} - 2\lambda\sqrt{12\varepsilon}M - 4\lambda^2\varepsilon \\
&\quad - 2M\sqrt{2\frac{\varepsilon}{\lambda}} - 2\varepsilon.
\end{aligned}$$

Since  $\lambda \geq 2$ , we have

$$f_\lambda(\boldsymbol{\theta}'_2, \boldsymbol{\theta}'_3) \geq f_\lambda(\mathbf{x}_\lambda, \mathbf{y}_\lambda) - 20M\lambda\sqrt{\varepsilon} - 10\lambda^2\varepsilon.$$

Let  $\mathbf{z} = \frac{\boldsymbol{\theta}'_2 + \boldsymbol{\theta}'_3}{2}$ . Lemma 9.4.1 shows that

$$\begin{aligned}
\max_{\mathbf{x} \in K_1 \cap K_2} \langle \mathbf{c}, \mathbf{x} \rangle &\leq \max_{\mathbf{x} \in K_1, \mathbf{y} \in K_2} f_\lambda(\mathbf{x}, \mathbf{y}) + \frac{M^2}{\lambda} \\
&\leq f_\lambda(\boldsymbol{\theta}'_2, \boldsymbol{\theta}'_3) + \frac{M^2}{\lambda} + 20M\lambda\sqrt{\varepsilon} + 10\lambda^2\varepsilon \\
&\leq \langle \mathbf{c}, \mathbf{z} \rangle + \frac{20M^2}{\lambda} + 20\lambda^3\varepsilon
\end{aligned}$$

because  $20M\lambda\sqrt{\varepsilon} \leq 10\frac{M^2}{\lambda} + 10\lambda^3\varepsilon$ . Furthermore, we have

$$\begin{aligned}
\|\mathbf{z} - \mathbf{x}_\lambda\|_2 + \|\mathbf{z} - \mathbf{y}_\lambda\|_2 &\leq \|\boldsymbol{\theta}'_2 - \mathbf{x}_\lambda\|_2 + \|\boldsymbol{\theta}'_3 - \mathbf{y}_\lambda\|_2 + \|\boldsymbol{\theta}'_2 - \boldsymbol{\theta}'_3\|_2 \\
&\leq 4\sqrt{2\lambda\varepsilon} + \sqrt{\frac{6M^2}{\lambda}}.
\end{aligned}$$

□

We now apply our cutting plane method to solve the optimization problem (9.9). First we show how to transform the optimization oracles for  $K_1$  and  $K_2$  to get a separation oracle for  $h_\lambda$ , with the

appropriate parameters.

**Lemma 9.4.3.** *Suppose we have a  $\varepsilon$ -optimization oracle for  $K_1$  and  $K_2$  for some  $0 < \varepsilon < 1$ . Then on the set  $\{\|\boldsymbol{\theta}\|_2 \leq D\}$ , we have a  $(O(\sqrt{\varepsilon\lambda D}), O(\sqrt{\varepsilon\lambda D}))$ -separation oracle for  $h_\lambda$  with time complexity  $\text{OO}_\varepsilon(K_1) + \text{OO}_\varepsilon(K_2)$ .*

*Proof.* Recall that the function  $h_\lambda$  is defined by

$$\begin{aligned} & h_\lambda(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) \\ &= \max_{\mathbf{x} \in K_1, \mathbf{y} \in K_2} \left( \left\langle \frac{\mathbf{c}}{2} + \lambda \boldsymbol{\theta}_1 + \frac{\boldsymbol{\theta}_2}{\lambda}, \mathbf{x} \right\rangle + \left\langle \frac{\mathbf{c}}{2} - \lambda \boldsymbol{\theta}_1 + \frac{\boldsymbol{\theta}_3}{\lambda}, \mathbf{y} \right\rangle + \frac{\lambda}{2} \|\boldsymbol{\theta}_1\|_2^2 + \frac{1}{2\lambda} \|\boldsymbol{\theta}_2\|_2^2 + \frac{1}{2\lambda} \|\boldsymbol{\theta}_3\|_2^2 \right) \\ &= \max_{\mathbf{x} \in K_1} \left\langle \frac{\mathbf{c}}{2} + \lambda \boldsymbol{\theta}_1 + \frac{\boldsymbol{\theta}_2}{\lambda}, \mathbf{x} \right\rangle + \max_{\mathbf{y} \in K_2} \left\langle \frac{\mathbf{c}}{2} - \lambda \boldsymbol{\theta}_1 + \frac{\boldsymbol{\theta}_3}{\lambda}, \mathbf{y} \right\rangle + \frac{\lambda}{2} \|\boldsymbol{\theta}_1\|_2^2 + \frac{1}{2\lambda} \|\boldsymbol{\theta}_2\|_2^2 + \frac{1}{2\lambda} \|\boldsymbol{\theta}_3\|_2^2. \end{aligned}$$

Lemma 9.2.6 shows how to compute the subgradient of functions of the form  $f(\mathbf{c}) = \max_{\mathbf{x} \in K} \langle \mathbf{c}, \mathbf{x} \rangle$  using the optimization oracle for  $K$ . The rest of the term are differentiable so its subgradient is just the gradient. Hence, by addition rule for subgradients (Fact 9.2.3), we have a  $O(\varepsilon\lambda)$ -subgradient oracle for  $f_\lambda$  using a  $O(\varepsilon)$ -optimization oracle for  $K_1$  and  $K_2$ . The result then follows from Lemma 9.2.4.  $\square$

**Theorem 9.4.4.** *Assume (9.10) and (9.11). Suppose that we have  $\varepsilon$ -optimization oracle for every  $\varepsilon > 0$ . For  $0 < \delta < 1$ , we can find  $\mathbf{z} \in \mathbb{R}^n$  such that*

$$\max_{\mathbf{x} \in K_1 \cap K_2} \langle \mathbf{c}, \mathbf{x} \rangle \leq \delta + \langle \mathbf{c}, \mathbf{z} \rangle$$

and  $\|\mathbf{z} - \mathbf{x}\|_2 + \|\mathbf{z} - \mathbf{y}\|_2 \leq \delta$  for some  $\mathbf{x} \in K_1$  and  $\mathbf{y} \in K_2$  in time

$$O\left(n(\text{OO}_\eta(K_1) + \text{OO}_\eta(K_2)) \log\left(\frac{nM}{\delta}\right) + n^3 \log^{O(1)}\left(\frac{nM}{\delta}\right)\right)$$

where  $\eta = \Omega\left(\left(\frac{\delta}{nM}\right)^{O(1)}\right)$ .

*Proof.* Setting  $\lambda = \frac{40M^2}{\delta^2}$  and  $\varepsilon = \frac{\delta^7}{10^7 M^6}$  in Lemma 9.4.2 we see that so long as we obtain any approximate solution  $(\boldsymbol{\theta}'_1, \boldsymbol{\theta}'_2, \boldsymbol{\theta}'_3)$  such that

$$h_\lambda(\boldsymbol{\theta}'_1, \boldsymbol{\theta}'_2, \boldsymbol{\theta}'_3) \leq \min_{(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) \in \Omega} h_\lambda(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) + \varepsilon,$$

then we obtain the point we want. To apply Theorem 9.3.2, we use

$$\tilde{h}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) = \begin{cases} h_\lambda(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) & \text{if } (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) \in \Omega \\ +\infty & \text{else} \end{cases}.$$

Lemma 9.4.3 shows that for any  $\gamma > 0$  we can obtain a  $(\gamma, \gamma)$ -separation oracle of  $h_\lambda(\boldsymbol{\theta})$  by using sufficiently accurate optimization oracles. Since  $\Omega$  is just a product of  $\ell^2$  balls, we can produce a separating hyperplane easily when  $(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) \notin \Omega$ . Hence, we can obtain a  $(\gamma, \gamma)$ -separation oracle of  $\tilde{h}(\boldsymbol{\theta})$ . For simplicity, we use  $\boldsymbol{\theta}$  to represent  $(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3)$ . Note that  $B_\infty(2M) \supseteq \Omega$  and therefore we can apply Theorem 9.3.2 with  $R = 2M$  to compute  $\boldsymbol{\theta}'$  such

$$\tilde{h}(\boldsymbol{\theta}') - \min_{\boldsymbol{\theta} \in \Omega} \tilde{h}(\boldsymbol{\theta}) \leq \gamma + \alpha \left( \max_{\boldsymbol{\theta} \in \Omega} \tilde{h}(\boldsymbol{\theta}) - \min_{\boldsymbol{\theta} \in \Omega} \tilde{h}(\boldsymbol{\theta}) \right)$$

in time  $O\left(n\text{SO}_{\gamma, \gamma} \log\left(\frac{n\kappa}{\alpha}\right) + n^3 \log^{O(1)}\left(\frac{n\kappa}{\alpha}\right)\right)$  where  $\gamma = \Omega\left(\alpha \text{MinWidth}(\Omega)/n^{O(1)}\right) = \Omega\left(\alpha M/n^{O(1)}\right)$

and  $\kappa = \frac{2M}{\text{MinWidth}(\Omega)} = O(1)$ . Using  $\lambda \geq 1$  and  $M \geq 1$ , we have

$$\max_{\boldsymbol{\theta} \in \Omega} \tilde{h}(\boldsymbol{\theta}) - \min_{\boldsymbol{\theta} \in \Omega} \tilde{h}(\boldsymbol{\theta}) \leq O(\lambda M^2) \leq O\left(\frac{M^4}{\delta^2}\right).$$

Setting  $\alpha = \Theta\left(\frac{\delta^9}{M^{10}}\right)$  with some small enough constant, we have that we can find  $\boldsymbol{\theta}'$  such that

$$\begin{aligned} h_\lambda(\boldsymbol{\theta}') &\leq \min_{\boldsymbol{\theta} \in P} h_\lambda(\boldsymbol{\theta}) + \gamma + \alpha O\left(\frac{M^4}{\delta^2}\right) \\ &= \min_{\boldsymbol{\theta} \in P} h_\lambda(\boldsymbol{\theta}) + O\left(\frac{\delta^7}{M^6}\right) \\ &= \min_{\boldsymbol{\theta} \in P} h_\lambda(\boldsymbol{\theta}) + \varepsilon \end{aligned}$$

in time  $O\left(n\text{SO}_{\gamma,\gamma} \log\left(\frac{nM}{\delta}\right) + n^3 \log^{O(1)}\left(\frac{nM}{\delta}\right)\right)$  where  $\gamma = \Omega\left(\left(\frac{\delta}{nM}\right)^{O(1)}\right)$ . Lemma 9.4.3 shows that the cost of  $(\gamma, \gamma)$ -separation oracle is just  $O(\text{OO}_\eta(K_1) + \text{OO}_\eta(K_2))$  where  $\eta = \Omega\left(\left(\frac{\delta}{nM}\right)^{O(1)}\right)$ .  $\square$

*Remark 9.4.5.* Note that the algorithm does not promise that we obtain a point close to  $K_1 \cap K_2$ . It only promises to give a point that is close to both some point in  $K_1$  and some point in  $K_2$ . It appears to the authors that a further assumption is needed to get a point close to  $K_1 \cap K_2$ . For example, if  $K_1$  and  $K_2$  are two almost parallel lines, it would be difficult to get an algorithm that does not depend on the angle. However, as far as we know, most algorithms tackling this problem are pseudo-polynomial and have polynomial dependence on the angle. Our algorithm depends on the logarithmic of the angle which is useful for combinatorial problems.

This reduction is very useful for problems in many areas including linear programming, semi-definite programming and algorithmic game theory. In the remainder of this section we demonstrate its power by applying it to classical combinatorial problems.

There is however one issue with applying our cutting plane algorithm to these problems. As with other convex optimization methods, only an approximately optimal solution is found. On the other hand, typically an exact solution is insisted in combinatorial optimization. To overcome this gap, we introduce the following lemma which (1) transforms the objective function so that there is only one optimal solution and (2) shows that an approximate solution is close to the optimal solution whenever it is unique. As we shall see in the next two subsections, this allows us to round an approximate solution to an optimal one.

**Lemma 9.4.6.** *Given a linear program  $\min_{\mathbf{Ax} \geq \mathbf{b}} \mathbf{c}^T \mathbf{x}$  where  $\mathbf{x}, \mathbf{c} \in \mathbb{Z}^n$ ,  $\mathbf{b} \in \mathbb{Z}^m$  and  $\mathbf{A} \in \mathbb{Z}^{m \times n}$ . Suppose  $\{\mathbf{Ax} \geq \mathbf{b}\}$  is an integral polytope (i.e. all extreme points are integral) contained in the set  $\{\|\mathbf{x}\|_\infty \leq M\}$ . Then we can find a random cost vector  $\mathbf{z} \in \mathbb{Z}^n$  with  $\|\mathbf{z}\|_\infty \leq O(n^2 M^2 \|\mathbf{c}\|_\infty)$  such that with constant probability,  $\min_{\mathbf{Ax} \geq \mathbf{b}} \mathbf{z}^T \mathbf{x}$  has an unique minimizer  $\mathbf{x}^*$  and this minimizer is one of the minimizer(s) of  $\min_{\mathbf{Ax} \geq \mathbf{b}} \mathbf{c}^T \mathbf{x}$ . Furthermore, if there is an interior point  $\mathbf{y}$  such that  $\mathbf{z}^T \mathbf{y} < \min_{\mathbf{Ax} \geq \mathbf{b}} \mathbf{z}^T \mathbf{x} + \delta$ , then  $\|\mathbf{y} - \mathbf{x}^*\|_\infty \leq 2nM\delta$ .*

*Proof.* The first part of the lemma follows by randomly perturbing the cost vector  $\mathbf{c}$ . We consider a new cost vector  $\mathbf{z} = 100n^2M^2\mathbf{c} + \mathbf{r}$  where each coordinate of  $\mathbf{r}$  is sampled randomly from  $\{0, 1, \dots, 10nM\}$ . [145, Lem 4] shows that the linear program  $\min_{\mathbf{Ax} \geq \mathbf{b}} \mathbf{z}^T \mathbf{x}$  has a unique minimizer with constant probability. Furthermore, it is clear that the minimizer of  $\min_{\mathbf{Ax} \geq \mathbf{b}} \mathbf{z}^T \mathbf{x}$  is a minimizer of  $\min_{\mathbf{Ax} \geq \mathbf{b}} \mathbf{c}^T \mathbf{x}$  (as  $r_i \ll 100n^2M^2|c_i|$ ).

Now we show the second part of the lemma. Given an interior point  $\mathbf{y}$  of the polytope  $\{\mathbf{Ax} \geq \mathbf{b}\}$ , we can write  $\mathbf{y}$  as a convex combination of the vertices of  $\{\mathbf{Ax} \geq \mathbf{b}\}$ , i.e.  $\mathbf{y} = \sum t_i \mathbf{v}_i$ . Note that

$\mathbf{z}^T \mathbf{y} = \sum t_i \mathbf{z}^T \mathbf{v}_i$ . If all  $\mathbf{v}_i$  are not the minimizer, then  $\mathbf{z}^T \mathbf{v}_i \geq \text{OPT} + 1$  and hence  $\mathbf{z}^T \mathbf{y} \geq \text{OPT} + 1$  which is impossible. Hence, we can assume that  $\mathbf{v}_1$  is the minimizer. Hence,  $\mathbf{z}^T \mathbf{v}_i = \text{OPT}$  if  $i = 1$  and  $\mathbf{z}^T \mathbf{v}_i \geq \text{OPT} + 1$  otherwise. We then have  $\mathbf{z}^T \mathbf{y} \geq \text{OPT} + (1 - t_1)$  which gives  $1 - t_1 < \delta$ . Finally, the claim follows from  $\|\mathbf{y} - \mathbf{v}_1\|_\infty \leq \sum_{i \neq 1} t_i \|\mathbf{v}_i - \mathbf{v}_1\|_\infty \leq 2nM\delta$ .  $\square$

### ■ 9.4.2 Matroid Intersection

Let  $M_1 = (E, \mathcal{I}_1)$  and  $M_2 = (E, \mathcal{I}_2)$  be two matroids sharing the same ground set. In this section we consider the weighted matroid intersection problem

$$\min_{S \in \mathcal{I}_1 \cap \mathcal{I}_2} \mathbf{w}(S).$$

where  $\mathbf{w} \in \mathbb{R}^E$  and  $w(S) \stackrel{\text{def}}{=} \sum_{e \in S} w_e$ .

For any matroid  $M = (E, \mathcal{I})$ , it is well known that the polytope of all independent sets has the following description [80]:

$$\text{conv}(\mathcal{I}_1) = \{\mathbf{x} \in \mathbb{R}^E \text{ s.t. } 0 \leq x(S) \leq r(S) \text{ for all } S \subseteq E\} \quad (9.20)$$

where  $r$  is the rank function for  $M$ , i.e.  $r(S)$  is the size of the largest independent set that is a subset of  $S$ . Furthermore, the polytope of the matroid intersection satisfies  $\text{conv}(\mathcal{I}_1 \cap \mathcal{I}_2) = \text{conv}(\mathcal{I}_1) \cap \text{conv}(\mathcal{I}_2)$ .

It is well known that the optimization problem

$$\min_{S \in \mathcal{I}_1} w(S) \text{ and } \min_{S \in \mathcal{I}_2} w(S)$$

can be solved efficiently by the greedy method. Given a matroid (polytope), the greedy method finds a maximum weight independent subset by maintaining a candidate independent subset  $S$  and iteratively attempts to add new element to  $S$  in descending weight. A element  $i$  is added to  $S$  if  $S \cup \{i\}$  is still independent. A proof of this algorithm is well-known and can be found in any standard textbook on combinatorial optimization.

Clearly, the greedy method can be implemented by  $O(n)$  calls to the independence oracle (also called membership oracle). For rank oracle, it requires  $O(r \log n)$  calls by finding the next element to add via binary search. Therefore, we can apply Theorem 9.4.4 to get the following result (note that this algorithm is the fastest if  $r$  is close to  $n$  for the independence oracle).

**Theorem 9.4.7.** *Suppose that the weights  $\mathbf{w}$  are integer with  $\|\mathbf{w}\|_\infty \leq M$ . Then, we can find*

$$S \in \arg \min_{S \in \mathcal{I}_1 \cap \mathcal{I}_2} w(S)$$

*in time  $O(n \text{GO} \log(nM) + n^3 \log^{O(1)}(nM))$  where GO is the cost of greedy method for  $\mathcal{I}_1$  and  $\mathcal{I}_2$ .*

*Proof.* Applying Lemma 9.4.6, we can find a new cost  $\mathbf{z}$  such that

$$\min_{\mathbf{x} \in \text{conv}(\mathcal{I}_1) \cap \text{conv}(\mathcal{I}_2)} \mathbf{z}^T \mathbf{x}$$

has a unique solution. Note that for any  $\mathbf{x} \in \text{conv}(\mathcal{I}_1)$ , we have  $\|\mathbf{x}\|_\infty \leq 1$ . Hence, applying theorem 9.4.4, we can find  $\mathbf{q}$  such that  $\mathbf{q}^T \mathbf{z} \leq \text{OPT} + \varepsilon$  and  $\|\mathbf{q} - \mathbf{x}\|_2 + \|\mathbf{q} - \mathbf{y}\|_2 \leq \varepsilon$  for some  $\mathbf{x} \in \text{conv}(\mathcal{I}_1)$  and  $\mathbf{y} \in \text{conv}(\mathcal{I}_2)$ . Using (9.20), we have the coordinate wise minimum of  $\mathbf{x}, \mathbf{y}$ , i.e.  $\min\{\mathbf{x}, \mathbf{y}\}$ , is in  $\text{conv}(\mathcal{I}_1) \cap \text{conv}(\mathcal{I}_2)$ . Since  $\|\mathbf{q} - \min\{\mathbf{x}, \mathbf{y}\}\|_2 \leq \|\mathbf{q} - \mathbf{x}\|_2 + \|\mathbf{q} - \mathbf{y}\|_2 \leq \varepsilon$ , we have

$$(\min\{\mathbf{x}, \mathbf{y}\})^T \mathbf{z} \leq \text{OPT} + nM\varepsilon.$$

Hence, we have a feasible point  $\min\{\mathbf{x}, \mathbf{y}\}$  which has value close to optimal and Lemma 9.4.6 shows that  $\|\min(\mathbf{x}, \mathbf{y}) - \mathbf{s}\|_\infty \leq 2n^2M^2\varepsilon$  where  $\mathbf{s}$  is the optimal solution. Hence, we have  $\|\mathbf{q} - \mathbf{s}\|_\infty \leq 2n^2M^2\varepsilon + \varepsilon$ . Picking  $\varepsilon = \frac{1}{6n^2M^2}$ , we have  $\|\mathbf{q} - \mathbf{s}\|_\infty < \frac{1}{2}$  and hence, we can get the optimal solution by rounding to the nearest integer.

Since optimization over  $\mathcal{I}_1$  and  $\mathcal{I}_2$  involves applying greedy method on certain vectors, it takes only  $O(\text{GO})$  time. Theorem 9.4.4 shows it only takes  $O\left(n\text{GO} \log(nM) + n^3 \log^{O(1)}(nM)\right)$  in finding such  $\mathbf{q}$ .  $\square$

This gives the following corollary.

**Corollary 9.4.8.** *We have  $O(n \min(n\mathcal{T}_{\text{ind}}, r\mathcal{T}_{\text{rank}}) \log(nM) + n^3 \log^{O(1)}(nM))$  time algorithm for weighted matroid intersection. Here  $\mathcal{T}_{\text{ind}}$  is the time needed to check if a subset is independent, and  $\mathcal{T}_{\text{rank}}$  is the time needed to compute the rank of a given subset.*

*Proof.* By Theorem 9.4.7, it suffices to show that the optimization oracle for the matroid polytope can be implemented in  $O(n\mathcal{T}_{\text{ind}})$  and  $O(r\mathcal{T}_{\text{rank}} \log n)$  time. This is simply attained by the well-known greedy algorithm which iterates through all the positively-weighted elements in decreasing order, and adds an element to our candidate independent set whenever possible.

For the independence oracle, this involves one oracle call for each element. On the other hand, for the rank oracle, we can find the next element to add by binary search which takes time  $O(\mathcal{T}_{\text{rank}} \log n)$ . Since there are at most  $r$  elements to add, we have the desired running time.  $\square$

### ■ 9.4.3 Submodular Flow

Let  $G = (V, E)$  be a directed graph with  $|E| = m$ , let  $f$  be a submodular function on  $\mathbb{R}^V$  with  $|V| = n$ ,  $f(\emptyset) = 0$  and  $f(V) = 0$ , and let  $A$  be the incidence matrix of  $G$ . In this section we consider the submodular flow problem

$$\begin{aligned} & \text{Minimize} && \langle \mathbf{c}, \varphi \rangle && (9.21) \\ & \text{subject to} && l(e) \leq \varphi(e) \leq u(e) \quad \forall e \in E \\ & && x(v) = (A\varphi)(v) \quad \forall v \in V \\ & && \sum_{v \in S} x(v) \leq f(S) \quad \forall S \subseteq V \end{aligned}$$

where  $\mathbf{c} \in \mathbb{Z}^E$ ,  $l \in \mathbb{Z}^E$ ,  $u \in \mathbb{Z}^E$  where  $C = \|\mathbf{c}\|_\infty$  and  $U = \max(\|u\|_\infty, \|l\|_\infty, \max_{S \subseteq V} |f(S)|)$ . Here  $\mathbf{c}$  is the cost on edges,  $\varphi$  is the flow on edges,  $l$  and  $u$  are lower and upper bounds on the amount of flow on the edges, and  $x(v)$  is the net flow out of vertex  $v$ . The submodular function  $f$  upper bounds the total net flow out of any subset  $S$  of vertices by  $f(S)$ .

**Theorem 9.4.9.** *Suppose that the cost vector  $\mathbf{c}$  is integer weight with  $\|\mathbf{c}\|_\infty \leq C$  and the capacity vector and the submodular function satisfy  $U = \max(\|u\|_\infty, \|l\|_\infty, \max_{S \subseteq V} |f(S)|)$ . Then, we can solve the submodular flow problem (9.21) in time  $O\left(n^2 \text{EO} \log(mCU) + n^3 \log^{O(1)}(mCU)\right)$  where EO is the cost of function evaluation.*

*Proof.* First, we can assume  $l(e) \leq u(e)$  for every edge  $e$ , otherwise, the problem is infeasible. Now, we apply a similar transformation in [110] to modify the graph. We create a new vertex  $v_0$ . For every vertex  $v$  in  $V$ , we create an edge from  $v_0$  to  $v$  with capacity lower bounded by 0, upper bounded by  $4nU$ , and with cost  $2mCU$ . Edmonds and Giles showed that the submodular flow polytope is integral [83]. Hence, there is an integral optimal flow on this new graph. If the optimal flow passes through the



newly created edge, then it has cost at least  $2mCU - mCU$  because the cost of all other edges in total has at least  $-mCU$ . That means the optimal flow has the cost larger than  $mCU$  which is impossible. So the optimal flow does not use the newly created edges and vertex and hence the optimal flow in the new problem gives the optimal solution of the original problem. Next, we note that for any  $\varphi$  on the original graph such that  $l(e) \leq \varphi(e) \leq u(e)$ , we can send suitable amount of flow from  $v_0$  to  $v$  to make  $\varphi$  feasible. Hence, this modification makes the feasibility problem trivial.

Lemma 9.4.6 shows that we can assume the new problem has an unique solution and it only blows up  $C$  by a  $(mU)^{O(1)}$  factors.

Note that the optimal value is an integer and its absolute value at most  $mCU$ . By binary search, we can assume we know the optimal value  $\text{OPT}$ . Now, we reduce the problem to finding a feasible  $\varphi$  with  $\{\langle d, \varphi \rangle \leq \text{OPT} + \varepsilon\}$  with  $\varepsilon$  determined later. Let  $P_\varepsilon$  be the set of such  $\varphi$ . Note that  $P_\varepsilon = K_{1,\varepsilon} \cap K_{2,\varepsilon}$  where

$$K_{1,\varepsilon} = \left\{ x \in \mathbb{R}^V \text{ such that } \begin{array}{l} l(e) \leq \varphi(e) \leq u(e) \quad \forall e \in E \\ x(v) = (A\varphi)(v) \quad \forall v \in V \quad \text{for some } \varphi \\ \langle d, \varphi \rangle \leq \text{OPT} + \varepsilon \end{array} \right\},$$

$$K_{2,\varepsilon} = \left\{ y \in \mathbb{R}^V \text{ such that } \begin{array}{l} \sum_{v \in S} y(v) \leq f(S) \quad \forall S \subseteq V, \sum_{v \in V} y(v) = f(V) \end{array} \right\}.$$

Note that the extra condition  $\sum_v y(v) = f(V)$  is valid because  $\sum_v y(v) = \sum_v (A\varphi)(v) = 0$  and  $f(V) = 0$ , and  $K_{1,\varepsilon}$  has radius bounded by  $O((mCU)^{O(1)})$  and  $K_{2,\varepsilon}$  has radius bounded by  $O(nU)$ . Furthermore, for any vector  $c \in \mathbb{R}^V$ , we note that

$$\begin{aligned} \max_{x \in K_{1,\varepsilon}} \langle c, x \rangle &= \max_{l \leq \varphi \leq u, \langle d, \varphi \rangle \leq \text{OPT} + \varepsilon, x = A\varphi} \langle c, x \rangle \\ &= \max_{l \leq \varphi \leq u, \langle d, \varphi \rangle \leq \text{OPT} + \varepsilon} \langle c, A\varphi \rangle \\ &= \max_{l \leq \varphi \leq u, \langle d, \varphi \rangle \leq \text{OPT} + \varepsilon} \langle A^T c, \varphi \rangle. \end{aligned}$$

To solve this problem, again we can do a binary search on  $\langle d, \varphi \rangle$  and reduce the problem to

$$\max_{l \leq \varphi \leq u, \langle d, \varphi \rangle = K} \langle A^T c, \varphi \rangle$$

for some value of  $K$ . Since  $A^T c$  is fixed, this is a linear program with only the box constraints and an extra equality constraint. Hence, it can be solved in nearly linear time [164, Thm 17, ArXiv v1]. As the optimization oracle for  $K_{1,\varepsilon}$  involves only computing  $A^T c$  and solving this simple linear program, it takes only  $O(n^2 \log^{O(1)}(mCU/\varepsilon))$  time. On the other hand, since  $K_{2,\varepsilon}$  is just a base polyhedron, the optimization oracle for  $K_{2,\varepsilon}$  can be done by greedy method and only takes  $O(n\text{EO})$  time.

Applying Theorem 9.4.4, we can find  $q$  such that  $\|q - x\|_2 + \|q - y\|_2 \leq \delta$  for some  $x \in K_{1,\varepsilon}$ ,  $y \in K_{2,\varepsilon}$  and  $\delta$  to be chosen later. According to the definition of  $K_{1,\varepsilon}$ , there is  $\varphi$  such that  $l(e) \leq \varphi(e) \leq u(e)$  and  $x(v) = (A\varphi)(v)$  for all  $v$  and  $\langle d, \varphi \rangle \leq \text{OPT} + \varepsilon$ . Since  $\|y - x\|_2 \leq 2\delta$ , that means  $|y(v) - (A\varphi)(v)| \leq 2\delta$  for all  $v$ .

- Case 1) If  $y(v) \geq (A\varphi)(v)$ , then we can replace  $y(v)$  by  $(A\varphi)(v)$ , note that  $y$  is still in  $K_{2,\varepsilon}$  because of the submodular constraints.
- Case 2) If  $y(v) \leq (A\varphi)(v)$ , then we can send a suitable amount of flow from  $v_0$  to  $v$  to make  $\varphi$  feasible  $y(v) \leq (A\varphi)(v)$ .

Note that under this modification, we increased the objective value by  $(\delta n)(2mCU)$  because the new edge cost  $2mCU$  per unit of flow. Hence, we find a flow  $\varphi$  which is feasible in new graph with objective

value  $\varepsilon + (\delta n)(2mCU)$  far from optimum value. By picking  $\delta = \frac{1}{2mnCU}$ , we have the value  $2\varepsilon$  far from OPT. Now, we use Lemma 9.4.6 to shows that when  $\varepsilon$  is small enough, i.e,  $\frac{1}{(mCU)^c}$  for some constant  $c$ , then we can guarantee that  $\|y - x^*\|_\infty \leq \frac{1}{4}$  where  $x^*$  is the optimal demand. Now, we note that  $\|q - y\|_2 \leq \delta$  and we note that we only modify  $y$  by a small amount, we in fact have  $\|q - x^*\|_\infty < \frac{1}{2}$ . Hence, we can read off the solution  $x^*$  by rounding  $q$  to the nearest integer. Note that we only need to solve the problem  $K_{1,\varepsilon} \cap K_{2,\varepsilon}$  to  $\frac{1}{(mCU)^{\Theta(1)}}$  accuracy and the optimization oracle for  $K_{1,\varepsilon}$  and  $K_{2,\varepsilon}$  takes time  $O(n^2 \log^{O(1)}(mCU))$  and  $O(nEO)$  respectively. Hence, Theorem 9.4.4 shows that it takes  $O\left(n^2EO \log(mCU) + n^3 \log^{O(1)}(mCU)\right)$  time to find  $x^*$  exactly.

After getting  $x^*$ , one can find  $\varphi^*$  by solving a min cost flow problem using interior point method (Chapter 6), which takes  $O(m\sqrt{n} \log^{O(1)}(mCU))$  time.  $\square$

### ■ 9.4.4 Affine Subspace of Convex Set

In this section, we give another example about using optimization oracle directly via regularization. We consider the following optimization problem

$$\max_{\mathbf{x} \in K \text{ and } \mathbf{Ax}=\mathbf{b}} \langle \mathbf{c}, \mathbf{x} \rangle \quad (9.22)$$

where  $\mathbf{x}, \mathbf{c} \in \mathbb{R}^n$ ,  $K$  is a convex subset of  $\mathbb{R}^n$ ,  $\mathbf{A} \in \mathbb{R}^{r \times n}$  and  $\mathbf{b} \in \mathbb{R}^m$ . We suppose that  $r \ll n$  and thus, the goal of this subsection is to show how to obtain an algorithm takes only  $\tilde{O}(r)$  many iterations. To do this, we assume a slightly stronger optimization oracle for  $K$ :

**Definition 9.4.10.** Given a convex set  $K$  and  $\delta > 0$ . A  $\delta$ -2nd-order-optimization oracle for  $K$  is a function on  $\mathbb{R}^n$  such that for any input  $\mathbf{c} \in \mathbb{R}^n$  and  $\lambda > 0$ , it outputs  $\mathbf{y}$  such that

$$\max_{\mathbf{x} \in K} (\langle \mathbf{c}, \mathbf{x} \rangle - \lambda \|\mathbf{x}\|^2) \leq \delta + \langle \mathbf{c}, \mathbf{y} \rangle - \lambda \|\mathbf{y}\|^2.$$

We denote by  $\text{OO}_{\delta,\lambda}^{(2)}(K)$  the time complexity of this oracle.

The strategy for solving this problem is very similar to the intersection problem and hence some details are omitted.

**Theorem 9.4.11.** Assume that  $\max_{\mathbf{x} \in K} \|\mathbf{x}\|_2 < M$ ,  $\|\mathbf{b}\|_2 < M$ ,  $\|\mathbf{c}\|_2 < M$ ,  $\|\mathbf{A}\|_2 < M$  and  $\lambda_{\min}(\mathbf{A}) > 1/M$ . Assume that  $K \cap \{\mathbf{Ax} = \mathbf{b}\} \neq \emptyset$  and we have  $\varepsilon$ -2nd-order-optimization oracle for every  $\varepsilon > 0$ . For  $0 < \delta < 1$ , we can find  $\mathbf{z} \in K$  such that

$$\max_{\mathbf{x} \in K \text{ and } \mathbf{Ax}=\mathbf{b}} \langle \mathbf{c}, \mathbf{x} \rangle \leq \delta + \langle \mathbf{c}, \mathbf{z} \rangle$$

and  $\|\mathbf{Az} - \mathbf{b}\|_2 \leq \delta$ . This algorithm takes time

$$O\left(r \text{OO}_{\eta,\lambda}^{(2)}(K) \log\left(\frac{nM}{\delta}\right) + r^3 \log^{O(1)}\left(\frac{nM}{\delta}\right)\right)$$

where  $r$  is the number of rows in  $\mathbf{A}$ ,  $\eta = \left(\frac{\delta}{nM}\right)^{\Theta(1)}$  and  $\lambda = \left(\frac{\delta}{nM}\right)^{\Theta(1)}$ .

*Proof.* The proof is based on the minimax problem

$$\text{OPT}_\lambda \stackrel{\text{def}}{=} \min_{\|\boldsymbol{\eta}\|_2 \leq 1/\lambda} \max_{\mathbf{x} \in K} \langle \mathbf{c}, \mathbf{x} \rangle + \langle \boldsymbol{\eta}, \mathbf{Ax} - \mathbf{b} \rangle - \frac{1}{\lambda} \|\mathbf{x}\|_2^2$$

where  $\lambda = \left(\frac{\delta}{nM}\right)^c$  for some large constant  $c$ . We note that

$$\begin{aligned} \text{OPT}_\lambda &= \max_{\mathbf{x} \in K} \min_{\|\boldsymbol{\eta}\|_2 \leq 1/\lambda} \langle \mathbf{c}, \mathbf{x} \rangle + \langle \boldsymbol{\eta}, \mathbf{A}\mathbf{x} - \mathbf{b} \rangle - \frac{1}{\lambda} \|\mathbf{x}\|_2^2 \\ &= \max_{\mathbf{x} \in K} \langle \mathbf{c}, \mathbf{x} \rangle - \frac{1}{\lambda} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2 - \frac{1}{\lambda} \|\mathbf{x}\|_2^2. \end{aligned}$$

Since  $\lambda_{\min}(\mathbf{A}) > 1/M$  and the set  $K$  is bounded by  $M$ , one can show that the saddle point  $(\mathbf{x}^*, \boldsymbol{\eta}^*)$  of the minimax problem gives a good enough solution  $\mathbf{x}$  for the original problem for large enough constant  $c$ .

For any  $\boldsymbol{\eta}$ , we define

$$\mathbf{x}_\boldsymbol{\eta} = \arg \max_{\mathbf{x} \in K} \langle \mathbf{c}, \mathbf{x} \rangle + \langle \boldsymbol{\eta}, \mathbf{A}\mathbf{x} - \mathbf{b} \rangle - \frac{1}{\lambda} \|\mathbf{x}\|_2^2.$$

Since the problem is strongly concave in  $\mathbf{x}$ , one can prove that

$$\|\mathbf{x}_\boldsymbol{\eta} - \mathbf{x}^*\|_2 \leq \left(\frac{nM}{\delta}\right)^{O(1)} \|\boldsymbol{\eta} - \boldsymbol{\eta}^*\|_2.$$

Hence, we can first find an approximate minimizer of the function  $f(\boldsymbol{\eta}) = \max_{\mathbf{x} \in K} \langle \mathbf{c}, \mathbf{x} \rangle + \langle \boldsymbol{\eta}, \mathbf{A}\mathbf{x} - \mathbf{b} \rangle - \frac{1}{\lambda} \|\mathbf{x}\|_2^2$  and use the oracle to find  $\mathbf{x}_\boldsymbol{\eta}$ .

To find an approximate minimizer of  $f$ , we note that the subgradient of  $f$  can be found using the optimization oracle similar to Theorem 9.4.4. Hence, the result follows from our cutting plane method and the fact that  $\boldsymbol{\eta} \in \mathbb{R}^r$ .  $\square$

*Remark 9.4.12.* In Chapter 6, we considered the special case  $K = \{\mathbf{x} : 0 \leq x_i \leq 1\}$  and showed that it can be solved in  $\tilde{O}(\sqrt{r})$  iterations using interior point methods. This gives the current fastest algorithm for the maximum flow problem on directed weighted graphs. This result generalizes that result to any convex set  $K$  but with  $\tilde{O}(r)$  iterations. This suggests the following open problem: under what condition on  $K$  can one optimize linear functions over affine subspaces of  $K$  with  $r$  constraints in  $\tilde{O}(\sqrt{r})$  iterations?



---

---

## Submodular Minimization In Nearly-Cubic Time

### ■ 10.1 Introduction

Submodular functions and submodular function minimization (SFM) are fundamental to the field of combinatorial optimization. Examples of submodular functions include graph cut functions, set coverage function, and utility functions from economics. Since the seminal work by Edmonds in 1970 [82], submodular functions and the problem of minimizing such functions (i.e. submodular function minimization) have served as a popular modeling and optimization tool in various fields such as theoretical computer science, operations research, game theory, and most recently, machine learning. Given its prevalence, fast algorithms for SFM are of immense interest both in theory and in practice.

Throughout this chapter, we consider the standard formulation of SFM: we are given a submodular function  $f$  defined over the subsets of a  $n$ -element ground set. The values of  $f$  are integers, have absolute value at most  $M$ , and are evaluated by querying an oracle that takes time  $\text{EO}$ . Our goal is to produce an algorithm that solves this SFM problem, i.e. finds a minimizer of  $f$ , while minimizing both the number of oracle calls made and the total running time.

We provide new  $O(n^2 \log nM \cdot \text{EO} + n^3 \log^{O(1)} nM)$  and  $O(n^3 \log^2 n \cdot \text{EO} + n^4 \log^{O(1)} n)$  time algorithms for SFM. These algorithms improve upon the previous fastest weakly and strongly polynomial time algorithms for SFM which had a running time of  $O((n^4 \cdot \text{EO} + n^5) \log M)$  [118] and  $O(n^5 \cdot \text{EO} + n^6)$  [216] respectively. Consequently, we improve the running times in both regimes by roughly a factor of  $O(n^2)$ .

Both of our algorithms bear resemblance to the classic approach of Grötschel, Lovász and Schrijver [110, 109] using the Lovász extension. In fact our weakly polynomial time algorithm directly uses the Lovász extension as well as the results of Chapter 9 to achieve these results. Our strongly polynomial time algorithm also uses the Lovász extension, along with more modern tools from the past 15 years.

At a high level, our strongly polynomial algorithms apply our cutting plane method in conjunction with techniques originally developed by Iwata, Fleischer, and Fujishige (IFF) [121]. Our cutting plane method is performed for enough iterations to sandwich the feasible region in a narrow strip from which useful structural information about the minimizers can be deduced. Our ability to derive the new information hinges on a significant extension of IFF techniques.

Over the past few decades, SFM has drawn considerable attention from various research communities, most recently in machine learning [24, 149]. Given this abundant interest in SFM, we hope that our ideas will be of value in various practical applications. Indeed, one of the critiques against existing theoretical algorithms is that their running time is too slow to be practical. Our contribution, on the contrary, shows that this school of algorithms can actually be made fast theoretically and we hope it may potentially be competitive against heuristics which are more commonly used.

### ■ 10.1.1 Previous Work

Here we provide a brief survey of the history of algorithms for SFM. For a more comprehensive account of the rich history of SFM, we refer the readers to recent surveys [189, 120].

The first weakly and strongly polynomial time algorithms for SFM were based on the ellipsoid method [142] and were established in the foundational work of Grötschel, Lovász and Schrijver in 1980's [110, 109]. Their work was complemented by a landmark paper by Cunningham in 1985 which provided a pseudopolynomial algorithm that followed a flow-style algorithmic framework [61]. His tools foreshadowed much of the development in SFM that would take place 15 years later. Indeed, modern algorithms synthesize his framework with inspirations from various max flow algorithms.

The first such “flow style” strongly polynomial algorithms for SFM were discovered independently in the breakthrough papers by Schrijver [231] and Iwata, Fleischer, and Fujishige (IFF) [121]. Schrijver's algorithm has a running of  $O(n^8 \cdot \text{EO} + n^9)$  and borrows ideas from the push-relabel algorithms [105, 68] for the maximum flow problem. On the other hand, IFF's algorithm runs in time  $O(n^7 \log n \cdot \text{EO})$  and  $O(n^5 \cdot \text{EO} \log M)$ , and applies a flow-scaling scheme with the aid of certain proximity-type lemmas as in the work of Tardos [247]. Their method has roots in flow algorithms such as [117, 106].

Subsequent work on SFM provided algorithms with considerably faster running time by extending the ideas in these two “genesis” papers [231, 121] in various novel directions [264, 87, 118, 216, 125]. Currently, the fastest weakly and strongly polynomial time algorithms for SFM have a running time of  $O((n^4 \cdot \text{EO} + n^5) \log M)$  [118] and  $O(n^5 \cdot \text{EO} + n^6)$  [216] respectively. Despite this impressive track record, the running time has not been improved in the last eight years.

We remark that all of the previous algorithms for SFM proceed by maintaining a convex combination of  $O(n)$  BFS's of the base polyhedron, and incrementally improving it in a relatively local manner. As we shall discuss in Section 10.1.2, our algorithms do not explicitly maintain a convex combination. This may be one of the fundamental reasons why our algorithms achieve a faster running time.

Finally, beyond the distinction between weakly and strongly polynomial time algorithms for SFM, there has been interest in another type of SFM algorithm, known as fully combinatorial algorithms in which only additions and subtractions are permitted. Previous such algorithms include [125, 118, 119]. We do not consider such algorithms in the remainder of this chapter and leave it as an open question if it is possible to turn our algorithms into fully combinatorial ones.

### ■ 10.1.2 Our Results and Techniques

In this chapter, we show how to improve upon the previous best known running times for SFM by a factor of  $O(n^2)$  in both the strongly and weakly polynomial regimes. In Table 10.1 summarizes the running time of the previous algorithms as well as the running times of the fastest algorithms presented in this chapter.

Both our weakly and strongly polynomial algorithms for SFM utilize a convex relaxation of the submodular function, called the *Lovász extension*. Our algorithms apply our cutting plane method from Chapter 8 using a separation oracle given by the subgradient of the Lovász extension. To the best of the author's knowledge, Grötschel, Lovász and Schrijver were the first to formulate this convex optimization framework for SFM [110, 109].

For weakly polynomial algorithms, our contribution is two-fold. First, we show that cutting plane methods such as Vaidya's [253] can be applied to SFM to yield faster algorithms. Second, as our cutting plane method, Theorem 9.3.2, improves upon previous cutting plane algorithms and consequently the running time for SFM as well. This gives a running time of  $O(n^2 \log nM \cdot \text{EO} + n^3 \log^{O(1)} nM)$ , an improvement over the previous best algorithm by Iwata [118] by a factor of almost  $O(n^2)$ .

Our strongly polynomial algorithms, on the other hand, require substantially more innovation. We

Years	Authors	Running times	Remarks
1981,1988	Grötschel, Lovász,	$\tilde{O}(n^5 \cdot \text{EO} + n^7)$ [189]	first weakly and strongly
1985	Cunningham [61]	$O(Mn^6 \log nM \cdot \text{EO})$	first pseudopoly
2000	Schrijver [231]	$O(n^8 \cdot \text{EO} + n^9)$	first combin. strongly
2000	Iwata, Fleischer,	$O(n^5 \cdot \text{EO} \log M)$ $O(n^7 \log n \cdot \text{EO})$	first combin. strongly
2000	Iwata, Fleischer [87]	$O(n^7 \cdot \text{EO} + n^8)$	
2003	Iwata [118]	$O((n^4 \cdot \text{EO} + n^5) \log M)$ $O((n^6 \cdot \text{EO} + n^7) \log n)$	current best weakly
2003	Vygen [264]	$O(n^7 \cdot \text{EO} + n^8)$	
2007	Orlin [216]	$O(n^5 \cdot \text{EO} + n^6)$	current best strongly
2009	Iwata, Orlin [125]	$O((n^4 \cdot \text{EO} + n^5) \log nM)$ $O((n^5 \cdot \text{EO} + n^6) \log n)$	
-	This Chapter	$O(n^2 \log nM \cdot \text{EO} + n^3 \log^{O(1)} nM)$ $O(n^3 \log^2 n \cdot \text{EO} + n^4 \log^{O(1)} n)$	

**Table 10.1:** Algorithms for submodular function minimization. Note that some of these algorithms were published in both conferences and journals, in which case the year we provided is the earlier one.

first begin with a very simple geometric argument that SFM can be solved in  $O(n^3 \log n \cdot \text{EO})$  oracle calls (but in exponential time). This proof only uses Grunbaum’s Theorem from convex geometry and is completely independent from the rest of this chapter. It was the starting point of our method and suggests that a running time of  $\tilde{O}(n^3 \cdot \text{EO} + n^{O(1)})$  for submodular minimization is in principle achievable.

To make this existence result algorithmic, we first run cutting plane, Theorem 8.3.26, for enough iterations such that we compute either a minimizer or a set  $P$  containing the minimizers that fits within in a narrow strip. This narrow strip consists of the intersection of two approximately parallel hyperplanes. If our narrow strip separates  $P$  from one of the faces  $x_i = 0$ ,  $x_i = 1$ , we can effectively eliminate the element  $i$  from our consideration and reduce the dimension of our problem by 1. Otherwise a pair of elements  $p, q$  can be identified for which  $q$  is guaranteed to be in any minimizer containing  $p$  (but  $p$  may not be contained in a minimizer). Our first algorithm deduces only one such pair at a time. This technique immediately suffices to achieve a  $\tilde{O}(n^4 \cdot \text{EO} + n^5)$  time algorithm for SFM (See Section 10.4.3). We then improve the running time to  $\tilde{O}(n^3 \cdot \text{EO} + n^4)$  by showing how to deduce many such pairs simultaneously. Similar to past algorithms, this structural information is deduced from a point in the so-called base polyhedron (See Section 10.2).

Readers well-versed in SFM literature may recognize that our strongly polynomial algorithms are reminiscent of the scaling-based approach first used by IFF [121] and later in [118, 125]. While both approaches share the same skeleton, there are differences as to how structural information about minimizers is deduced. A comparison of our algorithms and previous ones are presented in Section 10.5.

Finally, there is one more crucial difference between these algorithms which we believe is responsible for much of our speedup. One common feature shared by all the previous algorithms is that they maintain a convex combination of  $O(n)$  BFS’s of the base polyhedron, and incrementally improve on it by introducing new BFS’s by making local changes to existing ones. Our algorithms, on the other hand, choose new BFS’s by the cutting plane method. Because of this, our algorithm considers the geometry of the existing BFS’s where each of them has influences over the choice of the next BFS. In some sense, our next BFS is chosen in a more “global” manner.

### ■ 10.1.3 Organization

The rest of this chapter is organized as follows. We first begin with a gentle introduction to submodular functions in Section 10.2. In Section 10.3, we apply our cutting plane method to SFM to obtain a faster weakly polynomial algorithms. In Section 10.4 we then present our results for achieving better strongly polynomial algorithms, where a warm-up  $\tilde{O}(n^4 \cdot \text{EO} + n^5)$  algorithm is given before the full-fledged  $\tilde{O}(n^3 \cdot \text{EO} + n^4)$  algorithm. Finally, we end this chapter with a discussion and comparison between our algorithms and previous ones in Section 10.5.

We note that there are a few results in this chapter that can be read fairly independently of the rest of this chapter. In Theorem 10.3.2 we show how Vaidya's algorithm can be applied to SFM to obtain a faster weakly polynomial running time. Also in Theorem 10.4.1 we present a simple geometric argument that SFM can be solved with  $O(n^3 \log n \cdot \text{EO})$  oracle calls but with exponential time. These results can be read with only a working knowledge of the Lovász extension of submodular functions.

## ■ 10.2 Preliminaries

Here we introduce background on submodular function minimization (SFM) and notation that we use throughout this chapter. Our exposition is kept to a minimal amount sufficient for our purposes. We refer interested readers to the extensive survey by McCormick [189] for further intuition.

### ■ 10.2.1 Submodular Function Minimization

Throughout the rest of this chapter, let  $V = \{1, \dots, n\} = [n]$  denote a ground set and let  $f : 2^V \rightarrow \mathbb{Z}$  denote a *submodular function* defined on subsets of this ground set. We use  $V$  and  $[n]$  interchangeably and let  $[0] \stackrel{\text{def}}{=} \emptyset$ . We abuse notation by letting  $S + i \stackrel{\text{def}}{=} S \cup \{i\}$  and  $S - i \stackrel{\text{def}}{=} S \setminus \{i\}$  for an element  $i \in V$  and a set  $S \subseteq 2^V$ . Formally, we call a function submodular if it obeys the following property of *diminishing marginal differences*:

**Definition 10.2.1** (Submodularity). A function  $f : 2^V \rightarrow \mathbb{Z}$  is *submodular* if  $f(T + i) - f(T) \leq f(S + i) - f(S)$  for any  $S \subseteq T$  and  $i \in V \setminus T$ .

For convenience we assume without loss of generality that  $f(\emptyset) = 0$  by replacing  $f(S)$  by  $f(S) - f(\emptyset)$  for all  $S$ . We also let  $M \stackrel{\text{def}}{=} \max_{S \in 2^V} |f(S)|$ .

The central goal of this chapter is to design algorithms for SFM, i.e. computing the minimizer of  $f$ . We call such an algorithm *strongly polynomial* if its running time depends only polynomially on  $n$  and  $\text{EO}$ , the time needed to compute  $f(S)$  for a set  $S$ , and we call such an algorithm *weakly polynomial* if it also depends polylogarithmically on  $M$ .

### ■ 10.2.2 Lovász Extension

Our new algorithms for SFM all consider a convex relaxation of a submodular function, known as the Lovász extension, and then carefully apply our cutting plane methods to it. Here we formally introduce the Lovász extension and present basic facts that we use throughout this chapter.

The Lovász extension of  $\hat{f} : [0, 1]^n \rightarrow \mathbb{R}$  of our submodular function  $f$  is defined for all  $\mathbf{x}$  by

$$\hat{f}(\mathbf{x}) \stackrel{\text{def}}{=} \mathbb{E}_{t \sim [0, 1]} [f(\{i : x_i \geq t\})],$$

where  $t \sim [0, 1]$  is drawn uniformly at random from  $[0, 1]$ . The Lovász extension allows us to reduce SFM to minimizing a convex function defined over the interior of the hypercube. Below we state that the Lovász extension is a convex relaxation of  $f$  and that it can be evaluated efficiently.



**Theorem 10.2.2.** *The Lovász extension  $\hat{f}$  satisfies the following properties:*

1.  $\hat{f}$  is convex and  $\min_{\mathbf{x} \in [0,1]^n} \hat{f}(\mathbf{x}) = \min_{S \subseteq [n]} f(S)$ ;
2.  $f(S) = \hat{f}(I_S)$ , where  $I_S$  is the characteristic vector for  $S$ , i.e.  $I_S(i) = \begin{cases} 1 & \text{if } i \in S \\ 0 & \text{if } i \notin S \end{cases}$ ;
3. If  $S$  is a minimizer of  $f$ , then  $I_S$  is a minimizer of  $\hat{f}$ ;
4. Suppose  $x_1 \geq \dots \geq x_n \geq x_{n+1} \stackrel{\text{def}}{=} 0$ , then

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^n f([i])(x_i - x_{i+1}) = \sum_{i=1}^n (f([i]) - f([i-1]))x_i.$$

*Proof.* See [109] or any standard textbook on combinatorial optimization, e.g. [232].  $\square$

Next we show that we can efficiently compute a subgradient of the Lovász or alternatively, a separating hyperplane for the set of minimizers of our submodular function  $f$ . First we remind the reader of the definition of a separation oracle, and then we prove the necessary properties of the hyperplane, Theorem 10.2.4.

**Definition 10.2.3** (separation oracle, Definition 2.3.5 restated for Lovász extension). Given a point  $\bar{x}$  and a convex function  $\hat{f}$  over a convex set  $P$ ,  $\mathbf{a}^T \mathbf{x} \leq b$  is a separating hyperplane if  $\mathbf{a}^T \bar{x} \geq b$  and any minimizer  $x^*$  of  $\hat{f}$  over  $P$  satisfies  $\mathbf{a}^T x^* \leq b$ .

**Theorem 10.2.4.** *Given a point  $\bar{x} \in [0,1]^n$  assume without loss of generality (by re-indexing the coordinates) that  $\bar{x}_1 \geq \dots \geq \bar{x}_n$ . Then the following inequality is a valid separating hyperplane for  $\mathbf{x}$  and  $f$ :*

$$\sum_{i=1}^n (f([i]) - f([i-1]))x_i \leq \hat{f}(\bar{x})$$

*i.e., it satisfies the following:*

1. (separating)  $\bar{x}$  lies on  $\sum_{i=1}^n (f([i]) - f([i-1]))x_i \leq \hat{f}(\bar{x})$ .
2. (valid) For any  $\mathbf{x}$ , we have  $\sum_{i=1}^n (f([i]) - f([i-1]))x_i \leq \hat{f}(\mathbf{x})$ . In particular,  $\sum_{i=1}^n (f([i]) - f([i-1]))x_i^* \leq \hat{f}(\bar{x})$  for any minimizer  $\mathbf{x}^*$ , i.e. the separating hyperplane does not cut out any minimizer.

Moreover, such a hyperplane can be computed with  $n$  oracle calls to  $f$  and in time  $O(n \cdot EO + n^2)$ .

*Proof.* Note that by Theorem 10.2.2 we have that  $\sum_{i \in [n]} (f([i]) - f([i-1]))x_i = \hat{f}(\bar{x})$  and thus the hyperplane satisfies the separating condition. Moreover, clearly computing it only takes time  $O(n \cdot EO + n^2)$  as we simply need to sort the coordinates and evaluate  $f$  at  $n$  points, i.e. each of the  $[i]$ . All that remains is to show that the hyperplane satisfies the valid condition.

Let  $L^{(t)} \stackrel{\text{def}}{=} \{i : x_i \geq t\}$ . Recall that  $\hat{f}(\mathbf{x}) = \mathbb{E}_{t \sim [0,1]} [f(L_t)]$ . Thus  $\hat{f}(\mathbf{x})$  can be written as a convex combination  $\hat{f}(\mathbf{x}) = \sum_t \alpha_t f(L^{(t)})$ , where  $\alpha_t \geq 0$  and  $\sum_t \alpha_t = 1$ . However, by diminishing marginal

differences we see that for all  $t$

$$\begin{aligned} \sum_{i \in [n]} (f([i]) - f([i-1])) (I_{L^{(t)}})_i &= \sum_{i \in L^{(t)}} (f([i]) - f([i-1])) \\ &\leq \sum_{i \in L^{(t)}} (f([i] \cap L^{(t)}) - f([i-1] \cap L^{(t)})) \\ &= f(L^{(t)}) - f(\emptyset) = f(L^{(t)}) \end{aligned}$$

and therefore since  $\sum_t \alpha_t I_{L^{(t)}} = \mathbf{x}$  we have

$$\sum_{i \in [n]} (f[i] - f([i-1])) x_i = \sum_t \alpha_t \sum_{i=1}^n (f([i]) - f([i-1])) (I_{L^{(t)}})_i \leq \sum_t \alpha_t f(L^{(t)}) = \hat{f}(\mathbf{x}).$$

□

### ■ 10.2.3 Polyhedral Aspects of SFM

Here we provide a natural primal dual view of SFM that we use throughout the analysis. We provide a dual convex optimization program to minimizing the Lovász extension and provide several properties of these programs. We believe the material in this section helps crystallize some of the intuition behind our algorithm and we make heavy use of the notation presented in this section. However, we will not need to appeal to the strong duality of these programs in our proofs.

Consider the following primal and dual programs, where we use the shorthands  $y(S) = \sum_{i \in S} y_i$  and  $y_i^- = \min\{0, y_i\}$ . Here the primal constraints are often called the *base polyhedron*  $\mathcal{B}(f) \stackrel{\text{def}}{=} \{\mathbf{y} \in \mathbb{R}^n : y(S) \leq f(S) \forall S \subsetneq V, y(V) = f(V)\}$  and the dual program directly corresponds to minimizing the Lovász extension and thus  $f$ .

Primal	Dual
$\max y^-(V)$ $y(S) \leq f(S) \forall S \subsetneq V$ $y(V) = f(V)$	$\min \hat{f}(\mathbf{x})$ $0 \leq \mathbf{x} \leq 1$

**Theorem 10.2.5.**  $\mathbf{h}$  is a basic feasible solution (BFS) of the base polyhedron  $\mathcal{B}(f)$  if and only if

$$h_i = f(\{v_1, \dots, v_i\}) - f(\{v_1, \dots, v_{i-1}\})$$

for some permutation  $v_1, \dots, v_n$  of the ground set  $V$ . We call  $v_1, \dots, v_n$  the defining permutation of  $\mathbf{h}$ . We call  $v_i$  precedes  $v_j$  for  $i < j$ .

This theorem gives a nice characterization of the BFS's of  $\mathcal{B}(f)$ . It also gives the key observation underpinning our approach: **the coefficients of each separating hyperplane in Theorem 10.2.4 precisely corresponds to a primal BFS (Theorem 10.2.5)**. Our analysis relies heavily on this connection. We re-state Theorem 10.2.4 in the language of BFS.

**Lemma 10.2.6.** We have  $\mathbf{h}^T \mathbf{x} \leq \hat{f}(\mathbf{x})$  for any  $\mathbf{x} \in [0, 1]^n$  and BFS  $\mathbf{h}$ .

*Proof.* Any BFS is given by some permutation. Thus this is just Theorem 10.2.4 in disguise. □

We also note that since the objective function of the primal program is non-linear, we cannot say that the optimal solution to the primal program is a BFS. Instead we only know that it is a convex combination of the BFS's that satisfy the following property. A proof can be found in any standard textbook on combinatorial optimization.

**Theorem 10.2.7.** *The above primal and dual programs have no duality gap. Moreover, there always exists a primal optimal solution  $\mathbf{y} = \sum_k \lambda^{(k)} \mathbf{h}^{(k)}$  with  $\sum_k \lambda^{(k)} = 1$  (a convex combination of BFS  $\mathbf{h}^{(k)}$ ) s.t. any  $i$  with  $y_i < 0$  precedes any  $j$  with  $y_j > 0$  in the defining permutation for each BFS  $\mathbf{h}^{(k)}$ .*

Our algorithms will maintain collections of BFS and use properties of  $\mathbf{h} \in \mathcal{B}(f)$ , i.e. convex combination of BFS. To simplify our analysis at several points we will want to assume that such a vector  $\mathbf{h} \in \mathcal{B}(f)$  is *non-degenerate*, meaning it has both positive and negative entries. Below, we prove that such degenerate points in the base polytope immediately allow us to trivially solve the SFM problem.

**Lemma 10.2.8** (Degenerate Hyperplanes). *If  $\mathbf{h} \in \mathcal{B}(f)$  is non-negative then  $\emptyset$  is a minimizer of  $f$  and if  $\mathbf{h}$  is non-positive then  $V$  is a minimizer of  $f$ .*

*Proof.* While this follows immediately from Theorem 10.2.7, for completeness we prove this directly. Let  $S \in 2^V$  be arbitrary. If  $\mathbf{h} \in \mathcal{B}(f)$  is non-negative then by the we have

$$f(S) \geq \mathbf{h}(S) = \sum_{i \in S} h_i \geq 0 = f(\emptyset).$$

On the other hand if  $\mathbf{h}$  is non-positive then by definition we have

$$f(S) \geq \mathbf{h}(S) = \sum_{i \in S} h_i \geq \sum_{i \in V} h_i = \mathbf{h}(V) = f(V).$$

□

## ■ 10.3 Improved Weakly Polynomial Algorithms for SFM

In this section we show how our cutting plane method can be used to obtain a  $O(n^2 \log nM \cdot \text{EO} + n^3 \log^{O(1)} nM)$  time algorithm for SFM. Our main result in this section is the following theorem, which shows how directly applying our results from earlier chapters to minimize the Lovász extension yields the desired running time.

**Theorem 10.3.1.** *We have an  $O(n^2 \log nM \cdot \text{EO} + n^3 \log^{O(1)} nM)$  time algorithm for submodular function minimization.*

*Proof.* We apply Theorem 9.3.2 to the Lovász extension  $\hat{f} : [0, 1]^n \rightarrow \mathbb{R}$  with the separation oracle given by Theorem 10.2.4.  $\hat{f}$  fulfills the requirement on the domain as its domain  $\Omega = [0, 1]^n$  is symmetric about the point  $(1/2, \dots, 1/2)$  and has exactly  $2n$  constraints.

In the language of Theorem 9.3.2, our separation oracle is a  $(0, 0)$ -separation oracle with  $\eta = 0$  and  $\delta = 0$ .

We first show that  $\delta = 0$ . Firstly, our separating hyperplane can be written as

$$\sum_{i=1}^n (f([i]) - f([i-1]))x_i \leq \hat{f}(\bar{\mathbf{x}}) = \sum_{i=1}^n (f([i]) - f([i-1]))\bar{x}_i,$$

where the equality follows from Theorem 10.2.2. Secondly, for any  $\mathbf{x}$  with  $\hat{f}(\mathbf{x}) \leq \hat{f}(\bar{\mathbf{x}})$  we have by Theorem 10.2.4 that

$$\sum_{i=1}^n (f([i]) - f([i-1]))x_i \leq \hat{f}(\mathbf{x}) \leq \hat{f}(\bar{\mathbf{x}})$$

which implies that  $\mathbf{x}$  is not cut away by the hyperplane.

Next we show that  $\eta = 0$ . Our separating hyperplane induces a valid halfspace whenever it is not nonzero, i.e.  $f([i]) \neq f([i-1])$  for some  $i$ . In the case that it is zero  $f([i]) = f([i-1]) \forall i$ , by the same argument above, we have  $\hat{f}(\bar{x}) = \sum_{i=1}^n (f([i]) - f([i-1]))\bar{x}_i = 0$  and

$$\hat{f}(\mathbf{x}) \geq \sum_{i=1}^n (f([i]) - f([i-1]))x_i = 0 = \hat{f}(\bar{x}).$$

In other words,  $\bar{x}$  is an exact minimizer, i.e.  $\eta = 0$ .

Note that  $|\hat{f}(\mathbf{x})| = |\mathbb{E}_{t \sim [0,1]}[f(\{i : x_i \geq t\})]| \leq M$  as  $M = \max_S |f(S)|$ . Now plugging in  $\alpha = \frac{1}{4M}$  in the guarantee of Theorem 8.3.26, we can find a point  $x^*$  such that

$$\begin{aligned} \hat{f}(x^*) - \min_{\mathbf{x} \in [0,1]^n} \hat{f}(\mathbf{x}) &\leq \frac{1}{4M} \left( \max_{\mathbf{x} \in [0,1]^n} \hat{f}(\mathbf{x}) - \min_{\mathbf{x} \in [0,1]^n} \hat{f}(\mathbf{x}) \right) \\ &\leq \frac{1}{4M} (2M) \\ &< 1 \end{aligned}$$

We claim that  $\min_{t \in [0,1]} f(\{i : x_i^* \geq t\})$  is minimum. To see this, recall from 10.2.2 that  $\hat{f}$  has an integer minimizer and hence  $\min_{\mathbf{x} \in [0,1]^n} \hat{f}(\mathbf{x}) = \min_S f(S)$ . Moreover,  $\hat{f}(x^*)$  is a convex combination of  $f(\{i : x_i^* \geq t\})$  which gives

$$1 > \hat{f}(x^*) - \min_{\mathbf{x} \in [0,1]^n} \hat{f}(\mathbf{x}) = \hat{f}(x^*) - \min_S f(S) \geq \min_{t \in [0,1]} f(\{i : x_i^* \geq t\}) - \min_S f(S).$$

Since  $f$  is integer-valued, we must then have  $\min_{t \in [0,1]} f(\{i : x_i^* \geq t\}) = \min_S f(S)$  as desired. Since our separation oracle can be computed by  $n$  oracle calls and runs in time  $O(n \cdot \text{EO} + n^2)$ , by Theorem 9.3.2 the overall running time is then  $O(n^2 \log nM \cdot \text{EO} + n^3 \log^{O(1)} nM)$  as claimed.  $\square$

Needless to say the proof above completely depends on Theorem 9.3.2. We remark that one can use the Vaidya's cutting plane instead of ours to get a time complexity  $O(n^2 \log nM \cdot \text{EO} + n^{\omega+1} \log^{O(1)} n \cdot \log M)$ . There is actually an alternate argument that gives a time complexity of  $O(n^2 \log M \cdot \text{EO} + n^{O(1)} \cdot \log M)$ . Thus it requires slightly fewer oracle calls at the expense of slower running time. A proof is offered in this section, which can be skipped without any risk of discontinuation. This proof relies the following cutting plane method.

**Theorem 10.3.2** ([36]). *Given any convex set  $K \subset [0,1]^n$  with a separation oracle of cost  $SO$ , in time  $O(kSO + kn^{O(1)})$  one can find either find a point  $\mathbf{x} \in K$  or find a polytope  $P$  such that  $K \subset P$  and the volume of  $K$  is at most  $(\frac{2}{3})^k$ .*

The Theorem allows us to decrease the volume of the feasible region by a factor of  $(\frac{2}{3})^k$  after  $k$  iterations. Similar to above, we apply cutting plane to minimize  $\hat{f}$  over the hypercube  $[0,1]^n$  for  $O(n \log M)$  iterations, and outputs *any* integral point in the remaining feasible region  $P$ .

**Lemma 10.3.3.** *Let  $x^*$  achieve the minimum function value  $\hat{f}(x^*)$  among the points used to query the separation oracle. Then*

1.  $x^* \in P^{(k)}$ , the current feasible region.
2. Any  $\mathbf{x}$  with  $\hat{f}(\mathbf{x}) \leq \hat{f}(x^*)$  belongs to  $P^{(k)}$ .
3. suppose  $x_{i_1}^* \geq \dots \geq x_{i_n}^*$  and let  $S_j = \{i_1, \dots, i_j\}$ . Then  $S_l \in \arg \min_{S_j} f(S_j)$  also belongs to  $P^{(k)}$ .

*Proof.* For any separating hyperplane  $\mathbf{h}^T x \leq \hat{f}(\bar{x})$  given by  $\bar{x}$ , we have by Lemma 10.2.6 that  $\mathbf{h}^T x^* \leq \hat{f}(x^*)$ . Since  $\hat{f}(x^*)$  is the minimum among all  $\hat{f}(\bar{x})$ ,  $\mathbf{h}^T x^* \leq \hat{f}(\bar{x})$  and hence  $x^*$  is not removed by any new separating hyperplane. In other words,  $x^* \in P^{(k)}$ . The argument for (2) is analogous.

For (3), recall that by the definition of Lovász extension  $\hat{f}(x^*)$  is a convex combination of  $f(S_j)$  and thus the indicator variable  $I_{S_l}$  for  $S_l$  satisfies  $f(I_{S_l}) \leq \hat{f}(x^*)$ . By Lemma 10.2.6 again, this implies  $\mathbf{h}^T I_{S_l} \leq f(I_{S_l}) \leq \hat{f}(x^*) \leq \hat{f}(\bar{x})$  for any separating hyperplane  $\mathbf{h}^T x \leq \hat{f}(\bar{x})$ .  $\square$

**Theorem 10.3.4.** *Suppose that we run Cutting Plane in Theorem 10.3.2 for  $O(n \log M)$  iterations. Then  $S_l$  from the last lemma also minimizes  $f$ .*

*Proof.* We use the notations from the last lemma. After  $k = Kn \log_{2/3} M$  iterations, the volume of the feasible region  $P^{(k)}$  is at most  $1/M^{Kn}$ . By the last lemma,  $I_{S_l} \in P^{(k)}$ .

Suppose for the sake of contradiction that  $S$  minimizes  $f$  but  $f(S) < f(S_l)$ . Since  $f$  is integer-valued,  $f(S) + 1 \leq f(S_l)$ . Let  $r \stackrel{\text{def}}{=} 1/6M$ . Consider the set  $B \stackrel{\text{def}}{=} \{\mathbf{x} : 0 \leq x_i \leq r \forall i \notin S, 1 - r \leq x_i \leq 1 \forall i \in S\}$ . We claim that for  $\mathbf{x} \in B$ ,

$$\hat{f}(\mathbf{x}) \leq f(S) + 1.$$

To show this, note that  $f(\{i : x_i \geq t\}) = f(S)$  for  $r < t \leq 1 - r$  as  $x_i \leq r$  for  $i \notin S$  and  $x_i \geq 1 - r$  for  $i \in S$ . Now using conditional probability and  $|f(T)| \leq M$  for any  $T$ ,

$$\begin{aligned} \hat{f}(\mathbf{x}) &= \mathbb{E}_{t \sim [0,1]}[f(\{i : x_i \geq t\})] \\ &= (1 - 2r) \mathbb{E}[f(\{i : x_i \geq t\}) | r < t \leq 1 - r] + \\ &\quad r (\mathbb{E}[f(\{i : x_i \geq t\}) | 0 \leq t \leq r] + \mathbb{E}[f(\{i : x_i \geq t\}) | 1 - r \leq t \leq 1]) \\ &= (1 - r) f(S) + r (\mathbb{E}[f(\{i : x_i \geq t\}) | 0 \leq t \leq r] + \mathbb{E}[f(\{i : x_i \geq t\}) | 1 - r \leq t \leq 1]) \\ &\leq (1 - 2r) f(S) + 2rM \\ &\leq f(S) + 4rM \\ &\leq f(S) + 1 \end{aligned}$$

But now  $B \subseteq P^{(k)}$  as  $\hat{f}(\mathbf{x}) \leq f(S) + 1 \leq f(S_l)$  and by (2) of the last lemma. This would lead to a contradiction since

$$\text{vol}(B) = \frac{1}{(6M)^n} > \frac{1}{M^{Kn}} \geq \text{vol}(P^{(k)})$$

for sufficiently large  $K$ .  $\square$

**Corollary 10.3.5.** *There is an  $O(n^2 \log M \cdot \text{EO} + n^{O(1)} \log M)$  time algorithm for submodular function minimization.*

*Proof.* This simply follows from the last lemma, Theorem 10.3.2, and the fact that our separation oracle runs in time  $O(n \cdot \text{EO} + n^2)$ .  $\square$

Curiously, we obtained  $O(\log M)$  rather than  $O(\log nM)$  as in our algorithm. We leave it as an open problem whether one can slightly improve our running time to  $O(n^2 \log M \cdot \text{EO} + n^3 \log^{O(1)} n \cdot \log M)$ . The rest of this chapter is devoted to obtaining better strongly polynomial running time.

## ■ 10.4 Improved Strongly Polynomial Algorithms for SFM

In this section we show how our cutting plane method can be used to obtain a  $\tilde{O}(n^3 \cdot \text{EO} + n^4)$  time algorithm for SFM, which improves over the currently fastest  $O(n^5 \cdot \text{EO} + n^6)$  time algorithm by Orlin.

### ■ 10.4.1 Improved Oracle Complexity

We first present a simple geometric argument that  $f$  can be minimized with just  $O(n^3 \log n \cdot \text{EO})$  oracle calls. While this is our desired query complexity (and it improves upon the previous best known bounds by a factor of  $O(n^2)$ ) unfortunately the algorithm runs in exponential time. Nevertheless, it does provide some insight into how our more efficient algorithms should proceed and it alone, does suggest that information theoretically,  $O(n^3 \log n \cdot \text{EO})$  calls suffice to solve SFM. In the rest of this chapter, we combine this insight with some of the existing SFM tools developed over the last decade to get improved polynomial time algorithms.

**Theorem 10.4.1.** *Submodular functions can be minimized with  $O(n^3 \log n \cdot \text{EO})$  oracle calls.*

*Proof.* We use the cutting plane method in Theorem 10.3.2 with the separation oracle given by Theorem 10.2.4. This method reduce the volume of the feasible region by a factor of  $(\frac{2}{3})^k$  after  $k$  iterations if the optimal has not found yet.

Now, we argue that after  $O(n \log n)$  iterations of this procedure we have either found a minimizer of  $f$  or we have enough information to reduce the dimension of the problem by 1. To see this, first note that if the separation oracle ever returns a degenerate hyperplane, then by Lemma 10.2.8 then either  $\emptyset$  or  $V$  is the minimizer, which we can determine in time  $O(\text{EO} + n)$ . Otherwise, after  $100n \log n$  iterations, our feasible region  $P$  must have a volume of at most  $1/n^{10n}$ . In this case, we claim that the remaining integer points in  $P$  all lie on a hyperplane. This holds, as if this was not the case, then there is a simplex  $\Delta$ , with integral vertices  $v_0, v_1, \dots, v_n$ , contained in  $P$ . But then

$$\text{vol}(P) \geq \text{vol}(\Delta) = \frac{1}{n!} |\det(v_1 - v_0 \ v_2 - v_0 \ \dots \ v_n - v_0)| \geq \frac{1}{n!}$$

where the last inequality holds since the determinant of an integral matrix is integral, yielding a contradiction.

In other words after  $O(n \log n)$  iterations, we have reduced the dimension of all viable solutions by at least 1. Thus, we can recurse by applying the cutting plane method to the lower dimensional feasible region, i.e.  $P$  is (replaced by) the convex combination of all the remaining integer points. There is a minor technical issue we need to address as our space is now lower dimensional and the starting region is not necessarily the hypercube anymore and the starting volume is not necessarily equal to 1.

We argue that the starting volume is bounded by  $n^{O(n)}$ . If this is indeed the case, then our previous argument still works as the volume goes down by a factor of  $1/n^{O(n)}$  in  $O(n \log n)$  iterations.

Let  $v \in P$  be an integer point. Now the  $\dim(P)$ -dimensional ball of radius  $\sqrt{n}$  centered at  $v$  must contain all the other integer points in  $P$  as any two points of  $\{0, 1\}^n$  are at most  $\sqrt{n}$  apart. Thus the volume of  $P$  is bounded by the volume of the ball which is  $n^{O(n)}$ . Now to get the volume down to  $1/n^{10n}$ , the number of iterations is still  $O(n \log n)$ .

In summary, we have reduced our dimension by 1 using  $O(n \log n)$  iterations which requires  $O(n^2 \log n \cdot \text{EO})$  oracle calls (as each separating hyperplane is computed with  $n \cdot \text{EO}$  oracle calls). This can happen at most  $n$  times. The overall query complexity is then  $O(n^3 \log n \cdot \text{EO})$ .

Note that the minimizer  $\mathbf{x}$  obtained may not be integral. This is not a problem as the definition of Lovász extension implies that if  $\hat{f}(\mathbf{x})$  is minimal, then  $f(\{i : x_i \geq t\})$  is minimal for any  $t \in [0, 1]$ .

We remark that this algorithm does not have a polynomial runtime. Even though all the integral vertices of  $P$  lie on a hyperplane, the best way we know of that identifies it takes exponential time by checking for all the integer points  $\{0, 1\}^n$ .  $\square$

*Remark 10.4.2.* Note that this algorithm works for minimizing any convex function over the hypercube that obtains its optimal value at a vertex of the hypercube. Formally, our proof of Theorem 10.4.1

holds whenever a function  $f : 2^V \rightarrow \mathbb{R}^n$  admits a convex relaxation  $\hat{f}$  with the following properties:

1. For every  $S \subseteq V$ ,  $\hat{f}(I_S) = f(S)$ .
2. Every  $\hat{f}(\mathbf{x})$  can be written as a convex combination  $\sum_{S \in \mathcal{S}} \alpha_S f(S)$ , where  $\sum \alpha_S = 1$ ,  $|\mathcal{S}| = O(n)$ , and  $S$  can be computed without any oracle call.
3. A subgradient  $\partial \hat{f}(\mathbf{x})$  of  $\hat{f}$  at any point  $\mathbf{x} \in [0, 1]^n$  can be computed with  $O(n \cdot \text{EO})$  oracle calls.

In this case, the proof of Theorem 10.4.1, implies that  $\hat{f}$  and  $f$  can be minimized with  $O(n^3 \log n \cdot \text{EO})$  oracle calls by using the separating hyperplane  $\partial \hat{f}(\bar{\mathbf{x}})^T (\mathbf{x} - \bar{\mathbf{x}}) \leq 0$ .

## ■ 10.4.2 Technical Tools

To improve upon the running time of the algorithm in the previous section, we use more structure of our submodular function  $f$ . Rather than merely showing that we can decrease the dimension of our SFM problem by 1 we show how we can reduce the degrees of freedom of our problem in a more principled way. In Section 10.4.2.1 we formally define the abstraction we use for this and discuss how to change our separation oracle to accommodate this abstraction, and in Section 10.4.2.2 we show how we can deduce these constraints. These tools serve as the foundation for the faster strongly polynomial time SFM algorithms we present in Section 10.4.3 and Section 10.4.4.

### 10.4.2.1 SFM over Ring Family

For the remainder of this chapter we consider a more general problem than SFM in which we wish to compute a minimizer of our submodular function  $f$  over a ring family of the ground set  $V = [n]$ . A ring family  $\mathcal{F}$  is a collection of subsets of  $V$  such that for any  $S_1, S_2 \in \mathcal{F}$ , we have  $S_1 \cup S_2, S_1 \cap S_2 \in \mathcal{F}$ . Thus SFM corresponds to the special case where  $\mathcal{F}$  consists of every subset of  $V$ . This generalization has been considered before in the literature and was essential to the IFF algorithm.

It is well known that any ring family  $\mathcal{F}$  over  $V$  can be represented by a directed graph  $D = (V, A)$  where  $S \in \mathcal{F}$  iff  $S$  contains all of the descendants of any  $i \in S$ . An equivalent definition is that for any arc  $(i, j) \in A$ ,  $i \in S$  implies  $j \in S$ . It is customary to assume that  $A$  is acyclic as any (directed) cycle of  $A$  can be contracted (see section 10.4.3.1).

We denote by  $R(i)$  the set of descendants of  $i$  (including  $i$  itself) and  $Q(i)$  the set of ancestors of  $i$  (including  $i$  itself). Polyhedrally, an arc  $(i, j) \in A$  can be encoded as the constraint  $x_i \leq x_j$  as shown by the next lemma.

**Lemma 10.4.3.** *Let  $\mathcal{F}$  be a ring family over  $V$  and  $D = (V, A)$  be its directed acyclic graph representation. Suppose  $f : V \rightarrow \mathbb{R}$  is submodular with Lovász extension  $\hat{f}$ . Then the characteristic vector  $I_S$  of any minimizer  $S = \arg \min_{S \in \mathcal{F}} f(S)$  over  $\mathcal{F}$  is also the solution to*

$$\begin{aligned} \min \hat{f}(\mathbf{x}) \\ x_i \leq x_j \forall (i, j) \in A \\ 0 \leq \mathbf{x} \leq 1 \end{aligned} \tag{10.1}$$

*Proof.* Let  $x^*$  be a minimizer, and  $L^{(t)} = \{i : x_i^* \geq t\}$ . It is easy to check that the indicator variable  $I_{L^{(t)}}$  satisfies (10.1) since  $x^*$  does. Moreover, recall that  $\hat{f}(x^*) = \mathbb{E}_{t \sim [0, 1]} [f(L_t)]$ . Thus  $\hat{f}(x^*)$  can be written as a convex combination  $\hat{f}(x^*) = \sum_t \alpha_t f(L^{(t)}) = \sum_t \alpha_t \hat{f}(I_{L^{(t)}})$ , where  $\alpha_t > 0$  and  $\sum_t \alpha_t = 1$ . Thus all such  $\hat{f}(I_{L^{(t)}})$  are minimal, i.e. (10.1) has no “integrality gap”.  $\square$

We also modify our separation oracle to accommodate for this generalization as follows. Before doing so we need a definition which relates our BFS to the ring family formalism.

**Definition 10.4.4.** A permutation  $(v_1, \dots, v_n)$  of  $V$  is said to be *consistent* with an arc  $(i, j)$  if  $j$  precedes  $i$  in  $(v_1, \dots, v_n)$ . Similarly, a BFS of the base polyhedron is consistent with  $(i, j)$  if  $j$  precedes  $i$  in its defining permutation.  $(v_1, \dots, v_n)$  (or a BFS) is consistent with  $A$  if it is consistent with every  $(i, j) \in A$ .

Readers may find it helpful to keep in mind the following picture which depicts the relative positions between  $R(i), i, Q(i)$  in the defining permutation of  $\mathbf{h}$  that is consistent with  $A$ :

$$\dots\dots R(i) \setminus \{i\} \dots\dots i \dots\dots Q(i) \setminus \{i\} \dots\dots$$

In Theorem 10.2.4, given  $\bar{x} \in [0, 1]^n$  our separating hyperplane is constructed by sorting the entries of  $\bar{x}$ . This hyperplane is associated with some BFS  $\mathbf{h}$  of the base polyhedron. As we shall see towards the end of the section, we would like  $\mathbf{h}$  to be consistent with every arc  $(i, j) \in A$ .

This task is easy initially as  $\bar{x}$  satisfies  $x_i \leq x_j$  for  $(i, j) \in A$  for the starting polytope of (10.1). If  $x_i < x_j$ , nothing special has to be done as  $j$  must precede  $i$  in the ordering. On the other hand, whenever  $x_i = x_j$ , we can always break ties by ranking  $j$  ahead of  $i$ .

However, a technical issue arises due to the fact that our cutting plane algorithm may drop constraints from the current feasible region  $P$ . In other words,  $\bar{x}$  may violate  $x_i \geq 0$ ,  $x_j \leq 1$  or  $x_i \leq x_j$  if it is ever dropped. Fortunately this can be fixed by reintroducing the constraint. We summarize the modification needed in the pseudocode below and formally show that it fulfills our requirement.

---

**Algorithm 25:** Modified Separation Oracle

---

**Input:**  $\bar{x} \in \mathbb{R}^n$  and the set of arcs  $A$   
**if**  $\bar{x}_i < 0$  *for some*  $i$  **then**  
  | **Output:**  $x_i \geq 0$   
**else if**  $\bar{x}_j > 1$  *for some*  $j$  **then**  
  | **Output:**  $x_j \leq 1$   
**else if**  $\bar{x}_i > \bar{x}_j$  *for some*  $(i, j) \in A$  **then**  
  | **Output:**  $x_i \leq x_j$   
**else**  
  | Let  $i_1, \dots, i_n$  be a permutation of  $V$  such that  $\bar{x}_{i_1} \geq \dots \geq \bar{x}_{i_n}$  and for all  $(i, j) \in A$ ,  $j$   
  | precedes  $i$  in  $i_1, \dots, i_n$ .  
  | **Output:**  $\mathbf{h}^T \mathbf{x} \leq \hat{f}(\bar{x})$ , where  $\mathbf{h}$  is the BFS defined by the permutation  $i_1, \dots, i_n$ .

---

**Lemma 10.4.5.** *Our modified separation oracle returns either some BFS  $\mathbf{h} = 0$  or a valid separating hyperplane, i.e.*

1.  $\bar{x}$  either lies on the separating hyperplane or is cut away by it.
2. Any minimizer of (10.1) is not cut away by the separating hyperplane.

Such a hyperplane can be computed with  $n$  oracle calls to  $f$  and in time  $O(n \cdot EO + n^2)$ .

*Proof.* If we get  $x_i \geq 0$ ,  $x_j \leq 1$  or  $x_i \leq x_j$  (if loop or the first two else loops), then clearly  $\bar{x}$  is cut away by it and any minimizer must of course satisfy  $x_i \geq 0$ ,  $x_j \leq 1$  and  $x_i \leq x_j$  as they are the constraints in (10.1). This proves (1) and (2) for the case of getting  $x_i \geq 0$ ,  $x_j \leq 1$  or  $x_i \leq x_j$ .



Thus it remains to consider the case  $\mathbf{h}^T \mathbf{x} \leq \hat{f}(\bar{x})$  (last else loop). First of all,  $\bar{x}$  lies on it as  $\hat{f}(\bar{x}) = \mathbf{h}^T \bar{x}$ . This proves (1). For (2), we have from Lemma 10.2.6 that  $\mathbf{h}^T \mathbf{x} \leq \hat{f}(\mathbf{x})$ . If  $x^*$  is a minimizer of (10.1), we must then have  $\mathbf{h}^T x^* \leq \hat{f}(x^*) \leq \hat{f}(\bar{x})$  as  $\bar{x}$  is also feasible for (10.1).

Finally we note that the running time is self-evident.  $\square$

We stress again that the main purpose of modifying our separation oracle is to ensure that any BFS  $\mathbf{h}$  used to define a new separating hyperplane must be consistent with every  $(i, j) \in A$ .

### 10.4.2.2 Identifying New Valid Arcs

The reason for considering the ring family generalization of SFM is that our algorithms (and some previous algorithms too) work by adding new arcs to our digraph  $D$ . This operation yields a strongly polynomial algorithm since there are only  $2 \cdot \binom{n}{2}$  possible arcs to add. Of course, a new arc  $(i, j)$  is valid only if  $i \in S_{\min} \implies j \in S_{\min}$  for some minimizer  $S_{\min}$ . Here we show how to identify such valid arcs by extracting information from certain nice elements of the base polyhedron.

This is guaranteed by the next four lemmas, which are stated in a way different from previous works e.g. our version is extended to the ring family setting. This is necessary as our algorithms require a more general formulation. We also give a new polyhedral proof, which is mildly simpler than the previous combinatorial proof. On the other hand, Lemma 10.4.10 is new and unique to our work. It is an important ingredient of our  $\tilde{O}(n^3 \cdot \text{EO} + n^4)$  time algorithm.

Recall that each BFS of the base polyhedron is defined by some permutation of the ground set elements.

First, we prove the following two lemmas which show that should we ever encounter a non-degenerate point in the base polytope with a coordinate of very large value, then we can immediately conclude that that coordinate must be or must not be in solution to SFM over the ring family.

**Lemma 10.4.6.** *If  $\mathbf{y} \in \mathcal{B}(f)$  is non-degenerate and satisfies  $y_i > -(n-1) \min_j y_j$ , then  $i$  is not in any minimizer of  $f$  (over the ring family  $A$ ).*

*Proof.* We proceed by contradiction and suppose that  $S$  is a minimizer of  $f$  that contains  $i$ . Now since  $\mathbf{y}$  is non-degenerate we know that  $\min_j y_j \leq 0$  and by the definition of  $\mathbf{y}$  we have the following contradiction

$$0 < y_i + (n-1) \min_j y_j \leq \sum_{j \in S} y_j = \mathbf{y}(S) \leq f(S) \leq f(\emptyset) = 0.$$

$\square$

**Lemma 10.4.7.** *If  $\mathbf{y} \in \mathcal{B}(f)$  is non-degenerate and satisfies  $y_i < -(n-1) \max_j y_j$ , then  $i$  is in every minimizer of  $f$  (over the ring family  $A$ ).*

*Proof.* We proceed by contradiction and suppose that  $S$  is a minimizer of  $f$  that does not contain  $i$ . Now since  $\mathbf{y}$  is non-degenerate we know that  $\max_j y_j \geq 0$  and therefore

$$\sum_{j \in [n]} y_j = y_i + \sum_{j \in S} y_j + \sum_{j \in V - (S+i)} y_j < -(n-1) \max_j y_j + \sum_{j \in S} y_j + (|V| - |S| - 1) \max_j y_j \leq \sum_{j \in S} y_j.$$

However by the definition of  $\mathbf{y}$  we have

$$\sum_{j \in S} y_j = \mathbf{y}(S) \leq f(S) \leq f(V) = \sum_{j \in [n]} y_j.$$

Thus we have a contradiction and the result follows.  $\square$

Now we are ready to present conditions under which a new valid arc can be added. We begin with a simple observation. Let  $\mathbf{upper}(i) \stackrel{\text{def}}{=} f(R(i)) - f(R(i) - i)$  and  $\mathbf{lower}(i) \stackrel{\text{def}}{=} f(V \setminus Q(i) + i) - f(V \setminus Q(i))$ . As the names suggest, they bound the value of  $h_i$  for any BFS used.

**Lemma 10.4.8.** *For any BFS  $\mathbf{h}$  used to construct a separating hyperplane given by our modified separation oracle, we have  $\mathbf{lower}(i) \leq h_i \leq \mathbf{upper}(i)$ .*

*Proof.* Note that by Lemma 10.4.5,  $\mathbf{h}$  is consistent with every  $(j_1, j_2) \in A$  and hence  $i$  must precede  $Q(i)$  and be preceded by  $R(i)$ . Let  $S$  be the set of elements preceding  $i$  in the defining permutation of  $\mathbf{h}$ . Then  $h_i = f(S + i) - f(S) \leq f(R(i)) - f(R(i) - i)$  because of diminishing return and  $R(i) - i \subseteq S$ . The lower bound follows from the same argument as  $Q(i) - i$  comes after  $i$ , and so  $Q(i) \subseteq V \setminus S$ .  $\square$

In the following two lemmas, we show that if  $\mathbf{upper}(i)$  is ever sufficiently positive or  $\mathbf{lower}(i)$  is sufficiently negative, then we find a new arc.

While these lemmas may appear somewhat technical but actually has an intuitive interpretation. Suppose an element  $p$  is in a minimizer  $S_{\min}$  of  $f$  over the ring family  $D$ . Then  $R(p)$  must also be part of  $S_{\min}$ . Now if  $f(R(p))$  is very large relative to  $f(R(p) - p)$ , there should be some element  $q \in S_{\min} \setminus R(p)$  compensating for the discrepancy. The lemma says that such an element  $q$  can in fact be found efficiently.

**Lemma 10.4.9** (new arc). *Let  $\mathbf{y} = \sum_k \lambda^{(k)} \mathbf{y}^{(k)}$  be a non-degenerate convex combination of  $O(n)$  base polyhedron BFS's  $\mathbf{y}^{(k)}$  which are consistent with every arc  $(i, j) \in A$ . If some element  $p$  satisfies  $\mathbf{upper}(p) > n^4 \max y_j$ , then we can find, using  $O(n \cdot EO)$  oracle calls and  $O(n^2)$  time, some  $q \notin R(p)$  such that the arc  $(p, q)$  is valid, i.e. if  $p$  is in a minimizer, then so is  $q$ .*

*Proof.* If  $\max y_j < 0$  then we are immediately done by Lemma 10.2.8. We assume  $\max y_j \geq 0$  in the proof. For all  $k$  let  $\mathbf{y}'^{(k)}$  be the BFS obtained by taking the defining permutation of  $\mathbf{y}^{(k)}$  and moving  $R(p)$  to the front while preserving the relative ordering of  $R(p)$  within each permutation). Furthermore, let  $\mathbf{y}' \stackrel{\text{def}}{=} \sum_k \lambda^{(k)} \mathbf{y}'^{(k)}$ . Then since  $y'_p = f(R(p)) - f(R(p) - p) = \mathbf{upper}(p)$  we have  $\mathbf{upper}(p) = y'_p = f(R(p)) - f(R(p) - p)$ . Moreover,

$$y'_j \geq y_j \quad \forall j \in R(p) \quad \text{and} \quad y'_j \leq y_j \quad \forall j \notin R(p) \quad (10.2)$$

by diminishing marginal return.

Now, suppose  $p$  is in a minimizer  $S_{\min}$ . Then  $R(p) \subseteq S_{\min}$  by definition. We then define  $f'(S) = f(S \cup R(p))$  for  $S \subseteq V \setminus R(p)$ . It can be checked readily that  $f'$  is submodular and  $S_{\min} \setminus R(p)$  is a minimizer of  $f'$  (over the corresponding ring family). Note that now  $\mathbf{y}'_{V \setminus R(p)}$  (the restriction of  $\mathbf{y}'$  to  $V \setminus R(p)$ ) is a convex combination of the BFS's of the base polyhedron  $\mathcal{B}(f')$  of  $f'$ . We shall show that  $\mathbf{y}'_{V \setminus R(p)}$  has the desired property in Lemma 10.4.7.

Note that  $y'(V \setminus R(p) + p) \leq y(V \setminus R(p) + p)$  since

$$y'(V \setminus R(p) + p) = y'(V) - y'(R(p) - p) = y(V) - y'(R(p) - p) \leq y(V) - y(R(p) - p) = y(V \setminus R(p) + p).$$

But now since  $\mathbf{y}$  is non-degenerate  $\max_j y_j \geq 0$  and therefore

$$\begin{aligned} y'(V \setminus R(p)) &\leq y(V \setminus R(p) + p) - y'_p \\ &= y(V \setminus R(p) + p) - (f(R(p)) - f(R(p) - p)) \\ &\leq n \max y_j - (f(R(p)) - f(R(p) - p)) \\ &< (n - n^4) \max y_j \end{aligned} \quad (10.3)$$

Therefore by the Pigeonhole Principle some  $q \notin R(p)$  must satisfy

$$\begin{aligned} y'_q &< ((n - n^4) \max y_j) / (n - 1) \\ &= -(n^3 + n^2 + n) \max y_j \\ &\leq -(n^3 + n^2 + n) \max_{j \notin R(p)} y_j \\ &\leq -(n^3 + n^2 + n) \max_{j \notin R(p)} y'_j \quad \text{by (10.2)} \end{aligned}$$

By Lemma 10.4.7, this  $q$  must be in any minimizer of  $f'$ . In other words, whenever  $p$  is in a minimizer of  $f$ , then so is  $q$ .

Note however that computing all  $\mathbf{y}'$  would take  $O(n^2)$  oracle calls in the worst case as there are  $O(n)$   $\mathbf{y}'^{(k)}$ 's. We use the following trick to identify some  $q$  with  $y'_q < -(n - 1) \max y_j$  using just  $O(n)$  calls. The idea is that we actually only want to have sufficient decreases in  $y'(V \setminus R(p))$  which can be accomplished by having a large corresponding decrease in some  $\mathbf{y}'^{(k)}$ .

For each  $k$ , by the same argument above (see (10.3))

$$y'^{(k)}(V \setminus R(p)) - y^{(k)}(V \setminus R(p)) \leq y_p^{(k)} - (f(R(p)) - f(R(p) - p)) \quad (10.4)$$

The “weighted decrease”  $\lambda^{(k)} (y_p^{(k)} - (f(R(p)) - f(R(p) - p)))$  for  $\mathbf{y}'^{(k)}$  sum up to

$$\sum \lambda^{(k)} (y_p^{(k)} - (f(R(p)) - f(R(p) - p))) = y_p - (f(R(p)) - f(R(p) - p)) < (1 - n^4) \max y_j$$

Thus by the Pigeonhole Principle, some  $l$  will have

$$\lambda^{(l)} (y_p^{(l)} - (f(R(p)) - f(R(p) - p))) < ((1 - n^4) \max y_j) / O(n) < -n^2 \max y_j.$$

For this  $\mathbf{y}^{(l)}$  we compute  $\mathbf{y}''$ . We show that  $\mathbf{y}'' = \lambda^{(l)} \mathbf{y}'^{(l)} + \sum_{k \neq l} \lambda^{(k)} \mathbf{y}^{(k)}$  has the same property as  $\mathbf{y}'$  above.

$$\begin{aligned} y''(V \setminus R(p)) &= \lambda^{(l)} y'^{(l)}(V \setminus R(p)) + \sum_{k \neq l} \lambda^{(k)} y^{(k)}(V \setminus R(p)) \\ &= y(V \setminus R(p)) + \lambda^{(l)} (y'^{(l)}(V \setminus R(p)) - y^{(l)}(V \setminus R(p))) \\ &\leq y(V \setminus R(p)) + \lambda^{(l)} (y_p^{(l)} - (f(R(p)) - f(R(p) - p))) \quad \text{by (10.4)} \\ &< (n - 1) \max y_j - n^2 \max y_j \\ &< (n - n^2) \max y_j \end{aligned}$$

Then some  $q \in V \setminus R(p)$  must satisfy

$$y''_q < \frac{n - n^2}{n - 1} \max y_j = -n \max y_j$$

That is, the arc  $(p, q)$  is valid. This takes  $O(n)$  oracle calls as given  $\mathbf{y} = \sum_k \lambda^{(k)} \mathbf{y}^{(k)}$ , computing  $\mathbf{y}''$  requires knowing only  $f(R(p))$ ,  $f(R(p) - p)$ , and  $\mathbf{y}'^{(l)}$  which can be computed from  $\mathbf{y}^{(l)}$  with  $n$  oracle calls. The runtime is  $O(n^2)$  which is needed for computing  $\mathbf{y}''$ .  $\square$

**Lemma 10.4.10.** *Let  $\mathbf{y} = \sum_k \lambda^{(k)} \mathbf{y}^{(k)}$  be a non-degenerate convex combination of base polyhedron BFS  $\mathbf{y}^{(k)}$  which is consistent with every arc  $(i, j) \in A$ . If  $\mathbf{lower}(p) < n^4 \min y_j$ , then we can find, using  $O(n \cdot EO)$  oracle calls and  $O(n^2)$  time, some  $q \notin Q(p)$  such that the arc  $(q, p)$  is valid, i.e. if  $p$  is not in a minimizer, then  $q$  is not either.*

*Proof.* It is possible to follow the same recipe in the proof of Lemma 10.4.9 but using Lemma 10.4.6 instead of Lemma 10.4.7. Here we offer a proof which directly invokes Lemma 10.4.7 on a different submodular function.

Let  $g$  be defined by  $g(S) \stackrel{\text{def}}{=} f(V \setminus S)$  for any  $S$ , and  $A_g$  be the set of arcs obtained by reversing the directions of the arcs of  $A$ . Consider the problem of minimizing  $g$  over the ring family  $A_g$ . Using subscripts to avoid confusion with  $f$  and  $g$ , e.g.  $R_g(i)$  is the set of descendants of  $i$  w.r.t.  $A_g$ , it is not hard to verify the following:

- $g$  is submodular
- $R_g(i) = Q_f(i)$
- $g(R_g(p)) - g(R_g(p) - p) = -(f(V \setminus Q_f(p) + p) - f(V \setminus Q_f(p)))$
- $-\mathbf{y}^{(k)}$  is a BFS of  $\mathcal{B}(g)$  if and only if  $\mathbf{y}^{(k)}$  is a BFS of  $\mathcal{B}(f)$
- $\max(-y_j) = -\min y_j$

By using the above correspondence and applying Lemma 10.4.9 to  $g$  and  $A_g$ , we can find, using  $O(n)$  oracle calls and  $O(n^2)$  time, some  $q \notin R_g(p) = Q(p)$  such that the arc  $(p, q)$  is valid for  $g$  and  $A_g$ . In other words, the reverse  $(q, p)$  will be valid for  $f$  and  $A$ .  $\square$

These lemmas lay the foundation of our algorithm. They suggests that if the positive entries of a point in the base polyhedron are small relative to some  $\text{upper}(p) = f(R(p)) - f(R(p) - p)$ , a new arc  $(p, q)$  can be added to  $A$ . This can be seen as a robust version of Lemma 10.2.8.

Finally, we end the section with a technical lemma that will be used crucially for both of our algorithms. The importance of it would become obvious when it is invoked in our analyses.

**Lemma 10.4.11.** *Let  $\mathbf{h}''$  denote a convex combination of two vectors  $\mathbf{h}$  and  $\mathbf{h}'$  in the base polyhedron, i.e.  $\mathbf{h}'' = \lambda\mathbf{h} + (1 - \lambda)\mathbf{h}'$  for some  $\lambda \in [0, 1]$ . Further suppose that*

$$\|\mathbf{h}''\|_2 \leq \alpha \min \{ \lambda\|\mathbf{h}\|_2, (1 - \lambda)\|\mathbf{h}'\|_2 \}$$

for some  $\alpha \leq \frac{1}{2\sqrt{n}}$ . Then for  $p = \arg \max_j (\max\{\lambda|h_j|, (1 - \lambda)|h'_j|\})$  we have

$$\text{lower}(p) \leq -\frac{1}{2\alpha\sqrt{n}} \cdot \|\mathbf{h}''\|_\infty \quad \text{and} \quad \text{upper}(p) \geq \frac{1}{2\alpha\sqrt{n}} \cdot \|\mathbf{h}''\|_\infty \quad .$$

*Proof.* Suppose without loss of generality that  $\lambda|h_p| \geq (1 - \lambda)|h'_p|$ . Then by assumptions we have

$$\|\mathbf{h}''\|_\infty \leq \|\mathbf{h}''\|_2 \leq \alpha \cdot \min \{ \lambda\|\mathbf{h}\|_2, (1 - \lambda)\|\mathbf{h}'\|_2 \} \leq \alpha\sqrt{n} |\lambda h_p| \quad .$$

However, since  $\alpha \leq \frac{1}{2\sqrt{n}}$  we see that

$$|\lambda h_p + (1 - \lambda)h'_p| \leq \|\mathbf{h}''\|_\infty \leq \alpha\sqrt{n} |\lambda h_p| \leq \frac{1}{2} |\lambda h_p| \quad .$$

Consequently,  $\lambda h_p$  and  $(1 - \lambda)h'_p$  have opposite signs and  $|(1 - \lambda)h'_p| \geq \frac{1}{2} |\lambda h_p|$ . We then have,

$$\text{lower}(p) \leq \min \{ h_p, h'_p \} \leq \min \{ \lambda h_p, (1 - \lambda)h'_p \} \leq -\frac{1}{2} |\lambda h_p| \leq -\frac{1}{2\alpha\sqrt{n}} \|\mathbf{h}''\|_\infty$$

and

$$\text{upper}(p) \geq \max \{ h_p, h'_p \} \geq \max \{ \lambda h_p, (1 - \lambda)h'_p \} \geq \frac{1}{2} |\lambda h_p| \geq \frac{1}{2\alpha\sqrt{n}} \|\mathbf{h}''\|_\infty \quad .$$

$\square$

### ■ 10.4.3 $\tilde{O}(n^4 \cdot \text{EO} + n^5)$ Time Algorithm

Here we present a  $\tilde{O}(n^4 \cdot \text{EO} + n^5)$  time, i.e. strongly polynomial time algorithm, for SFM. We build upon the algorithm achieved in the section to achieve a faster running time in Section 10.4.4.

Our new algorithm combines the existing tools for SFM developed over the last decade with our cutting plane method. While there are certain similarities with previous algorithms (especially [118, 125, 121]), our approach significantly departs from all the old approaches in one important aspect.

All of the previous algorithms actively maintain a point in the base polyhedron and represent it as a *convex combination* of BFS's. At each step, a new BFS may enter the convex combination and an old BFS may exit. Our algorithm, on the other hand, maintains only a *collection* of BFS's (corresponding to our separating hyperplanes), rather than an explicit convex combination. A "good" convex combination is computed from the *collection* of BFS's only after running Cutting Plane for enough iterations. We believe that this crucial difference is the fundamental reason which offers the speedup. This is achieved by the Cutting Plane method which considers the *geometry* of the collection of BFS's. On the other hand, considering only a convex combination of BFS's effectively narrows our sight to only one *point* in the base polyhedron.

#### Overview

Now we are ready to describe our strongly polynomial time algorithm. Similar to the weakly polynomial algorithm, we first run our cutting plane for enough iterations on the initial feasible region  $\{\mathbf{x} \in [0, 1]^n : x_i \leq x_j \forall (i, j) \in A\}$ , after which a pair of approximately parallel supporting hyperplanes  $F_1, F_2$  of width  $1/n^{\Theta(1)}$  can be found. Our strategy is to write  $F_1$  and  $F_2$  as a nonnegative combination of the facets of remaining feasible region  $P$ . This combination is made up of newly added separating hyperplanes as well as the inequalities  $x_i \geq 0$ ,  $x_j \leq 1$  and  $x_i \leq x_j$ . We then argue that one of the following updates can be done:

- Collapsing:  $x_i = 0$ ,  $x_j = 1$  or  $x_i = x_j$
- Adding a new arc  $(i, j)$ :  $x_i \leq x_j$  for some  $(i, j) \notin A$

The former case is easy to handle by elimination or contraction. If  $x_i = 0$ , we simply eliminate  $i$  from the ground set  $V$ ; and if  $x_i = 1$ , we redefine  $f$  so that  $f(S) = f(S + i)$  for any  $S \subseteq V - i$ .  $x_i = x_j$  can be handled in a similar fashion. In the latter case, we simply add the arc  $(i, j)$  to  $A$ . We then repeat the same procedure on the new problem.

Roughly speaking, our strongly polynomial time guarantee follows as eliminations and contractions can happen at most  $n$  times and at most  $2 \cdot \binom{n}{2}$  new arcs can be added. While the whole picture is simple, numerous technical details come into play in the execution. We advise readers to keep this overview in mind when reading the subsequent sections.

#### Algorithm

Our algorithm is summarized below. Again, we remark that our algorithm simply uses Theorem 10.4.12 regarding our cutting plane and is agnostic as to how the cutting plane works, thus it could be replaced with other methods, albeit at the expense of slower runtime.

1. Run cutting plane on (10.1) (Theorem 10.4.12 with  $\tau = \Theta(1)$ ) using our modified separation oracle (Section 10.4.2.1).
2. Identify a pair of "narrow" approximately parallel supporting hyperplanes or get some BFS  $\mathbf{h} = 0$  (in which case both  $\emptyset$  and  $V$  are minimizers).

3. Deduce from the hyperplanes some new constraint of the forms  $x_i = 0, x_j = 1, x_i = x_j$  or  $x_i \leq x_j$  (Section 10.4.3.2).
4. Consolidate  $A$  and  $f$  (Section 10.4.3.1).
5. Repeat by running our cutting plane method on (10.1) with updated  $A$  and  $f$ . (Note that Any previously found separating hyperplanes are discarded.)

We call step (1) a *phase* of cutting plane. The minimizer can be constructed by unraveling the recursion.

### 10.4.3.1 Consolidating $A$ and $f$

Here we detail how the set of valid arcs  $A$  and submodular function  $f$  should be updated once we deduce new information  $x_i = 0, x_i = 1, x_i = x_j$  or  $x_i \leq x_j$ . Recall that  $R(i)$  and  $Q(i)$  are the sets of descendants and ancestors of  $i$  respectively (including  $i$  itself). The changes below are somewhat self-evident, and are actually used in some of the previous algorithms so we only sketch how they are done without a detailed justification.

Changes to the digraph representation  $D$  of our ring family include:

- $x_i = 0$ : remove  $Q(i)$  from the ground set and all the arcs incident to  $Q(i)$
- $x_i = 1$ : remove  $R(i)$  from the ground set and all the arcs incident to  $R(i)$
- $x_i = x_j$ : contract  $i$  and  $j$  in  $D$  and remove any duplicate arcs
- $x_i \leq x_j$ : insert the arc  $(i, j)$  to  $A$
- For the last two cases, we also contract the vertices on a directed cycle of  $A$  until there is no more. Remove any duplicate arcs.

Here we can contract any cycle  $(i_1, \dots, i_k)$  because the inequalities  $x_{i_1} \leq x_{i_2}, \dots, x_{i_{k-1}} \leq x_{i_k}, x_{i_k} \leq x_{i_1}$  imply  $x_{i_1} = \dots = x_{i_k}$ .

Changes to  $f$ :

- $x_i = 0$ : replace  $f$  by  $f' : 2^{V \setminus Q(i)} \rightarrow \mathbb{R}, f'(S) = f(S)$  for  $S \subseteq V \setminus Q(i)$
- $x_i = 1$ : replace  $f$  by  $f' : 2^{V \setminus R(i)} \rightarrow \mathbb{R}, f'(S) = f(S \cup R(i))$  for  $S \subseteq V \setminus R(i)$
- $x_i = x_j$ : see below
- $x_i \leq x_j$ : no changes to  $f$  needed if it does not create a cycle in  $A$ ; otherwise see below
- Contraction of  $C = \{i_1, \dots, i_k\}$ : replace  $f$  by  $f' : 2^{V \setminus C} \rightarrow \mathbb{R}, f'(S) = f(S)$  for  $S \subseteq V \setminus C$  and  $f'(S) = f((S - l) \cup C)$  for  $S \ni l$

Strictly speaking, these changes are in fact *not* needed as they will automatically be taken care of by our cutting plane method. Nevertheless, performing them lends a more natural formulation of the algorithm and simplifies its description.

### 10.4.3.2 Deducing New Constraints $x_i = 0$ , $x_j = 1$ , $x_i = x_j$ or $x_i \leq x_j$

Here we show how to deduce new constraints through the result of our cutting plane method. This is the most important ingredient of our algorithm. As mentioned before, similar arguments were used first by IFF [121] and later in [118, 125]. There are however two important differences for our method:

- We maintain a collection of BFS's rather a convex combination; a convex combination is computed and needed only after each phase of cutting plane.
- As a result, our results are proved mostly *geometrically* whereas the previous ones were proved mostly *combinatorially*.

Our ability to deduce such information hinges on the power of the cutting plane method in Chapter 8. We re-state our main result Theorem 8.3.26 in the language of SFM. Note that Theorem 10.4.12 is formulated in a fairly general manner in order to accommodate for the next section. Readers may wish to think  $\tau = \Theta(1)$  for now.

**Theorem 10.4.12** (Theorem 8.3.26 restated for SFM). *For any  $\tau \geq 100$ , applying our cutting plane method, Theorem 10.4.12, to (10.1) with our modified separation oracle (or its variant in Section 10.4.4) with high probability in  $n$  either*

1. Finds a degenerate BFS  $\mathbf{h} \geq \vec{0}$  or  $\mathbf{h} \leq \mathbf{0}$ .
2. Finds a polytope  $P$  consisting of  $O(n)$  constraints which are our separating hyperplanes or the constraints in (10.1). Moreover,  $P$  satisfies the following inequalities

$$\mathbf{c}^T \mathbf{x} \leq M \quad \text{and} \quad \mathbf{c}'^T \mathbf{x} \leq M',$$

both of which are nonnegative combinations of the constraints of  $P$ , where

$$\|\mathbf{c} + \mathbf{c}'\|_2 \leq \min\{\|\mathbf{c}\|_2, \|\mathbf{c}'\|_2\}/n^{\Theta(\tau)} \quad \text{and} \quad |M + M'| \leq \min\{\|\mathbf{c}\|_2, \|\mathbf{c}'\|_2\}/n^{\Theta(\tau)}.$$

Furthermore, the algorithm runs in expected time  $O(n^2 \tau \log n \cdot EO + n^3 \tau^{O(1)} \log^{O(1)} n)$ .

*Proof.* In applying Theorem 10.4.12 we let  $K$  be the set of minimizers of  $f$  over the ring family and the box is the hypercube with  $R = 1$ . We run cutting plane with our modified separation oracle (Lemma 10.4.5). The initial polytope  $P^{(0)}$  can be chosen to be, say, the hypercube. If some separating hyperplane is degenerate, then we have the desired result (and know that either  $\emptyset$  or  $V$  is optimal). Otherwise let  $P$  be the current feasible region. Note that  $P \neq \emptyset$ , because our minimizers of  $\hat{f}$  are all in  $P^{(0)}$  and  $P^{(k)}$  as they are never cut away by the separating hyperplanes.

Let  $\mathcal{S}$  be the collection of inequalities (10.1) as well as the separating hyperplanes  $\mathbf{h}^T \mathbf{x} \leq \hat{f}(\bar{x}_h) = \mathbf{h}^T \bar{x}_h$  used. By Theorem 8.3.26, all of our minimizers will be contained in  $P$ , consisting of  $O(n)$  constraints  $\mathbf{A}\mathbf{x} \geq \mathbf{b}$ . Each such constraint  $\mathbf{a}_i^T \mathbf{x} \geq b_i$  is a scaling and shifting of some inequality  $\mathbf{p}_i^T \mathbf{x} \geq q_i$  in  $\mathcal{S}$ , i.e.  $\mathbf{a}_i = \mathbf{p}_i/\|\mathbf{p}_i\|_2$  and  $b_i \leq q_i/\|\mathbf{p}_i\|_2$ .

By taking  $\varepsilon = 1/n^{\Theta(\tau)}$  with sufficiently large constant in  $\Theta$ , our theorem certifies that  $P$  has a narrow width by  $\mathbf{a}_1$ , some nonnegative combination  $\sum_{i=2}^{O(n)} t_i \mathbf{a}_i$  and point  $\mathbf{x}_o \in P$  with  $\|\mathbf{x}_o\|_\infty \leq 3\sqrt{n}R = 3\sqrt{n}$  satisfying the following:

$$\left\| \mathbf{a}_1 + \sum_{i=2}^{O(n)} t_i \mathbf{a}_i \right\|_2 \leq 1/n^{\Theta(\tau)}$$

$$0 \leq \mathbf{a}_1^T \mathbf{x}_o - \mathbf{b}_1 \leq 1/n^{\Theta(\tau)}$$

$$0 \leq \left( \sum_{i=2}^{O(n)} t_i a_i \right)^T \mathbf{x}_o - \sum_{i=2}^{O(n)} t_i b_i \leq 1/n^{\Theta(\tau)}$$

We convert these inequalities to  $\mathbf{p}$  and  $q$ . Let  $t'_i \stackrel{\text{def}}{=} t_i \cdot \|\mathbf{p}_1\|_2 / \|\mathbf{p}_i\|_2 \geq 0$ .

$$\begin{aligned} \left\| \mathbf{p}_1 + \sum_{i=2}^{O(n)} t'_i \mathbf{p}_i \right\|_2 &\leq \|\mathbf{p}_1\|_2 / n^{\Theta(\tau)} \\ 0 \leq \mathbf{p}_1^T \mathbf{x}_o - q_1 &\leq \|\mathbf{p}_1\|_2 / n^{\Theta(\tau)} \\ 0 \leq \left( \sum_{i=2}^{O(n)} t'_i \mathbf{p}_i \right)^T \mathbf{x}_o - \sum_{i=2}^{O(n)} t'_i q_i &\leq \|\mathbf{p}_1\|_2 / n^{\Theta(\tau)} \end{aligned}$$

We claim that<sup>1</sup>  $\mathbf{c} = -\mathbf{p}_1$ ,  $M = -q_1$ ,  $\mathbf{c}' = -\sum_{i=2}^{O(n)} t'_i \mathbf{p}_i$ ,  $M' = -\sum_{i=2}^{O(n)} t'_i q_i$  satisfy our requirement.

We first show that  $\|\mathbf{c} + \mathbf{c}'\|_2 \leq \min\{\|\mathbf{c}\|_2, \|\mathbf{c}'\|_2\} / n^{\Theta(\tau)}$ . We have  $\|\mathbf{c} + \mathbf{c}'\|_2 \leq \|\mathbf{c}\|_2 / n^{\Theta(\tau)}$  from the first inequality. If  $\|\mathbf{c}\|_2 \leq \|\mathbf{c}'\|_2$  we are done. Otherwise, by triangle inequality

$$\|\mathbf{c}'\|_2 - \|\mathbf{c}\|_2 \leq \|\mathbf{c} + \mathbf{c}'\|_2 \leq \|\mathbf{c}\|_2 / n^{\Theta(\tau)} \implies 2\|\mathbf{c}\|_2 \geq \|\mathbf{c}'\|_2$$

and hence  $\|\mathbf{c} + \mathbf{c}'\|_2 \leq \|\mathbf{c}\|_2 / n^{\Theta(\tau)} \leq \|\mathbf{c}'\|_2 / 2n^{\Theta(\tau)} = \|\mathbf{c}'\|_2 / n^{\Theta(\tau)}$ .

We also need to prove  $|M + M'| \leq \min\{\|\mathbf{c}\|_2, \|\mathbf{c}'\|_2\} / n^{\Theta(\tau)}$ . Summing the second and third inequalities,

$$-\|\mathbf{c}\|_2 / n^{\Theta(\tau)} \leq (\mathbf{c} + \mathbf{c}')^T \mathbf{x}_o - (M + M') \leq 0$$

Recall that we have  $\|\mathbf{x}_o\|_\infty \leq 3\sqrt{n}$ . Then

$$\begin{aligned} |M + M'| &\leq |(\mathbf{c} + \mathbf{c}')^T \mathbf{x}_o - (M + M')| + |(\mathbf{c} + \mathbf{c}')^T \mathbf{x}_o| \\ &\leq \|\mathbf{c}\|_2 / n^{\Theta(\tau)} + 3\sqrt{n} \|\mathbf{c} + \mathbf{c}'\|_2 \\ &\leq \|\mathbf{c}\|_2 / n^{\Theta(\tau)} + 3\sqrt{n} \|\mathbf{c}\|_2 / n^{\Theta(\tau)} \\ &= \|\mathbf{c}\|_2 / n^{\Theta(\tau)} \end{aligned}$$

as desired. Our result then follows as we proved  $2\|\mathbf{c}'\|_2 \geq \|\mathbf{c}\|_2$ .

Finally, we have the desired runtime as our modified separation oracle runs in time  $O(n \cdot \text{EO} + n^2 \log^{O(1)} n)$ .  $\square$

Informally, the theorem above simply states that after  $O(n\tau \log n)$  iterations of cutting plane, the remaining feasible region  $P$  can be sandwiched between two approximately parallel supporting hyperplanes of width  $1/n^{O(\tau)}$ . A good intuition to keep in mind is that every  $O(n)$  iterations of cutting plane reduces the minimum width by a constant factor.

*Remark 10.4.13.* As shown in the proof of Theorem 10.4.12, one of the two approximately parallel hyperplanes can actually be chosen to be a constraint of our feasible region  $P$ . However we do not exploit this property as it does not seem to help us and would break the notational symmetry in  $\mathbf{c}$  and  $\mathbf{c}'$ .

## Setup

<sup>1</sup>Minus signs is needed because we express our inequalities as e.g.  $\mathbf{h}^T \mathbf{x} \leq \mathbf{h}^T \bar{\mathbf{x}}_h$  whereas in Theorem 8.3.26,  $\mathbf{a}_i^T \mathbf{x} \geq b_i$  is used. We apologize for the inconvenience.



In each phase, we run cutting plane using Theorem 10.4.12 with  $\tau = \Theta(1)$ . If some separating hyperplane used is degenerate, we have found the minimizer by Lemma 10.2.8.

Now assume none of the separating hyperplanes is degenerate. By Theorem 10.4.12,  $P$  is sandwiched by a pair of approximately parallel supporting hyperplanes  $F, F'$  which are of width  $1/10n^{10}$  apart. The width here can actually be  $1/n^c$  for any constant  $c$  by taking a sufficiently large constant in Theta.

Here, we show how to deduce from  $F$  and  $F'$  some  $x_i = 0, x_j = 1, x_i = x_j$ , or  $x_i \leq x_j$  constraint on the minimizers of  $f$  over the ring family. Let

$$\mathbf{c}^T \mathbf{x} = \sum c_i x_i \leq M \quad \text{and} \quad \mathbf{c}'^T \mathbf{x} = \sum c'_i x_i \leq M'$$

be the inequality for  $F$  and  $F'$  such that

$$|M + M'|, \|\mathbf{c} + \mathbf{c}'\|_2 \leq \text{gap}, \quad \text{where } \text{gap} \stackrel{\text{def}}{=} \frac{1}{10n^{10}} \min\{\|\mathbf{c}\|_2, \|\mathbf{c}'\|_2\}.$$

By the same theorem we can write  $\mathbf{c}^T \mathbf{x} \leq M$  as a nonnegative combination of the constraints for  $P$ . Recall that the constraints for  $P$  take on four different forms: (1)  $-x_i \leq 0$ ; (2)  $x_j \leq 1$ ; (3)  $-(x_j - x_i) \leq 0$ ; (4)  $\mathbf{h}^T \mathbf{x} = \sum h_i x_i \leq \hat{f}(\bar{x}_h)$ . Here the first three types are present initially whereas the last type is the separating hyperplane added. As alleged previously, the coefficient vector  $\mathbf{h}$  corresponds to a BFS of the base polyhedron for  $f$ . Our analysis crucially exploits this property.

Thus suppose  $\mathbf{c}^T \mathbf{x} = \sum_i c_i x_i \leq M$  is a nonnegative combination of our constraints with weights  $\alpha_i, \beta_j, \gamma_{ij}, \lambda_h \geq 0$ . The number of (positive)  $\alpha_i, \beta_j, \gamma_{ij}, \lambda_h$  is at most  $O(n)$ . Here we denote separating hyperplanes by  $\mathbf{h}^T \mathbf{x} \leq \hat{f}(\bar{x}_h)$ . Let  $H$  be the set of BFS's used to construct separating hyperplanes.

$$\mathbf{c}^T \mathbf{x} = - \sum_i \alpha_i x_i + \sum_j \beta_j x_j + \sum_{(i,j) \in A} \gamma_{ij} (x_i - x_j) + \sum_{h \in H} \lambda_h \mathbf{h}^T \mathbf{x} \quad \text{and} \quad M = \sum_j \beta_j + \sum_{h \in H} \lambda_h \hat{f}(\bar{x}_h). \quad (10.5)$$

Similarly, we write the inequality for  $F'$  as a nonnegative combination of the constraints for  $P$  and the number of (positive)  $\alpha'_i, \beta'_j, \gamma'_{ij}, \lambda'_h$  is  $O(n)$ :

$$\mathbf{c}'^T \mathbf{x} = - \sum_i \alpha'_i x_i + \sum_j \beta'_j x_j + \sum_{(i,j) \in A} \gamma'_{ij} (x_i - x_j) + \sum_{h \in H} \lambda'_h \mathbf{h}^T \mathbf{x} \quad \text{and} \quad M' = \sum_j \beta'_j + \sum_{h \in H} \lambda'_h \hat{f}(\bar{x}_h). \quad (10.6)$$

We also scale  $\mathbf{c}, \mathbf{c}', \alpha, \alpha', \beta, \beta', \gamma, \gamma', \lambda, \lambda'$  so that

$$\sum_{h \in H} (\lambda_h + \lambda'_h) = 1$$

as this does not change any of our preceding inequalities regarding  $F$  and  $F'$ .

Now that  $F, F'$  have been written as combinations of our constraints, we have gathered the necessary ingredients to derive our new arc. We first give a geometric intuition why we would expect to be able to derive a new constraint. Consider the nonnegative combination making up  $F$ . We think of the coefficient  $\beta_j$  as the contribution of  $x_j \leq 1$  to  $F$ . Now if  $\beta_j$  is very large,  $F$  is “very parallel” to  $x_j \leq 1$  and consequently  $F'$  would miss  $x_j = 0$  as the gap between  $F$  and  $F'$  is small.  $P$  would then miss  $x_j = 0$  too as it is sandwiched between  $F$  and  $F'$ . Similarly, a large  $\alpha_i$  and a large  $\gamma_{ij}$  would respectively imply that  $x_i = 1$  and  $(x_i = 0, x_j = 1)$  would be missed. The same argument works for  $F'$  as well.

But on the other hand, if the contributions from  $x_i \geq 0, x_j \leq 1, x_i \leq x_j$  to both  $F$  and  $F'$  are small, then the supporting hyperplanes  $\mathbf{c}^T \mathbf{x} \leq \dots$  and  $\mathbf{c}'^T \mathbf{x} \leq \dots$  would be mostly made up of

separating hyperplanes  $\mathbf{h}^T \mathbf{x} \leq \hat{f}(\bar{x}_h)$ . By summing up these separating hyperplanes (whose coefficients form BFS's), we would then get a point in the base polyhedron which is very close to the origin 0. Moreover, by Lemma 10.4.11 and Lemma 10.4.9 we should then be able to deduce some interesting information about the minimizer of  $f$  over  $D$ .

The rest of this section is devoted to realizing the vision sketched above. We stress that while the algebraic manipulations may be long, they are simply the execution of this elegant geometric picture.

Now, consider the following weighted sum of  $\mathbf{h}^T \mathbf{x} \leq \hat{f}(\bar{x}_h)$ :

$$\left( \sum_{h \in H} \lambda_h \mathbf{h}^T + \sum_{h \in H} \lambda'_h \mathbf{h}^T \right) \mathbf{x} = \sum_{h \in H} \lambda_h \mathbf{h}^T \mathbf{x} + \sum_{h \in H} \lambda'_h \mathbf{h}^T \mathbf{x} \leq \sum_{h \in H} \lambda_h \hat{f}(\bar{x}_h) + \sum_{h \in H} \lambda'_h \hat{f}(\bar{x}_h).$$

Observe that  $\sum_{h \in H} \lambda_h \mathbf{h}^T + \sum_{h \in H} \lambda'_h \mathbf{h}^T$  is in the base polyhedron since it is a convex combination of BFS  $\mathbf{h}$ . Furthermore, using (10.5) and (10.6) this can also be written as

$$\begin{aligned} \left( \sum_{h \in H} \lambda_h \mathbf{h}^T + \sum_{h \in H} \lambda'_h \mathbf{h}^T \right) \mathbf{x} &= \left( \mathbf{c}^T \mathbf{x} + \sum \alpha_i x_i - \sum \beta_j x_j + \sum_{(i,j) \in A} \gamma_{ij} (x_j - x_i) \right) \\ &\quad + \left( \mathbf{c}'^T \mathbf{x} + \sum \alpha'_i x_i - \sum \beta'_j x_j + \sum_{(i,j) \in A} \gamma'_{ij} (x_j - x_i) \right) \end{aligned} \quad (10.7)$$

and

$$\begin{aligned} \sum_{h \in H} \lambda_h \hat{f}(\bar{x}_h) + \sum_{h \in H} \lambda'_h \hat{f}(\bar{x}_h) &= \left( M - \sum \beta_j \right) + \left( M' - \sum \beta'_j \right) \\ &= (M + M') - \sum \beta_j - \sum \beta'_j \end{aligned}$$

Furthermore, we can bound  $\mathbf{c}^T \mathbf{x} + \mathbf{c}'^T \mathbf{x}$  by  $\mathbf{c}^T \mathbf{x} + \mathbf{c}'^T \mathbf{x} \geq -\|\mathbf{c} + \mathbf{c}'\|_1 \geq -\sqrt{n}\|\mathbf{c} + \mathbf{c}'\|_2 \geq -\sqrt{n}\text{gap}$  as  $\mathbf{x} \leq 1$ . Since  $M + M' \leq \text{gap}$ , we obtain

$$\begin{aligned} LHS &\stackrel{\text{def}}{=} \sum \alpha_i x_i + \sum \alpha'_i x_i - \sum \beta_j x_j - \sum \beta'_j x_j + \sum_{(i,j) \in A} \gamma_{ij} (x_j - x_i) + \sum_{(i,j) \in A} \gamma'_{ij} (x_j - x_i) \\ &\leq 2\sqrt{n}\text{gap} - \sum \beta_j - \sum \beta'_j \end{aligned}$$

Geometrically, the next lemma states that if the contribution from, say  $x_i \geq 0$ , to  $F$  is too large, then  $F'$  would be forced to miss  $x_i = 1$  because they are close to one another.

**Lemma 10.4.14.** *Suppose  $\mathbf{x}$  satisfies (10.1),  $LHS \leq 2\sqrt{n}\text{gap} - \sum \beta_j - \sum \beta'_j$  and  $\alpha_i, \beta_j, \gamma_{ij}, \alpha'_i, \beta'_j, \gamma'_{ij} \geq 0$ .*

1. If  $\alpha_i > 2\sqrt{n}\text{gap}$  or  $\alpha'_i > 2\sqrt{n}\text{gap}$ , then  $x_i < 1$ .
2. If  $\beta_j > 2\sqrt{n}\text{gap}$  or  $\beta'_j > 2\sqrt{n}\text{gap}$ , then  $x_j > 0$ .
3. If  $\gamma_{ij} > 2\sqrt{n}\text{gap}$  or  $\gamma'_{ij} > 2\sqrt{n}\text{gap}$ , then  $0 \leq x_j - x_i < 1$ .

*Proof.* We only prove it for  $\alpha_i, \beta_j, \gamma_{ij}$  as the other case follows by symmetry.

Using  $0 \leq x \leq 1$  and  $x_i \leq x_j$  for  $(i, j) \in A$ , we have  $LHS \geq \alpha_i x_i - \sum \beta_j - \sum \beta'_j$ . Hence  $\alpha_i x_i \leq 2\sqrt{n}\text{gap}$  and we get  $x_i < 1$  if  $\alpha_i > 2\sqrt{n}\text{gap}$ .

Similarly,  $LHS \geq -\beta_k x_k - \sum_{j \neq k} \beta_j - \sum \beta'_j$  which gives  $-\beta_k x_k \leq 2\sqrt{n}\text{gap} - \beta_k$ . Then  $x_k > 0$  if  $\beta_k > 2\sqrt{n}\text{gap}$ .

Finally,  $LHS \geq \gamma_{ij} (x_j - x_i) - \sum \beta_j - \sum \beta'_j$  which gives  $\gamma_{ij} (x_j - x_i) \leq 2\sqrt{n}\text{gap}$ . Then  $x_j - x_i < 1$  if  $\gamma_{ij} > 2\sqrt{n}\text{gap}$ . We have  $x_i \leq x_j$  since  $(i, j) \in A$ .  $\square$

So if either condition of Lemma 10.4.14 holds, we can set  $x_i = 0$  or  $x_j = 1$  or  $x_i = x_j$  since our problem (10.1) has an integral minimizer and any minimizer of  $\hat{f}$  is never cut away by Lemma 10.4.5. Consequently, in this case we can reduce the dimension by at least 1. From now on we may assume that

$$\max\{\alpha_i, \alpha'_i, \beta_j, \beta'_j, \gamma_{ij}, \gamma'_{ij}\} \leq 2\sqrt{n}\text{gap}. \quad (10.8)$$

Geometrically, (10.8) says that if the supporting hyperplanes are both mostly made up of the separating hyperplanes, then their aggregate contributions to  $F$  and  $F'$  should be small in absolute value.

The next lemma identifies some  $p \in V$  for which  $f(R(p)) - f(R(p) - p)$  is “big”. This prepares for the final step of our approach which invokes Lemma 10.4.9.

**Lemma 10.4.15.** *Let  $\mathbf{y} \stackrel{\text{def}}{=} \sum_{h \in H} \lambda_h \mathbf{h}$  and  $\mathbf{y}' \stackrel{\text{def}}{=} \sum_{h \in H} \lambda'_h \mathbf{h}$  and let  $p \in \arg \max_l \{\max\{|y_l|, |y'_l|\}\}$  then*

$$\text{upper}(p) \geq n^7 \|\mathbf{y} + \mathbf{y}'\|_\infty$$

assuming (10.8).

*Proof.* Recall that  $\|\mathbf{c} + \mathbf{c}'\|_2 \leq \text{gap}$  where  $\text{gap} = \frac{1}{10n^{10}} \min\{\|\mathbf{c}\|_2, \|\mathbf{c}'\|_2\}$ ,

$$\mathbf{c} = \mathbf{y} - \sum_i \alpha_i \mathbf{1}_i + \sum_j \beta_j \mathbf{1}_j + \sum_{(i,j)} \gamma_{ij} (\mathbf{1}_i - \mathbf{1}_j) \quad \text{and} \quad \mathbf{c}' = \mathbf{y}' - \sum_i \alpha'_i \mathbf{1}_i + \sum_j \beta'_j \mathbf{1}_j + \sum_{(i,j)} \gamma'_{ij} (\mathbf{1}_i - \mathbf{1}_j).$$

By (10.8) we know that  $\|\mathbf{c} - \mathbf{y}\|_2 \leq 4n^2 \text{gap} \leq \frac{4}{10n^8} \|\mathbf{c}\|_2$  and  $\|\mathbf{c}' - \mathbf{y}'\|_2 \leq 4n^2 \text{gap} \leq \frac{4}{10n^8} \|\mathbf{c}'\|_2$ . Consequently, by the triangle inequality we have that

$$\|\mathbf{y} + \mathbf{y}'\|_2 \leq \|\mathbf{c} + \mathbf{c}'\|_2 + \|\mathbf{c} - \mathbf{y}\|_2 + \|\mathbf{c}' - \mathbf{y}'\|_2 \leq 9n^2 \text{gap}$$

and

$$\|\mathbf{c}\|_2 \leq \|\mathbf{c} - \mathbf{y}\|_2 + \|\mathbf{y}\|_2 \leq \frac{4}{10n^8} \|\mathbf{c}\|_2 + \|\mathbf{y}\|_2 \quad \Rightarrow \quad \|\mathbf{c}\|_2 \leq 2\|\mathbf{y}\|_2$$

Similarly, we have that  $\|\mathbf{c}'\|_2 \leq 2\|\mathbf{y}'\|_2$ . Consequently since  $\text{gap} \leq \frac{1}{10n^{10}} \min\{\|\mathbf{c}\|_2, \|\mathbf{c}'\|_2\}$ , we have that

$$\|\mathbf{y} + \mathbf{y}'\|_2 \leq \frac{2}{n^8} \min\{\|\mathbf{y}\|_2, \|\mathbf{y}'\|_2\}$$

and thus, invoking Lemma 10.4.11 yields the result.  $\square$

We summarize the results in the lemma below.

**Corollary 10.4.16.** *Let  $P$  be the feasible region after running cutting plane on (10.1). Then one of the following holds:*

1. *We found a degenerate BFS and hence either  $\emptyset$  or  $V$  is a minimizer.*
2. *The integral points of  $P$  all lie on some hyperplane  $x_i = 0$ ,  $x_j = 1$  or  $x_i = x_j$  which we can find.*
3. *Let  $H$  be the collection of BFS's  $\mathbf{h}$  used to construct our separating hyperplanes for  $P$ . Then there is a convex combination  $\mathbf{y}$  of  $H$  such that  $n^4 |y_i| < \max_p \text{upper}(p)$  for all  $i$ .*

*Proof.* As mentioned before, (1) happens if some separating hyperplane is degenerate. We have (2) if one of the conditions in Lemma 10.4.14 holds. Otherwise,  $\mathbf{y} = \sum_{h \in H} \lambda_h \mathbf{h} + \sum_{h \in H} \lambda'_h \mathbf{h}$  is a candidate for Case 3 by Lemma 10.4.15.  $\square$

Let us revisit the conditions of Lemma 10.4.9 and explain that they are satisfied by Case 3 of the last lemma.

- $\mathbf{y}$  is a convex combination of at most  $O(n)$  BFS's. This holds in Case 3 since our current feasible region consists of only  $O(n)$  constraints thanks to the Cutting Plane method.
- Those BFS's must be consistent with every arc of  $A$ . This holds because Case 3 uses the BFS's for constructing our separating hyperplane. Our modified separation oracle guarantees that they are consistent with  $A$ .

Thus in Case 3 of the last corollary, Lemma 10.4.9 allows us to deduce a new constraint  $x_p \leq x_q$  for some  $q \notin R(p)$ .

### 10.4.3.3 Running Time

Here we bound the total running time of our algorithm and prove the following.

**Theorem 10.4.17.** *Our algorithm runs in time  $O(n^4 \log n \cdot EO + n^5 \log^{O(1)} n)$ .*

*Proof.* To avoid being repetitive, we appeal to Corollary 10.4.16. Each phase of cutting plane takes time  $O(n^2 \log n \cdot EO + n^3 \log^{O(1)} n)$  (Theorem 10.4.12 with  $\tau$  being a big constant. Given  $F$  and  $F'$  represented as a nonnegative combination of facets, we can check for the conditions in Lemma 10.4.14 in  $O(n)$  time as there are only this many facets of  $P$ . This settles Case 2 of Corollary 10.4.16. Finally, Lemma 10.4.9 tells us that we can find a new arc in  $O(n \cdot EO + n^2)$  time for Case 3 of Corollary 10.4.16. Our conclusion follows from the fact that we can get  $x_i = 0$ ,  $x_i = 1$ ,  $x_i = x_j$  at most  $n$  times and  $x_i \leq x_j$  at most  $O(n^2)$  times.  $\square$

### ■ 10.4.4 $\tilde{O}(n^3 \cdot EO + n^4)$ Time Algorithm

Here we show how to improve our running time for strongly polynomial SFM to  $\tilde{O}(n^3 \cdot EO + n^4)$ . Our algorithm can be viewed as an extension of the algorithm we presented in the previous Section 10.4.3. The main bottleneck of our previous algorithm was the time needed to identify a new arc, which cost us  $\tilde{O}(n^2 \cdot EO + n^3)$ . Here we show how to reduce our amortized cost for identifying a valid arc down to  $\tilde{O}(n \cdot EO + n^2)$  and thereby achieve our result.

The key observation we make to improve this running time is that our choice of  $p$  for adding an arc in the previous lemma can be relaxed.  $p$  actually need not be  $\arg \max_i \text{upper}(i)$ ; instead it is enough to have  $\text{upper}(p) > n^4 \max\{\alpha_i, \alpha'_i, \beta_j, \beta'_j, \gamma_{ij}, \gamma'_{ij}\}$ . For each such  $p$  a new constraint  $x_p \leq x_q$  can be identified via Lemma 10.4.9. So if there are many  $p$ 's satisfying this we will be able to obtain many new constraints and hence new valid arcs  $(p, q)$ .

On the other hand, the bound in Lemma 10.4.15 says that our point in the base polyhedron is small in *absolute* value. This is actually stronger than what we need in Lemma 10.4.9 which requires only its positive entries to be “small”. However as we saw in Lemma 10.4.10 we can generate a constraint of the form  $x_q \leq x_p$  whenever  $\text{lower}(p)$  is sufficiently negative.

Using this idea, we divide  $V$  into different buckets according to  $\text{upper}(p)$  and  $\text{lower}(p)$ . This will allow us to get a speedup for two reasons.

First, bucketing allows us to disregard unimportant elements of  $V$  during certain executions of our cutting plane method. If both  $\text{upper}(i)$  and  $\text{lower}(i)$  are small in absolute value, then  $i$  is essentially negligible because for a separating hyperplane  $\mathbf{h}^T \mathbf{x} \leq \hat{f}(\bar{\mathbf{x}})$ , any  $h_i \in [\text{lower}(i), \text{upper}(i)]$  small in absolute value would not really make a difference. We can then run our cutting plane algorithm only on those non-negligible  $i$ 's, thereby reducing our time complexity. Of course, whether  $h_i$  is small is something relative. This suggests that partitioning the ground set by the relative size of  $\text{upper}(i)$  and  $\text{lower}(i)$  is a good idea.

Second, bucketing allows us to ensure that we can always add an arc for many edges simultaneously. Recall that we remarked that all we want is  $n^{O(1)}|y_i| \leq \mathbf{upper}(p)$  for some  $\mathbf{y}$  in the base polyhedron. This would be sufficient to identify a new valid arc  $(p, q)$ . Now if the marginal differences  $\mathbf{upper}(p)$  and  $\mathbf{upper}(p')$  are close in value, knowing  $n^{O(1)}|y_i| \leq \mathbf{upper}(p)$  would effectively give us the same for  $p'$  for free. This suggests that elements with similar marginal differences should be grouped together.

The remainder of this section simply formalizes these ideas. In Section 10.4.4.1 we discuss how we partition the ground set  $V$ . In Section 10.4.4.2, we present our cutting plane method on a subset of the coordinates. Then in Section 10.4.4.3 we show how we find new arcs. Finally, in Section 10.4.4.4 we put all of this together to achieve our desired running time.

#### 10.4.4.1 Partitioning Ground Set into Buckets

We partition the ground set  $V$  into different buckets according to the values of  $\mathbf{upper}(i)$  and  $\mathbf{lower}(i)$ . This is reminiscent to Iwata-Orlin's algorithm [125] which considers elements with big  $\mathbf{upper}(i)$ . However they did not need to do bucketing by size or to consider  $\mathbf{lower}(i)$ , whereas these seem necessary for our algorithm.

Let  $N = \max_i \{\max\{\mathbf{upper}(i), -\mathbf{lower}(i)\}\}$  be the largest marginal difference in absolute value. By Lemma (10.4.8),  $N \geq 0$ . We partition our ground set  $V$  as follows:

$$B_1 = \{i : \mathbf{upper}(i) \geq N/n^{10} \text{ or } \mathbf{lower}(i) \leq -N/n^{10}\}$$

$$B_k = \{i \notin B_1 \cup \dots \cup B_{k-1} : \begin{aligned} &N/n^{10k} \leq \mathbf{upper}(i) < N/n^{10(k-1)} \\ &\text{or } -N/n^{10(k-1)} < \mathbf{lower}(i) \leq -N/n^{10k} \}, \quad k \geq 2 \end{aligned}$$

We call  $B_k$  *buckets*. Our buckets group elements by the values of  $\mathbf{upper}(i)$  and  $\mathbf{lower}(i)$  at  $1/n^{10}$  "precision". There are two cases.

- Case 1: the number of buckets is at most  $\log n^2$ , in which case  $\mathbf{upper}(i) > N/n^{O(\log n)}$  or  $\mathbf{lower}(i) < -N/n^{O(\log n)}$  for all  $i$ .
- Case 2: there is some  $k$  for which  $|B_1 \cup \dots \cup B_k| \geq |B_{k+1}|$ .

This is because if there is no such  $k$  in Case 2, then by induction each bucket  $B_{k+1}$  has at least  $2^k|B_1| \geq 2^k$  elements and hence  $k \leq \log n$ .

Case 1 is easier to handle, and is in fact a special case of Case 2. We first informally sketch the treatment for Case 1 which should shed some light into how we deal with Case 2.

We run Cutting Plane for  $O(n \log^2 n)$  iterations (i.e.  $\tau = \Theta(\log n)$ ). By Theorem 10.4.12, our feasible region  $P$  would be sandwiched by a pair of approximately parallel supporting hyperplanes of width at most  $1/n^{\Theta(\log n)}$ . Now proceeding as in the last section, we would be able to find some  $\mathbf{y}$  in the base polyhedron and some element  $p$  such that  $n^{\Theta(\log n)}|y_i| \leq \mathbf{upper}(p)$ . This gives

$$n^{\Theta(\log n)}|y_i| \leq \frac{\mathbf{upper}(p)}{n^{\Theta(\log n)}} \leq \frac{N}{n^{\Theta(\log n)}}.$$

Since  $\mathbf{upper}(i) > N/n^{\Theta(\log n)}$  or  $\mathbf{lower}(i) < -N/n^{\Theta(\log n)}$  for all  $i$  in Case 1, we can then conclude that some valid arc  $(i, q)$  or  $(q, i)$  can be added for every  $i$ . Thus we add  $n/2$  arcs simultaneously in one phase of the algorithm at the expense of blowing up the runtime by  $O(\log n)$ . This saves a factor of  $n/\log n$  from our runtime in the last section, and the amortized cost for an arc would then be  $\tilde{O}(n \cdot \text{EO} + n^2)$ .

---

<sup>2</sup>More precisely,  $B_k = \emptyset$  for  $k > \lceil \log n \rceil$ .

On the other hand, in Case 2 we have a “trough” at  $B_{k+1}$ . Roughly speaking, this trough is useful for acting as a soft boundary between  $B_1 \cup \dots \cup B_k$  and  $\bigcup_{l \geq k+2} B_l$ . Recall that we are able to “ignore”  $\bigcup_{l \geq k+2} B_l$  because their  $h_i$  is relatively small in absolute value. In particular, we know that for any  $p \in B_1 \cup \dots \cup B_k$  and  $i \in B_l$ , where  $l \geq k+2$ ,

$$\max\{\text{upper}(p), -\text{lower}(p)\} \geq n^{10} \max\{\text{upper}(i), -\text{lower}(i)\}.$$

This is possible because  $B_{k+1}$ , which is sandwiched in between, acts like a shield preventing  $B_l$  to “mess with”  $B_1 \cup \dots \cup B_k$ . This property comes at the expense of sacrificing  $B_{k+1}$  which must confront  $B_l$ .

Furthermore, we require that  $|B_1 \cup \dots \cup B_k| \geq |B_{k+1}|$ , and run Cutting Plane on  $B = (B_1 \cup \dots \cup B_k) \cup B_{k+1}$ . If  $|B_{k+1}| \gg |B_1 \cup \dots \cup B_k|$ , our effort would mostly be wasted on  $B_{k+1}$  which is sacrificed, and the amortized time complexity for  $B_1 \cup \dots \cup B_k$  would then be large.

Before discussing the algorithm for Case 2, we need some preparatory work.

#### 10.4.4.2 Separating Hyperplane: Project and Lift

Our speedup is achieved by running our cutting plane method on the projection of our feasible region onto  $B := (B_1 \cup \dots \cup B_k) \cup B_{k+1}$ . More precisely, we start by running our cutting plane on  $P^B = \{\mathbf{x} \in \mathbb{R}^B : \exists \mathbf{x}' \in \mathbb{R}^{\bar{B}} \text{ s.t. } (\mathbf{x}, \mathbf{x}') \text{ satisfies (10.1)}\}$ , which has a lower dimension. However, to do this, we need to specify a separation oracle for  $P^B$ . Here we make one of the most natural choices.

We begin by making an apparently immaterial change to our set of arcs  $A$ . Let us take the *transitive closure* of  $A$  by adding the arc  $(i, j)$  whenever there is a path from  $i$  to  $j$ . Clearly this would not change our ring family as a path from  $i$  to  $j$  implies  $j \in R(i)$ . Roughly speaking, we do this to handle pathological cases such as  $(i, k), (k, j) \in A, (i, j) \notin A$  and  $i, j \in B, k \notin B$ . Without introducing the arc  $(i, j)$ , we risk confusing a solution containing  $i$  but not  $j$  as feasible since we are restricting our attention to  $B$  and ignoring  $k \notin B$ .

**Definition 10.4.18.** Given a digraph  $D = (V, A)$ , the transitive closure of  $A$  is the set of arcs  $(i, j)$  for which there is a directed path from  $i$  to  $j$ . We say that  $A$  is *complete* if it is equal to its transitive closure.

Given  $\bar{x} \in [0, 1]^B$ , we define the completion of  $\bar{x}$  with respect to  $A$  as follows.

**Definition 10.4.19.** Given  $\bar{x} \in [0, 1]^B$  and a set of arcs  $A$ ,  $x^C \in [0, 1]^n$  is a completion of  $\bar{x}$  if  $x_B^C = \bar{x}$  and  $x_i^C \leq x_j^C$  for every  $(i, j) \in A$ . Here  $x_B^C$  denotes the restriction of  $x^C$  to  $B$ .

**Lemma 10.4.20.** Given  $\bar{x} \in [0, 1]^B$  and a complete set of arcs  $A$ , there is a completion of  $\bar{x}$  if  $\bar{x}_i \leq \bar{x}_j$  for every  $(i, j) \in A \cap (B \times B)$ . Moreover, it can be computed in  $O(n^2)$  time.

*Proof.* We set  $x_B^C = \bar{x}$ . For  $i \notin B$ , we set

$$x_i^C = \begin{cases} 1 & \text{if } \nexists j \in B \text{ s.t. } (i, j) \in A \\ \min_{(i, j) \in A, j \in B} x_j^C & \text{otherwise} \end{cases}$$

One may verify that  $x^C$  satisfies our requirement as  $A$  is complete. Computing each  $x_i^C$  takes  $O(n)$  time. Since  $|V \setminus B| = |\bar{B}| \leq n$ , computing the whole  $x^C$  takes  $O(n^2)$  time.  $\square$

This notion of completion is needed since our original separation oracle requires a full dimensional input  $\bar{x}$ . Now that  $\bar{x} \in \mathbb{R}^B$ , we need a way of extending it to  $\mathbb{R}^n$  while retaining the crucial property that  $\mathbf{h}$  is consistent with every arc in  $A$ .

**Algorithm 26:** Projected Separation Oracle

---

**Input:**  $\bar{x} \in \mathbb{R}^B$  and a complete set of arcs  $A$

**if**  $\bar{x}_i < 0$  **for some**  $i \in B$  **then**

  | **Output:**  $x_i \geq 0$

**else if**  $\bar{x}_j > 1$  **for some**  $j \in B$  **then**

  | **Output:**  $x_j \leq 1$

**else if**  $\bar{x}_i > \bar{x}_j$  **for some**  $(i, j) \in A \cap B^2$  **then**

  | **Output:**  $x_i \leq x_j$

**else**

  | Let  $x^C \in \mathbb{R}^n$  be a completion of  $\bar{x}$

  | Let  $i_1, \dots, i_n$  be a permutation of  $V$  such that  $x_{i_1}^C \geq \dots \geq x_{i_n}^C$  and for all  $(i, j) \in A$ ,  $j$  precedes  $i$  in  $i_1, \dots, i_n$ .

  | **Output:**  $\mathbf{h}_B^T \mathbf{x}_B = \sum_{i \in B} h_i x_i \leq \sum_{i \in B} h_i \bar{x}_i$ , where  $\mathbf{h}$  is the BFS defined by the permutation  $i_1, \dots, i_n$ .

---

Note that the runtime is still  $O(n \cdot \text{EO} + n^2 \log^{O(1)} n)$  as  $x^C$  can be computed in  $O(n^2)$  time by the last lemma.

We reckon that the hyperplane  $\mathbf{h}_B^T \mathbf{x}_B \leq \sum_{i \in B} h_i \bar{x}_i$  returned by the oracle is *not* a valid separating hyperplane (i.e. it may cut out the minimizers). Nevertheless, we will show that it is a decent “proxy” to the true separating hyperplane  $\mathbf{h}^T \mathbf{x} \leq \hat{f}(x^C) = \sum_{i \in V} h_i x_i^C$  and is good enough to serve our purpose of sandwiching the remaining feasible region in a small strip. To get a glimpse, note that the terms missing  $\mathbf{h}_B^T \mathbf{x}_B \leq \sum_{i \in B} h_i \bar{x}_i$  all involve  $h_i$  for  $i \notin B$ , which is “negligible” compared to  $B_1 \cup \dots \cup B_k$ .

One may try to make  $\mathbf{h}_B^T \mathbf{x}_B \leq \sum_{i \in B} h_i \bar{x}_i$  valid, say, by  $\mathbf{h}_B^T \mathbf{x}_B \leq \sum_{i \in B} h_i \bar{x}_i + \sum_{i \notin B} |h_i|$ . The problem is that such hyperplanes would not be separating for  $\bar{x}$  anymore as  $\mathbf{h}_B^T \bar{x} = \sum_{i \in B} h_i \bar{x}_i < \sum_{i \in B} h_i \bar{x}_i + \sum_{i \notin B} |h_i|$ . Consequently, we lose the width (or volume) guarantee of our cutting plane algorithm. Although this seems problematic, it is actually still possible to show a guarantee sufficient for our purpose as  $\sum_{i \notin B} |h_i|$  is relatively small. We leave it as a nontrivial exercise to interested readers.

In conclusion, it seems that one cannot have the best of both worlds: the hyperplane returned by the oracle cannot be simultaneously valid and separating.

**Algorithm**

We take  $k$  to be the first for which  $|B_1 \cup \dots \cup B_k| \geq |B_{k+1}|$ , i.e.  $|B_1 \cup \dots \cup B_l| < |B_{l+1}|$  for  $l \leq k-1$ . Thus  $k \leq \log n$ . Let  $b = |B|$ , and so  $|B_1 \cup \dots \cup B_k| \geq b/2$ . Case 1 is a special case by taking  $B = V$ .

Our algorithm is summarized below. Here  $A$  is always complete as  $A$  is replaced its transitive closure whenever a new valid arc is added.

1. Run Cutting Plane on  $P^B = \{x \in \mathbb{R}^B : \exists x' \in \mathbb{R}^{\bar{B}} \text{ s.t. } (x, x') \text{ satisfies (10.1)}\}$  with the new projected separation oracle.
2. Identify a pair of “narrow” approximately parallel supporting hyperplanes.
3. Deduce from the hyperplanes certain new constraints of the forms  $x_i = 0, x_j = 1, x_i = x_j$  or  $x_i \leq x_j$  by lifting separating hyperplanes back to  $\mathbb{R}^n$
4. Consolidate  $A$  and  $f$ . If some  $x_i \leq x_j$  added, replace  $A$  by its transitive closure.
5. Repeat Step 1 with updated  $A$  and  $f$ . (Any previously found separating hyperplanes are discarded.)

The minimizer can be constructed by unraveling the recursion.

First of all, to be able to run Cutting Plane on  $P^B$  we must come up with a polyhedral description of  $P^B$  which consists of just the constraints involving  $B$ . This is shown in the next lemma.

**Lemma 10.4.21.** *Let  $P^B = \{\mathbf{x} \in \mathbb{R}^B : \exists \mathbf{x}' \in \mathbb{R}^{\bar{B}} \text{ s.t. } (\mathbf{x}, \mathbf{x}') \text{ satisfies (10.1)}\}$ . Then*

$$P^B = \{\mathbf{x} \in \mathbb{R}^B : 0 \leq \mathbf{x} \leq 1, x_i \leq x_j \forall (i, j) \in A \cap (B \times B)\}$$

*Proof.* It is clear that  $P^B \subseteq \{\mathbf{x} \in \mathbb{R}^B : 0 \leq \mathbf{x} \leq 1, x_i \leq x_j \forall (i, j) \in A \cap (B \times B)\}$  as the constraints  $0 \leq x \leq 1, x_i \leq x_j \forall (i, j) \in A \cap (B \times B)$  all appear in (10.1).

Conversely, for any  $\mathbf{x} \in \mathbb{R}^B$  satisfying  $0 \leq \mathbf{x} \leq 1, x_i \leq x_j \forall (i, j) \in A \cap (B \times B)$ , we know there is some completion  $x^C$  of  $\mathbf{x}$  by Lemma 10.4.20 as  $A$  is complete. Now  $x^C$  satisfies (10.1) by definition, and hence  $\mathbf{x} \in P^B$ .  $\square$

The only place where we have really changed the algorithm is Step (3).

### 10.4.4.3 Deducing New Constraints $x_i = 0, x_j = 1, x_i = x_j$ or $x_i \leq x_j$

Our method will deduce one of the following:

- $x_i = 0, x_j = 1$  or  $x_i = x_j$
- for each  $p \in B_1 \cup \dots \cup B_k, x_p \leq x_q$  for some  $q \notin R(p)$  or  $x_p \geq x_q$  for some  $q \notin Q(p)$

Our argument is very similar to the last section's. Roughly speaking, it is the same argument but with "noise" introduced by  $i \notin B$ . We use extensively the notations from the last section.

Our main tool is again Theorem 10.4.12. Note that  $n$  should be replaced by  $b$  in the Theorem statement. We invoke it with  $\tau = k \log_b n = O(\log^2 n)$  (using  $k \leq \log n$ ) to get a width of  $1/b^{\Theta(\tau)} = 1/n^{\Theta(k)}$ . This takes time at most  $O(bn \log^2 n \cdot \text{EO} + bn^2 \log^{O(1)} n)$ . Again, this is intuitively clear as we run it for  $O(kb \log n)$  iterations, each of which takes time  $O(n \cdot \text{EO} + n^2 \log^{O(1)} n)$ .

After each phase of (roughly  $O(kb \log n)$  iterations) of Cutting Plane,  $P^B$  is sandwiched between a pair of approximately parallel supporting hyperplanes  $F$  and  $F'$  which have width  $1/n^{20k}$ . Let  $F$  and  $F'$  be

$$\mathbf{c}^T \mathbf{x}_B = \sum_{i \in B} c_i x_i \leq M, \quad \mathbf{c}'^T \mathbf{x}_B = \sum_{i \in B} c'_i x_i \leq M',$$

such that

$$|M + M'|, \|\mathbf{c} + \mathbf{c}'\|_2 \leq \text{gap}, \quad \text{where } \text{gap} = \frac{1}{n^{20k}} \min\{\|\mathbf{c}\|_2, \|\mathbf{c}'\|_2\}.$$

The rest of this section presents an execution of the ideas discussed above. All of our work is basically geared towards bringing the amortized cost for identifying a valid arc down to  $\tilde{O}(n \cdot \text{EO} + n^2)$ . Again, we can write these two constraints as a nonnegative combination. Here  $\bar{x}_h^C$  is the completion of the point  $\bar{x}_h$  used to construct  $\mathbf{h}_B^T \mathbf{x}_B \leq \mathbf{h}_B^T (\bar{x}_h^C)_B$ . (Recall that  $(\bar{x}_h^C)_B$  is the restriction of  $\bar{x}_h^C$  to  $B$ .)

$$\mathbf{c}^T \mathbf{x}_B = - \sum_{i \in B} \alpha_i x_i + \sum_{j \in B} \beta_j x_j + \sum_{(i,j) \in A \cap B^2} \gamma_{ij} (x_i - x_j) + \sum_{h \in H} \lambda_h \mathbf{h}_B^T \mathbf{x}_B \quad \text{and} \quad M = \sum_{j \in B} \beta_j + \sum_{h \in H} \lambda_h \mathbf{h}_B^T (\bar{x}_h^C)_B.$$

$$\mathbf{c}'^T \mathbf{x}_B = - \sum_{i \in B} \alpha'_i x_i + \sum_{j \in B} \beta'_j x_j + \sum_{(i,j) \in A \cap B^2} \gamma'_{ij} (x_i - x_j) + \sum_{h \in H} \lambda'_h \mathbf{h}_B^T \mathbf{x}_B \quad \text{and} \quad M' = \sum_{j \in B} \beta'_j + \sum_{h \in H} \lambda'_h \mathbf{h}_B^T (\bar{x}_h^C)_B.$$



As we have discussed, the problem is that the separating hyperplanes  $\mathbf{h}_B^T \mathbf{x}_B \leq \mathbf{h}_B^T (\bar{x}_h^C)_B$  are not actually valid. We can, however, recover their valid counterpart by lifting them back to  $\mathbf{h}^T \mathbf{x} \leq \mathbf{h}^T \bar{x}_h^C$ . The hope is that  $\mathbf{h}_B^T \mathbf{x}_B \leq \mathbf{h}_B^T (\bar{x}_h^C)_B$  and  $\mathbf{h}^T \mathbf{x} \leq \mathbf{h}^T \bar{x}_h^C$  are not too different so that the arguments will still go through. We show that this is indeed the case.

Again, we scale  $c, c', \alpha, \alpha', \beta, \beta', \gamma, \gamma', \lambda, \lambda'$  so that

$$\sum_{h \in H} (\lambda_h + \lambda'_h) = 1.$$

By adding all the constituent separating hyperplane inequalities, we get

$$\sum_{h \in H} \lambda_h \mathbf{h}^T \mathbf{x} + \sum_{h \in H} \lambda'_h \mathbf{h}^T \mathbf{x} \leq \sum_{h \in H} \lambda_h \mathbf{h}^T \bar{x}_h^C + \sum_{h \in H} \lambda'_h \mathbf{h}^T \bar{x}_h^C$$

Let

$$LHS \stackrel{\text{def}}{=} \sum \alpha_i x_i + \sum \alpha'_i x_i - \sum \beta_j x_j - \sum \beta'_j x_j + \sum \gamma_{ij} (x_j - x_i) + \sum \gamma'_{ij} (x_j - x_i).$$

Here we know that

$$\begin{aligned} \sum_{h \in H} \lambda_h \mathbf{h}^T \mathbf{x} + \sum_{h \in H} \lambda'_h \mathbf{h}^T \mathbf{x} &= LHS + (\mathbf{c} + \mathbf{c}')^T \mathbf{x}_B + \sum_{h \in H} \lambda_h \mathbf{h}_B^T \mathbf{x}_B + \sum_{h \in H} \lambda'_h \mathbf{h}_B^T \mathbf{x}_B \\ \sum_{h \in H} \lambda_h \mathbf{h}^T \bar{x}_h^C + \sum_{h \in H} \lambda'_h \mathbf{h}^T \bar{x}_h^C &= (M + M') + \sum_{h \in H} \lambda_h \mathbf{h}_B^T (\bar{x}_h^C)_B + \sum_{h \in H} \lambda'_h \mathbf{h}_B^T (\bar{x}_h^C)_B - \sum \beta_j - \sum \beta'_j \end{aligned}$$

Combining all yields

$$\begin{aligned} &LHS + (\mathbf{c} + \mathbf{c}')^T \mathbf{x}_B + \sum_{h \in H} \lambda_h \mathbf{h}_B^T \mathbf{x}_B + \sum_{h \in H} \lambda'_h \mathbf{h}_B^T \mathbf{x}_B \\ &\leq (M + M') + \sum_{h \in H} \lambda_h \mathbf{h}_B^T (\bar{x}_h^C)_B + \sum_{h \in H} \lambda'_h \mathbf{h}_B^T (\bar{x}_h^C)_B - \sum \beta_j - \sum \beta'_j. \end{aligned}$$

Here  $(\mathbf{c} + \mathbf{c}')^T \mathbf{x}_B$  can be bounded as before:  $(\mathbf{c} + \mathbf{c}')^T \mathbf{x}_B \geq -\sqrt{n} \|\mathbf{c} + \mathbf{c}'\|_2 \geq -\sqrt{n} \text{gap}$ . Since  $M + M' \leq \text{gap}$ , We then obtain

$$LHS + \sum_{h \in H} \lambda_h \mathbf{h}_B^T \mathbf{x}_B + \sum_{h \in H} \lambda'_h \mathbf{h}_B^T \mathbf{x}_B \leq 2\sqrt{n} \text{gap} + \sum_{h \in H} \lambda_h \mathbf{h}_B^T (\bar{x}_h^C)_B + \sum_{h \in H} \lambda'_h \mathbf{h}_B^T (\bar{x}_h^C)_B - \sum \beta_j - \sum \beta'_j$$

We should expect the contribution from  $\mathbf{h}_B$  to be small as  $h_i$  for  $i \notin B$  is small compared to  $B_1 \cup \dots \cup B_k$ . We formalize our argument in the next two lemmas.

**Lemma 10.4.22.** *We have  $\sum_{h \in H} \lambda_h \mathbf{h}_B^T (\bar{x}_h^C)_B + \sum_{h \in H} \lambda'_h \mathbf{h}_B^T (\bar{x}_h^C)_B \leq N/n^{10(k+1)-1}$ .*

*Proof.* We bound each component of  $\sum_{h \in H} \lambda_h \mathbf{h}_B^T (\bar{x}_h^C)_B + \sum_{h \in H} \lambda'_h \mathbf{h}_B^T (\bar{x}_h^C)_B$ . For  $i \in \bar{B}$ , we have  $\text{upper}(i) \leq N/n^{10(k+1)}$ . By Lemma 10.4.8  $h_i \leq \text{upper}(i)$ . Therefore,

$$\sum_{h \in H} \lambda_h \mathbf{h}_i^T (\bar{x}_h^C)_i + \sum_{h \in H} \lambda'_h \mathbf{h}_i^T (\bar{x}_h^C)_i \leq \left( \sum_{h \in H} \lambda_h + \sum_{h \in H} \lambda'_h \right) N/n^{10(k+1)} = N/n^{10(k+1)}.$$

Our result then follows since

$$\sum_{h \in H} \lambda_h \mathbf{h}_B^T (\bar{x}_h^C)_B + \sum_{h \in H} \lambda'_h \mathbf{h}_B^T (\bar{x}_h^C)_B = \sum_{i \in \bar{B}} \left( \sum_{h \in H} \lambda_h \mathbf{h}_i^T (\bar{x}_h^C)_i + \sum_{h \in H} \lambda'_h \mathbf{h}_i^T (\bar{x}_h^C)_i \right).$$

□

**Lemma 10.4.23.** We have  $\sum_{h \in H} \lambda_h \mathbf{h}_B^T \mathbf{x}_B + \sum_{h \in H} \lambda'_h \mathbf{h}_B^T \mathbf{x}_B \geq -N/n^{10(k+1)-1}$ .

*Proof.* The proof is almost identical to the last lemma except that we use  $h_i \geq \text{lower}(i)$  instead of  $h_i \leq \text{upper}(i)$ , and  $\text{lower}(i) \geq -N/n^{10(k+1)}$ . □

The two lemmas above imply that

$$LHS \leq 2\sqrt{n}\text{gap} - \sum \beta_j - \sum \beta_j + 2N/n^{10(k+1)-1} = \text{gap}' - \sum \beta_j - \sum \beta_j$$

where  $\text{gap}' = 2\sqrt{n}\text{gap} + 2N/n^{10(k+1)-1}$ .

**Lemma 10.4.24.** Suppose  $x$  satisfies (10.1) and  $LHS \leq \text{gap}' - \sum \beta_j - \sum \beta'_j$  with  $\alpha_i, \beta_j, \gamma_{ij}, \alpha'_i, \beta'_j, \gamma'_{ij} \geq 0$ .

1. If  $\alpha_i > \text{gap}'$  or  $\alpha'_i > \text{gap}'$ , then  $x_i < 1$ .
2. If  $\beta_j > \text{gap}'$  or  $\beta'_j > \text{gap}'$ , then  $x_j > 0$ .
3. If  $\gamma_{ij} > \text{gap}'$  or  $\gamma'_{ij} > \text{gap}'$ , then  $0 \leq x_j - x_i < 1$ .

*Proof.* The proof is exactly the same as Lemma 10.4.14 with  $2\sqrt{n}\text{gap}$  replaced by  $\text{gap}'$ . □

From now on we may assume that

$$\max\{\alpha_i, \alpha'_i, \beta_j, \beta'_j, \gamma_{ij}, \gamma'_{ij}\} \leq \text{gap}'. \quad (10.9)$$

**Lemma 10.4.25.** Let  $\mathbf{y} \stackrel{\text{def}}{=} \sum_{h \in H} \lambda_h \mathbf{h}$  and  $\mathbf{y}' \stackrel{\text{def}}{=} \sum_{h \in H} \lambda'_h \mathbf{h}$  and let  $p \in \arg \max_{l \in B} \{\max\{|y_l|, |y'_l|\}\}$  then

$$N \geq n^{10k+6} \|\mathbf{y}_B + \mathbf{y}'_B\|_\infty$$

assuming (10.9).

*Proof.* Recall that  $\|\mathbf{c} + \mathbf{c}'\|_2 \leq \text{gap} < \text{gap}'$  where  $\text{gap} = \frac{1}{n^{20k}} \min\{\|\mathbf{c}\|_2, \|\mathbf{c}'\|_2\}$  and  $\text{gap}' = 2\sqrt{n}\text{gap} + 2N/n^{10(k+1)-1}$ . Now there are two cases.

**Case 1:**  $2\sqrt{n}\text{gap} \geq 2N/n^{10(k+1)-1}$ . Then  $\text{gap}' \leq 4\sqrt{n}\text{gap}$  and we follow the same proof of Lemma 10.4.15. We have

$$\mathbf{c} = \mathbf{y}_B - \sum_i \alpha_i \mathbf{1}_i + \sum_j \beta_j \mathbf{1}_j + \sum_{(i,j)} \gamma_{ij} (\mathbf{1}_i - \mathbf{1}_j) \quad \text{and} \quad \mathbf{c}' = \mathbf{y}'_B - \sum_i \alpha'_i \mathbf{1}_i + \sum_j \beta'_j \mathbf{1}_j + \sum_{(i,j)} \gamma'_{ij} (\mathbf{1}_i - \mathbf{1}_j).$$

By (10.9) we know that  $\|\mathbf{c} - \mathbf{y}_B\|_2 \leq 4n^2 \text{gap}' \leq \frac{1}{n^{17k}} \|\mathbf{c}\|_2$  and  $\|\mathbf{c}' - \mathbf{y}'_B\|_2 \leq 4n^2 \text{gap}' \leq \frac{1}{n^{17k}} \|\mathbf{c}\|_2$ . Consequently, by the triangle inequality we have that

$$\|\mathbf{y}_B + \mathbf{y}'_B\|_2 \leq \|\mathbf{c} + \mathbf{c}'\|_2 + \|\mathbf{c} - \mathbf{y}_B\|_2 + \|\mathbf{c}' - \mathbf{y}'_B\|_2 \leq 9n^2 \text{gap}'$$

and

$$\|\mathbf{c}\|_2 \leq \|\mathbf{c} - \mathbf{y}_B\|_2 + \|\mathbf{y}_B\|_2 \leq \frac{1}{n^{17k}} \|\mathbf{c}\|_2 + \|\mathbf{y}_B\|_2 \quad \Rightarrow \quad \|\mathbf{c}\|_2 \leq 2\|\mathbf{y}_B\|_2$$

Similarly, we have that  $\|\mathbf{c}'\|_2 \leq 2\|\mathbf{y}'_B\|_2$ . Consequently since  $\text{gap}' \leq \frac{1}{n^{19k}} \min\{\|\mathbf{c}\|_2, \|\mathbf{c}'\|_2\}$ , we have that

$$\|\mathbf{y}_B + \mathbf{y}'_B\|_2 \leq \frac{18}{n^{17k}} \min\{\|\mathbf{y}_B\|_2, \|\mathbf{y}'_B\|_2\}$$

and thus, invoking Lemma 10.4.11 yields  $N \geq \text{upper}(p) \geq n^{16k} \|\mathbf{y}_B + \mathbf{y}'_B\|_\infty$ , as desired.

**Case 2:**  $2\sqrt{n}\text{gap} < 2N/n^{10(k+1)-1}$ . Then for any  $i \in B$ ,  $|c_i + c'_i| \leq \|\mathbf{c} + \mathbf{c}'\|_2 \leq \text{gap} < 2N/n^{10(k+1)-1}$ . Since

$$\mathbf{y}_B + \mathbf{y}'_B = (\mathbf{c} + \mathbf{c}') + \sum_i \alpha_i \mathbf{1}_i - \sum_j \beta_j \mathbf{1}_j - \sum_{(i,j)} \gamma_{ij} (\mathbf{1}_i - \mathbf{1}_j) + \sum_i \alpha'_i \mathbf{1}_i - \sum_j \beta'_j \mathbf{1}_j - \sum_{(i,j)} \gamma'_{ij} (\mathbf{1}_i - \mathbf{1}_j)$$

we have

$$\|\mathbf{y}_B + \mathbf{y}'_B\|_\infty \leq 2N/n^{10(k+1)-1} + 2n^{1.5}\text{gap}' \leq N/n^{10k+7}.$$

□

**Corollary 10.4.26.** *Let  $P$  be the feasible region after running Cutting Plane on (10.1) with the projected separation oracle. Then one of the following holds:*

1. We found a BFS  $\mathbf{h}$  with  $\mathbf{h}_B = 0$ .
2. The integral points of  $P$  all lie on some hyperplane  $x_i = 0, x_j = 1$  or  $x_i = x_j$ .
3. Let  $H$  be the collection of BFS's  $\mathbf{h}$  used to construct our separating hyperplanes for  $P$ . Then there is a convex combination  $\mathbf{y}$  of  $H$  such that for  $p \in B_1 \cup \dots \cup B_k$ , we have  $n^4|y_i| < \text{upper}(p)$  or  $\text{lower}(p) < -n^4|y_i|$  for all  $i$ .

*Proof.* As mentioned before, (1) happens if some separating hyperplane satisfies  $\mathbf{h}_B = 0$  when running cutting plane on the non-negligible coordinates. We have (2) if some condition in Lemma 10.4.24 holds. Otherwise, we claim  $\mathbf{y} = \sum_{\mathbf{h}} \lambda_{\mathbf{h}} \mathbf{h} + \sum_{\mathbf{h}'} \lambda'_{\mathbf{h}'} \mathbf{h}'$  is a candidate for Case 3.  $\mathbf{y}$  is a convex combination of BFS and by Lemma 10.4.25, for the big elements  $i \in B$  we have

$$|y_i| \leq N/n^{10k+6} \leq \frac{1}{n^4} \max\{\text{upper}(p), -\text{lower}(p)\}.$$

where the last inequality holds since for  $p \in B_1 \cup \dots \cup B_k$ ,  $\max\{\text{upper}(p), -\text{lower}(p)\} \geq N/n^{10k}$ .

On the other hand, for the small elements  $i \notin B$ ,  $|y_i| \leq N/n^{10(k+1)} \leq \frac{1}{n^4} \max\{\text{upper}(p), -\text{lower}(p)\}$  as desired. □

The gap is then smaller enough to add an arc for each  $p \in B_1 \cup \dots \cup B_k$  by Lemmas 10.4.9 and 10.4.10. Therefore we can add a total of  $|B_1 \cup \dots \cup B_k|/2 \geq b/4$  arcs with roughly  $O(kb \log n) = \tilde{O}(b)$  iterations of Cutting Plane, each of which takes  $\tilde{O}(n \cdot \text{EO} + n^2)$ . That is, the amortized cost for each arc is  $\tilde{O}(n \cdot \text{EO} + n^2)$ . We give a more formal time analysis in below but it should be somewhat clear why we have the desired time complexity.

**Lemma 10.4.27.** *Suppose there is a convex combination  $\mathbf{y}$  of  $H$  such that for  $p \in B_1 \cup \dots \cup B_k$ , we have  $n^4|y_i| < \text{upper}(p)$  or  $\text{lower}(p) < -n^4|y_i|$  for all  $i$ . Then we can identify at least  $b/4$  new valid arcs.*

*Proof.* We have  $|H| = O(n)$  since  $H$  is the set of BFS's used for the constraints of  $P$  which has  $O(n)$  constraints. By Lemmas 10.4.9 and 10.4.10, for  $p \in B_1 \cup \dots \cup B_k$  we can add a new valid arc  $(p, q)$  or  $(q, p)$ . However note that a new arc  $(p_1, p_2)$  may added twice by both  $p_1$  and  $p_2$ . Therefore the total number of new arcs is only at least  $|B_1 \cup \dots \cup B_k|/2 \geq b/4$ . □

#### 10.4.4.4 Running Time

Not much changes to the previous runtime analysis are needed. To avoid repetition, various details already present in the corresponding part of the last section are omitted. Recall  $k \leq \log n$ , and of course,  $b \leq n$ .

For each (roughly)  $O(kb \log n)$  iterations of Cutting Plane we either get  $x_i = 0, x_i = 1, x_i = x_j$  or  $b/4 x_i \leq x_j$ 's. The former can happen at most  $n$  times while in the latter case, the amortized cost of each arc is  $O(k \log n)$  iterations of Cutting Plane. In the worst case the overall number of iterations required is  $\tilde{O}(n^2)$ . Thus our algorithm has a runtime of  $\tilde{O}(n^3 \cdot \text{EO} + n^4)$  since each iteration is  $\tilde{O}(n \cdot \text{EO} + n^2)$  as shown below.

**Theorem 10.4.28.** *Our algorithm runs in time  $O(n^3 \log^2 n \cdot \text{EO} + n^4 \log^{O(1)} n)$ .*

*Proof.* We use Corollary 10.4.26. First we note that Case 1 can actually be integrated into Case 3 since  $\max\{\text{upper}(p), -\text{lower}(p)\} \geq N/n^{10k} = n^{10}N/n^{10(k+1)} \geq h_i$  for  $i \notin B$ .

As we have argued in the beginning of the last section, Theorem 10.4.12 with  $\tau = k \log_b n$  implies that the runtime for each phase is  $O(bn \log^2 n \cdot \text{EO} + bn^2 \log^{O(1)} n)$ . In each phase we either get  $x_i = 0, x_i = 1, x_i = x_j$  (Case 2) or  $b/4 x_i \leq x_j$ 's (Case 3), the latter of which follows from Corollary 10.4.26 and Lemma 10.4.27.

Case 2 can only happen  $n$  times. Thus the total cost is at most  $O(n^3 \log^2 n \cdot \text{EO} + n^4 \log^{O(1)} n)$ . The overhead cost is also small. Similar to before, given  $F$  and  $F'$  represented as a nonnegative combination of facets, we can check for the conditions in Lemma 10.4.24 in  $O(n)$  time as there are only this many facets of  $P$ . This settles Case 2.

For case 3 the amortized cost for each arc is  $O(n \log^2 n \cdot \text{EO} + n^2 \log^{O(1)} n)$ . Our desired runtime follows since there are only  $O(n^2)$  arcs to add. Unlike Case 2 some extra care is needed to handle the overhead cost. The time needed to deduce a new arc (applying Lemmas 10.4.9 and 10.4.10 to  $\mathbf{y}$  and  $p \in B_1 \cup \dots \cup B_k$ ) is still  $O(n \cdot \text{EO} + n^2)$ . But as soon as we get a new arc, we must update  $A$  to be its transitive closure so that it is still complete. Given  $A$  complete and a new arc  $(p, q) \notin A$ , we can simply add the arcs from the ancestors of  $p$  to  $q$  and from  $p$  to the descendants of  $q$ . There are at most  $O(n)$  arcs to add so this takes time  $O(n^2)$  per arc, which is okay.  $\square$

## ■ 10.5 Discussion and Comparison with Previous Algorithms

We compare and contrast our algorithms with the previous ones. We focus primarily on strongly polynomial time algorithms.

### Convex combination of BFS's

All of the previous algorithms maintain a convex combination of BFS's and iteratively improve over it to get a better primal solution. In particular, the new BFS's used are typically obtained by making local changes to existing ones. Our algorithms, on the other hand, considers the geometry of the existing BFS's. The weighted "influences"<sup>3</sup> then aggregately govern the choice of the next BFS. We believe that this is the main driving force for the speedup of our algorithms.

### Scaling schemes

Many algorithms for combinatorial problems are explicitly or implicitly scaling a potential function or a parameter. In this chapter, our algorithms in some sense aim to minimize the volume of the feasible region. Scaling schemes for different potential functions and parameters were also designed in previous works [121, 118, 125, 119]. All of these functions and parameters have an explicit form. On the contrary, our potential function is somewhat unusual in the sense that it has no closed form.

<sup>3</sup>In the terminology of Chapter 8, these weighted influences are the leverage scores.

### Deducing new constraints

As mentioned in the main text, our algorithms share the same skeleton and tools for deducing new constraints with [121, 118, 125, 119]. Nevertheless, there are differences in the way these tools are employed. Our algorithms proceed by invoking them in a geometric manner, whereas previous algorithms were mostly combinatorial.

### Big elements and bucketing

Our bucketing idea has roots in Iwata-Orlin's algorithm [125] but is much more sophisticated. For instance, it is sufficient for their algorithm to consider only big elements, i.e.  $\text{upper}(i) \geq N/n^{O(1)}$ . Our algorithm, on the other hand, must carefully group elements by the size of both  $\text{upper}(i)$  and  $\text{lower}(i)$ . The speedup appears impossible without these new ideas. We do however note that it is unfair to expect such a sophisticated scheme in Iwata-Orlin's algorithm as it would not lead to a speedup. In other words, their method is fully sufficient for their purposes, and the simplicity in their case is a virtue rather than a shortcoming.

## ■ 10.5.1 Open Problems

One natural open problem is improving our weakly polynomial algorithm to  $O(n^2 \log M \cdot \text{EO} + n^3 \log^{O(1)} n \cdot \log M)$  time. Our application of center of mass to SFM demonstrates that it should be possible.

For strongly polynomial algorithms, the existential result of Theorem 10.4.1 shows that SFM can be solved with  $O(n^3 \log n \cdot \text{EO})$  oracle calls. Unfortunately, our algorithm incurs an overhead of  $\log n$  as there can be as many as  $\log n$  buckets each time. One may try to remove this  $\log n$  overhead by designing a better bucketing scheme or arguing that more arcs can be added.

The other  $\log n$  overhead seem much trickier to remove. Our method currently makes crucial use of the tools developed by [121], where the  $\log n$  factors in the runtime seem inevitable. We suspect that our algorithm may have an analogue similar to [231, 216], which do not carry any  $\log n$  overhead in the running time.

Perhaps an even more interesting open problem is whether our algorithm is optimal (up to polylogarithmic factors). There are grounds for optimism. So far the best way of *certifying* the optimality of a given solution  $S \subseteq V$  is to employ duality and express some optimal solution to the base polyhedron as a convex combination of  $n + 1$  BFS's. This already takes  $n^2$  oracle calls as each BFS requires  $n$ . Thus one would expect the optimal number of oracle calls needed for SFM to be at least  $n^2$ . Our bound is not too far off from it, and anything strictly between  $n^2$  and  $n^3$  seems instinctively unnatural.



## First Order Method vs Cutting Plane Method

### ■ 11.1 Introduction

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a  $\beta$ -smooth and  $\alpha$ -strongly convex function. Thus, for any  $x, y \in \mathbb{R}^n$ , we have

$$f(x) + \nabla f(x)^\top (y - x) + \frac{\alpha}{2} |y - x|^2 \leq f(y) \leq f(x) + \nabla f(x)^\top (y - x) + \frac{\beta}{2} |y - x|^2.$$

Let  $\kappa = \beta/\alpha$  be its condition number. It is a one line calculation to verify that a step of gradient descent on  $f$  will decrease (multiplicatively) the squared distance to the optimum by  $1 - 1/\kappa$ .

In Section 11.3, we propose a new method, which can be viewed as some combination of gradient descent and the ellipsoid method, for which the squared distance to the optimum decreases at a rate of  $(1 - 1/\sqrt{\kappa})$  (and each iteration requires one gradient evaluation and two line-searches). This matches the optimal rate of convergence among the class of first-order methods, [203, 206].

In Section 11.4, we generalize this idea and show how to automatically combine different optimization algorithms. In particular, we show how to combine gradient descent and cutting plane methods to get an algorithm with convergence rate  $\min(1 - 1/\kappa, 1 - 1/O(n))$ . This is the first algorithm we know that can naturally achieve both  $\kappa$  dependence and  $n$  dependence guarantee.

In Section 11.5, we provide some numerical evidence that the new methods can be superior to Nesterov's accelerated gradient descent.

#### ■ 11.1.1 Preliminaries

We write  $|\cdot|$  for the Euclidean norm in  $\mathbb{R}^n$ , and  $B(x, r^2) := \{y \in \mathbb{R}^n : |y - x|^2 \leq r^2\}$  (note that the second argument is the radius squared). We define the map `line_search` :  $\mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}^n$  by

$$\text{line\_search}(x, y) = \arg \min_{t \in \mathbb{R}} f(x + t(y - x)),$$

and we denote

$$x^+ = x - \frac{1}{\beta} \nabla f(x), \text{ and } x^{++} = x - \frac{1}{\alpha} \nabla f(x).$$

Recall that by strong convexity one has

$$\forall y \in \mathbb{R}^n, f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{\alpha}{2} |y - x|^2,$$

which implies in particular:

$$x^* \in B\left(x^{++}, \frac{|\nabla f(x)|^2}{\alpha^2} - \frac{2}{\alpha}(f(x) - f(x^*))\right).$$

Furthermore recall that by smoothness one has  $f(x^+) \leq f(x) - \frac{1}{2\beta}|\nabla f(x)|^2$  which allows to *shrink* the above ball by a factor of  $1 - \frac{1}{\kappa}$  and obtain the following:

$$x^* \in B\left(x^{++}, \frac{|\nabla f(x)|^2}{\alpha^2} \left(1 - \frac{1}{\kappa}\right) - \frac{2}{\alpha}(f(x^+) - f(x^*))\right) \quad (11.1)$$

## ■ 11.2 Intuition

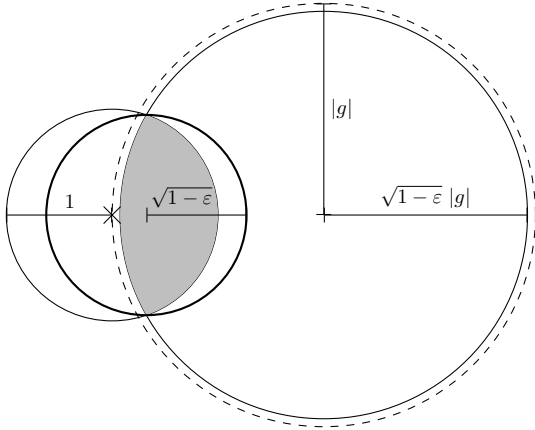


Figure 11-1: One ball shrinks.

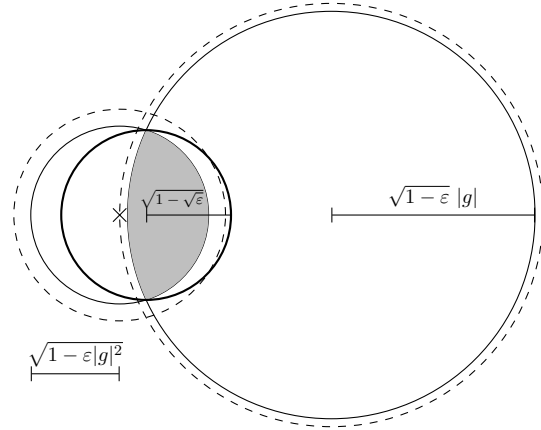


Figure 11-2: Two balls shrink.

The left diagram shows the intersection shrinks at the same rate if only one of the ball shrinks; the right diagram shows the intersection shrinks much faster if two balls shrink at the same absolute amount.

In Section 11.2.1 we describe a geometric alternative to gradient descent (with the same convergence rate) which gives the core of our new optimal method. Then in Section 11.2.2 we explain why one can expect to accelerate this geometric algorithm.

### ■ 11.2.1 A Suboptimal Algorithm

Assume that we are given a guarantee  $R_0 > 0$  on the distance from some point  $x_0$  to the optimum, that is  $x^* \in B(x_0, R_0^2)$ . Combining this original enclosing ball for  $x^*$  and the one obtained by (11.1) (with  $f(x^*) \leq f(x_0^+)$ ) one obtains

$$x^* \in B(x_0, R_0^2) \cap B\left(x_0 - \frac{1}{\alpha}\nabla f(x_0), \frac{|\nabla f(x_0)|^2}{\alpha^2} \left(1 - \frac{1}{\kappa}\right)\right).$$

If  $\frac{|\nabla f(x_0)|^2}{\alpha^2} \leq R_0^2(1 - \frac{1}{\kappa})$  then the second ball already shrinks by a factor of  $(1 - \frac{1}{\kappa})$ . In the other case when  $\frac{|\nabla f(x_0)|^2}{\alpha^2} > R_0^2(1 - \frac{1}{\kappa})$ , the center of the two balls are far apart and therefore there is a much smaller ball containing the intersection of two balls. Formally, it is an easy calculation to see that for any  $g \in \mathbb{R}^n$ ,  $\varepsilon \in (0, 1)$ , there exists  $x \in \mathbb{R}^n$  such that

$$B(0, 1) \cap B(g, |g|^2(1 - \varepsilon)) \subset B(x, 1 - \varepsilon). \quad (\text{Figure 11-1})$$

In particular the two above display implies that there exists  $x_1 \in \mathbb{R}^n$  such that

$$x^* \in B\left(x_1, R_0^2 \left(1 - \frac{1}{\kappa}\right)\right).$$



Denote by  $T$  the map from  $x_0$  to  $x_1$  defined implicitly above, and let  $(x_k)$  be defined by  $x_{k+1} = T(x_k)$ . Then we just proved

$$|x^* - x_k|^2 \leq \left(1 - \frac{1}{\kappa}\right)^k R_0^2.$$

In other words, after  $2\kappa \log(R_0/\varepsilon)$  iterations where each iteration cost one call to the gradient oracle) one obtains a point  $\varepsilon$ -close to the minimizer of  $f$ .

### ■ 11.2.2 Why One Can Accelerate

Assume now that we are give a guarantee  $R_0 > 0$  such that  $x^* \in \mathbb{B}(x_0, R_0^2 - \frac{2}{\alpha}(f(y) - f(x^*)))$  where  $f(x_0) \leq f(y)$  (say by choosing  $y = x_0$ ). Using the fact that  $f(x_0^+) \leq f(x_0) - \frac{1}{2\beta}|\nabla f(x_0)|^2 \leq f(y) - \frac{1}{2\alpha\kappa}|\nabla f(x_0)|^2$ , we obtain that

$$x^* \in \mathbb{B}\left(x_0, R_0^2 - \frac{|\nabla f(x_0)|^2}{\alpha^2\kappa} - \frac{2}{\alpha}(f(x_0^+) - f(x^*))\right)$$

which, intuitively, allows us the shrink the radius squared from  $R_0^2$  to  $R_0^2 - \frac{|\nabla f(x_0)|^2}{\alpha^2\kappa}$  using the local information at  $x_0$ . From (11.1), we have

$$x^* \in \mathbb{B}\left(x_0^{++}, \frac{|\nabla f(x_0)|^2}{\alpha^2} \left(1 - \frac{1}{\kappa}\right) - \frac{2}{\alpha}(f(x_0^+) - f(x^*))\right).$$

Now, intersecting the above two shrunk balls and using Lemma 11.3.2 (see below and also see Figure 11-2), we obtain that there is an  $x'_1$  such that

$$x^* \in \mathbb{B}\left(x'_1, R_0^2 \left(1 - \frac{1}{\sqrt{\kappa}}\right) - \frac{2}{\alpha}(f(x_0^+) - f(x^*))\right)$$

giving us an acceleration in shrinking of the radius. To carry the argument for the next iteration, we would have required that  $f(x'_1) \leq f(x_0^+)$  but it may not hold. Thus, we choose  $x_1$  by a line search

$$x_1 = \text{line\_search}(x'_1, x_0^+)$$

which ensures that  $f(x_1) \leq f(x_0^+)$ . To remedy the fact that the ball for the next iteration is centered at  $x'_1$  and not  $x_1$ , we observe that the line search also ensures that  $\nabla f(x_1)$  is perpendicular to the line going through  $x_1$  and  $x'_1$ . This geometric fact is enough for the algorithm to work at the next iteration as well. In the next section we describe precisely our proposed algorithm which is based on the above insights.

## ■ 11.3 An Optimal Algorithm

Let  $x_0 \in \mathbb{R}^n$ ,  $c_0 = x_0^{++}$ , and  $R_0^2 = \left(1 - \frac{1}{\kappa}\right) \frac{|\nabla f(x_0)|^2}{\alpha^2}$ . For any  $k \geq 0$  let

$$x_{k+1} = \text{line\_search}(c_k, x_k^+),$$

and  $c_{k+1}$  (respectively  $R_{k+1}^2$ ) be the center (respectively the squared radius) of the ball given by (the proof of) Lemma 11.3.2 which contains

$$\mathbb{B}\left(c_k, R_k^2 - \frac{|\nabla f(x_{k+1})|^2}{\alpha^2\kappa}\right) \cap \mathbb{B}\left(x_{k+1}^{++}, \frac{|\nabla f(x_{k+1})|^2}{\alpha^2} \left(1 - \frac{1}{\kappa}\right)\right).$$

**Algorithm 27:** Minimum Enclosing Ball of the Intersection to Two Balls

**Input:** a ball centered at  $x_A$  with radius  $R_A$  and a ball centered at  $x_B$  with radius  $R_B$ .

**if**  $|x_A - x_B|^2 \geq |R_A^2 - R_B^2|$  **then**

$$c = \frac{1}{2}(x_A + x_B) - \frac{R_A^2 - R_B^2}{2|x_A - x_B|^2}(x_A - x_B).$$

$$R^2 = R_B^2 - \frac{(|x_A - x_B|^2 + R_B^2 - R_A^2)^2}{4|x_A - x_B|^2}.$$

**else if**  $|x_A - x_B|^2 < R_A^2 - R_B^2$  **then**

$$c = x_B. \quad R = R_B.$$

**else**

$$c = x_A. \quad R = R_A.$$

**end**

**Output:** a ball centered at  $c$  with radius  $R$ .

**Algorithm 28:** Geometric Descent Method (GeoD)

**Input:** parameters  $\alpha$  and initial points  $x_0$ .

$$x_0^+ = \text{line\_search}(x_0, x_0 - \nabla f(x_0)).$$

$$c_0 = x_0 - \alpha^{-1} \nabla f(x_0).$$

$$R_0^2 = \frac{|\nabla f(x_0)|^2}{\alpha^2} - \frac{2}{\alpha} (f(x_0) - f(x_0^+)).$$

**for**  $i \leftarrow 1, 2, \dots$  **do**

**Combining Step:**

$$x_k = \text{line\_search}(x_{k-1}^+, c_{k-1}).$$

**Gradient Step:**

$$x_k^+ = \text{line\_search}(x_k, x_k - \nabla f(x_k)).$$

**Ellipsoid Step:**

$$x_A = x_k - \alpha^{-1} \nabla f(x_k). \quad R_A^2 = \frac{|\nabla f(x_k)|^2}{\alpha^2} - \frac{2}{\alpha} (f(x_k) - f(x_k^+)).$$

$$x_B = c_{k-1}. \quad R_B^2 = R_{k-1}^2 - \frac{2}{\alpha} (f(x_{k-1}^+) - f(x_k^+)).$$

Let  $B(c_k, R_k^2)$  is the minimum enclosing ball of  $B(x_A, R_A^2) \cap B(x_B, R_B^2)$ .

**end**

**Output:**  $x_T$ .

The formulas for  $c_{k+1}$  and  $R_{k+1}^2$  are given in Algorithm 27.

**Theorem 11.3.1.** For any  $k \geq 0$ , one has  $x^* \in B(c_k, R_k^2)$ ,  $R_{k+1}^2 \leq \left(1 - \frac{1}{\sqrt{\kappa}}\right) R_k^2$ , and thus

$$|x^* - c_k|^2 \leq \left(1 - \frac{1}{\sqrt{\kappa}}\right)^k R_0^2.$$

*Proof.* We will prove a stronger claim by induction that for each  $k \geq 0$ , one has

$$x^* \in B\left(c_k, R_k^2 - \frac{2}{\alpha} (f(x_k^+) - f(x^*))\right).$$

The case  $k = 0$  follows immediately by (11.1). Let us assume that the above display is true for some  $k \geq 0$ . Then using  $f(x^*) \leq f(x_{k+1}^+) \leq f(x_{k+1}) - \frac{1}{2\beta} |\nabla f(x_{k+1})|^2 \leq f(x_k^+) - \frac{1}{2\beta} |\nabla f(x_{k+1})|^2$ , one gets

$$x^* \in B\left(c_k, R_k^2 - \frac{|\nabla f(x_{k+1})|^2}{\alpha^2 \kappa} - \frac{2}{\alpha} (f(x_{k+1}^+) - f(x^*))\right).$$

Furthermore by (11.1) one also has

$$B\left(x_{k+1}^{++}, \frac{|\nabla f(x_{k+1})|^2}{\alpha^2} \left(1 - \frac{1}{\kappa}\right) - \frac{2}{\alpha} (f(x_{k+1}^+) - f(x^*))\right).$$

Thus it only remains to observe that the squared radius of the ball given by Lemma 11.3.2 which encloses the intersection of the two above balls is smaller than  $\left(1 - \frac{1}{\sqrt{\kappa}}\right) R_k^2 - \frac{2}{\alpha} (f(x_{k+1}^+) - f(x^*))$ . We apply Lemma 11.3.2 after moving  $c_k$  to the origin and scaling distances by  $R_k$ . We set  $\varepsilon = \frac{1}{\kappa}$ ,  $g = \frac{|\nabla f(x_{k+1})|}{\alpha}$ ,  $\delta = \frac{2}{\alpha} (f(x_{k+1}^+) - f(x^*))$  and  $a = x_{k+1}^{++} - c_k$ . The line search step of the algorithm implies that  $\nabla f(x_{k+1})^\top (x_{k+1} - c_k) = 0$  and therefore,  $|a| = |x_{k+1}^{++} - c_k| \geq |\nabla f(x_{k+1})|/\alpha = g$  and Lemma 11.3.2 applies to give the result.  $\square$

**Lemma 11.3.2.** *Let  $a \in \mathbb{R}^n$  and  $\varepsilon \in (0, 1), g \in \mathbb{R}_+$ . Assume that  $|a| \geq g$ . Then there exists  $c \in \mathbb{R}^n$  such that for any  $\delta \geq 0$ ,*

$$B(0, 1 - \varepsilon g^2 - \delta) \cap B(a, g^2(1 - \varepsilon) - \delta) \subset B(c, 1 - \sqrt{\varepsilon} - \delta).$$

*Proof.* First observe that if  $g^2 \leq 1/2$  then one can take  $c = a$  since  $\frac{1}{2}(1 - \varepsilon) \leq 1 - \sqrt{\varepsilon}$ . Thus we assume now that  $g^2 > 1/2$ , and note that we can also clearly assume that  $n = 2$ . Consider the segment joining the two points at the intersection of the two balls under consideration. We denote  $c$  for the point at the intersection of this segment and  $[0, a]$ , and  $x = |c|$  (that is  $c = x \frac{a}{|a|}$ ). A simple picture reveals that  $x$  satisfies

$$1 - \varepsilon g^2 - \delta - x^2 = g^2(1 - \varepsilon) - \delta - (|a| - x)^2 \Leftrightarrow x = \frac{1 + |a|^2 - g^2}{2|a|}.$$

When  $x \leq |a|$ , neither of the balls covers more than half of the other ball and hence the intersection of the two balls is contained in the ball  $B\left(x \frac{a}{|a|}, 1 - \varepsilon g^2 - \delta - x^2\right)$  (See figure 11-2). Thus it only remains to show that  $x \leq |a|$  and that  $1 - \varepsilon g^2 - \delta - x^2 \leq 1 - \sqrt{\varepsilon} - \delta$ . The first inequality is equivalent to  $|a|^2 + g^2 \geq 1$  which follows from  $|a|^2 \geq g^2 \geq 1/2$ . The second inequality to prove can be written as

$$\varepsilon g^2 + \frac{(1 + |a|^2 - g^2)^2}{4|a|^2} \geq \sqrt{\varepsilon},$$

which is straightforward to verify (recall that  $|a|^2 \geq g^2 \geq 1/2$ ).  $\square$

Algorithm 28 we give is more aggressive than Theorem 11.3.1, for instance, using line search instead of fixed step size. The correctness of this version follows from a similar proof as Theorem 11.3.1.

This algorithm does not require the smoothness parameter and the number of iterations; and it guarantees the function value is strictly decreasing. They are useful properties for machine learning applications because the only required parameter  $\alpha$  is usually given. Furthermore, we believe that the integration of zeroth and first order information about the function makes our new method particularly well-suited in practice.

## ■ 11.4 Politician

In this section, we discuss how to in general combine different optimization algorithms. In standard black-box convex optimization [201, 206, 40] first-order methods interact with an *oracle*: given a query point  $x$ , the oracle reports the value and gradient of the underlying objective function  $f$  at  $x$ . Here, we propose to replace the oracle by a *politician*. Instead of answering the original query  $x$  the politician changes the question and answers a new query  $y$  which is guaranteed to be better than the original

query  $x$  in the sense that  $f(y) \leq f(x)$ . The newly selected query  $y$  also depends on the history of queries that were made to the politician. Formally we introduce the following definition (for sake of simplicity we write  $\nabla f(x)$  for either a gradient or a subgradient of  $f$  at  $x$ ).

**Definition 11.4.1.** Let  $\mathcal{X} \subset \mathbb{R}^n$  and  $f : \mathcal{X} \rightarrow \mathbb{R}$ . A politician  $\Phi$  for  $f$  is a mapping from  $\mathcal{X} \times \cup_{k=0}^{\infty} (\mathcal{X} \times \mathbb{R} \times \mathbb{R}^n)^k$  to  $\mathcal{X}$  such that for any  $k \geq 0, x \in \mathcal{X}, h \in (\mathcal{X} \times \mathbb{R} \times \mathbb{R}^n)^k$  one has  $f(\Phi(x, h)) \leq f(x)$ . Furthermore when queried at  $x$  with history  $h$  a politician for  $f$  also output  $f(\Phi(x, h))$  and  $\nabla f(\Phi(x, h))$  (in order to not overload notation we do not include these outputs in the range of  $\Phi$ ).

Let us clarify the interaction of a first-order method with a politician. Note that we refer to the couple (first-order method, politician) as the *algorithm*. Let  $M : \cup_{k=0}^{\infty} (\mathcal{X} \times \mathbb{R} \times \mathbb{R}^n)^k \rightarrow \mathcal{X}$  be a first-order method and  $\Phi$  a politician for some function  $f : \mathcal{X} \rightarrow \mathbb{R}$ . The course of the algorithm  $(M, \Phi)$  then goes as follows: at iteration  $k + 1$  one first calculates the method's query point  $x_{k+1} = M(h_k)$  (with  $h_0 = \emptyset$ ), then one calculates the politician's new query point  $y_{k+1} = \Phi(x_{k+1}, h_k)$  and the first order information at this point  $(f(y_{k+1}), \nabla f(y_{k+1}))$ , and finally one updates the history with this new information  $h_{k+1} = (h_k, (y_{k+1}, f(y_{k+1}), \nabla f(y_{k+1})))$ . Note that a standard oracle simply corresponds to a politician  $\mathcal{O}$  for  $f$  such that  $\mathcal{O}(x, h) = (x, f(x), \nabla f(x))$  (in particular the algorithm  $(M, \mathcal{O})$  is the usual algorithm corresponding to the first-order method  $M$ ).

The philosophy of the above definition is that it gives in some sense an automatic way to combine different optimization algorithms. Say for example that we wish to combine the ellipsoid method with gradient descent. One way to do so is to design an “ellipsoidal politician”: the politician keeps track of a feasible ellipsoidal region based on the previously computed gradients, and when asked with the query  $x$  the politician chooses as a new query  $y$  the result of a line-search on the line between  $x$  and the center of current ellipsoid. Gradient descent with this ellipsoidal politician would then replace the step  $x \leftarrow x - \eta \nabla f(x)$  by  $x \leftarrow y - \eta \nabla f(y)$ . The hope is that in practice such a combination would integrate the fast incremental progress of gradient descent with the geometrical progress of the ellipsoid method.

In this section, we focus on unconstrained convex optimization. We are particularly interested in situations where calculating a (sub)gradient has superlinear complexity (i.e.,  $\gg n$ ) such as in logistic regression and semidefinite programming. In such cases it is natural to try to make the most out of the computed gradients by incorporating geometric reasoning (such as in the ellipsoid method). We do so by introducing the *geometric politician* (Section 11.4.2), which is based on a combination of the ideas in previous section with standard cutting plane machinery (through the notion of a “center” of a set, see Section 11.4.3). For a given first order method  $M$ , we denote by  $M+$  the algorithm obtained by running  $M$  with the geometric politician. We demonstrate empirically (Section 11.5) the effectiveness of the geometric politician on various standard first-order methods for convex optimization (gradient descent, Nesterov's accelerated gradient descent, non-linear conjugate gradient, BFGS). In particular we show that BFGS+ is a surprisingly robust and parameter-free algorithm with state of the art performance across a wide range of problems (both smooth and non-smooth).

### ■ 11.4.1 Affine Invariant Politician

As mentioned above we assume that the complexity of computing the map  $x \mapsto \nabla f(x)$  is superlinear. This implies that we can afford to have a politician such that the complexity of computing the map  $(x, h) \mapsto \Phi(x, h)$  is  $O(n \times \text{poly}(k))$  (we think of the number of iterations  $k$  as typically much smaller than the dimension  $n$ ). We show in this section that this condition is (essentially) automatically satisfied as long as the politician is affine invariant in the following sense (we use a slight abuse of language and refer to a map  $f \mapsto \Phi_f$ , where  $\Phi_f$  is a politician for  $f$ , as a politician):

**Definition 11.4.2.** A politician  $f \mapsto \Phi_f$  is called affine invariant if for any function  $f$  and any affine map  $T : \mathbb{R}^m \rightarrow \mathbb{R}^n$  such that  $T(x) = z + Lx$  for some matrix  $L$ ,  $k \geq 0$ ,  $x \in \mathbb{R}^m$ ,  $(y_i, v_i, g_i) \in \mathbb{R}^m \times \mathbb{R} \times \mathbb{R}^n$ , one has

$$T(\Phi_{f \circ T}(x, (y_i, v_i, L^\top g_i)_{i \in [k]})) = \Phi_f(T(x), (T(y_i), v_i, g_i)_{i \in [k]}).$$

We say that an affine invariant politician has cost  $\psi : \mathbb{N} \rightarrow \mathbb{N}$  if for any  $f : \mathbb{R}^k \rightarrow \mathbb{R}$  the map  $(x, h) \in \mathbb{R}^k \times (\mathbb{R}^k \times \mathbb{R} \times \mathbb{R}^k)^k \mapsto \Phi_f(x, h)$  can be computed in time  $\psi(k)$ .

**Proposition 11.4.3.** Let  $\Phi$  be an affine invariant politician with cost  $\psi$ . Then for any  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $(y_i, v_i, g_i) \in \mathbb{R}^n \times \mathbb{R} \times \mathbb{R}^n$ ,  $i \in [k]$  and  $x, y_i \in y_1 + \text{Span}(g_1, \dots, g_k)$  one can compute  $\Phi_f(x, (y_i, v_i, g_i)_{i \in [k]}) \in \mathbb{R}^n$  in time  $\psi(k) + O(nk^2)$ .

*Proof.* Let  $G$  be the  $n \times k$  matrix with  $i^{\text{th}}$  column given by  $g_i$ . We consider the  $QR$  decomposition of  $G$  which can be computed in time  $O(nk^2)$ , that is  $Q$  is an  $n \times k$  matrix and  $R$  a  $k \times k$  matrix such that  $G = QR$  and  $Q^\top Q = I_k$ . Let  $T$  be the affine map defined by  $T = y_1 + Q$ . Note that since  $x \in y_1 + \text{Span}(g_1, \dots, g_k)$  one has  $x = T(Q^\top(x - y_1))$  (and similarly for  $y_i$ ). Thus by affine invariance one has

$$\begin{aligned} \Phi_f(x, (y_i, v_i, g_i)) &= \Phi_f(T(Q^\top(x - y_1)), (T(Q^\top(y_i - y_1)), v_i, g_i)) \\ &= y_1 + Q\Phi_{f \circ T}(Q^\top(x - y_1), (Q^\top(y_i - y_1), v_i, R_i)), \end{aligned}$$

where  $R_i$  is the  $i^{\text{th}}$  column of  $R$ . Furthermore by definition of the cost  $\psi$  and since  $f \circ T$  is defined on  $\mathbb{R}^k$  we see that this last quantity can be computed in time  $\psi(k) + O(nk^2)$ , thus concluding the proof.  $\square$

The above proposition shows that with an affine invariant politician and a first order method  $M$  verifying for any  $(y_i, v_i, g_i)_{i \in [k]} \in (\mathbb{R}^n \times \mathbb{R} \times \mathbb{R}^n)^k$ ,

$$M((y_i, v_i, g_i)_{i \in [k]}) \in y_1 + \text{Span}(y_1, \dots, y_k, g_1, \dots, g_k),$$

one can run  $k$  steps of the corresponding algorithm in time  $O(nk^2 + k\psi(k))$  plus the time to compute the  $k$  function values and gradients of the underlying function  $f$  to be optimized. Note that one gets a time of  $O(nk^2)$  instead of  $O(nk^3)$  as one can store the  $QR$  decomposition from one step to the next, and updating the decomposition only cost  $O(nk)$ .

## ■ 11.4.2 Geometric Politician

We describe in this section the *geometric politician* which is based on ideas developed in Section 11.3. A key observation in that Section is that if  $f$  is a  $\alpha$ -strongly convex function minimized at  $x^*$  then one has for any  $x$ ,

$$\left\| x^* - x - \frac{1}{\alpha} \nabla f(x) \right\|^2 \leq \frac{\|\nabla f(x)\|^2}{\alpha^2} - \frac{2}{\alpha} (f(x) - f(x^*)).$$

This motivates the following definition:

$$B(x, \alpha, \text{fval}) := \left\{ z \in \mathbb{R}^n : \left\| z - x - \frac{1}{\alpha} \nabla f(x) \right\|^2 \leq \frac{\|\nabla f(x)\|^2}{\alpha^2} - \frac{2}{\alpha} (f(x) - \text{fval}) \right\}.$$

In particular given the first order information at  $y_1, \dots, y_k$  one knows that the optimum  $x^*$  lies in the region  $R_k \subset \mathbb{R}^n$  defined by

$$R_k = \bigcap_{i \in [k]} B(y_i, \alpha, \text{fval}) \text{ where } \text{fval} = \min_{i \in [k]} f(y_i). \quad (11.2)$$

**Algorithm 29:** Geometric Politician

**Parameter:** An upper bound on the strong convexity parameter  $\alpha$ . (Can be  $+\infty$ .)

**Input:** Query  $x$ , past queries and the corresponding first order information

$$(y_i, f(y_i), \nabla f(y_i))_{i \in [k]}.$$

Let  $f_{\text{val}} = \min_{i \in [k]} f(y_i)$  and the feasible region  $R_k(\alpha) = \bigcap_{i \in [k]} B(y_i, \alpha, f_{\text{val}})$ .

**if**  $R_k(\alpha) = \emptyset$  **then**

    Let  $\alpha$  be the largest number such that  $R_k(\alpha) \neq \emptyset$ .

$\alpha \leftarrow \alpha/4$ .

**end**

Let  $y_{k+1} = \arg \min_{y \in \{tx + (1-t)c(R_k(\alpha)), t \in \mathbb{R}\}} f(y)$

where  $c(R_k(\alpha))$  is the volumetric center of  $R_k(\alpha)$  (see Section 11.4.3).

**Output:**  $y_{k+1}$ ,  $f(y_{k+1})$  and  $\nabla f(y_{k+1})$ .

Now suppose that given this first order information at  $y_1, \dots, y_k$  the first order method asks to query  $x$ . How should we modify this query in order to take into account the geometric information that  $x^* \in R_k$ ? First observe that for any  $z$ ,  $B(z, \alpha, f_{\text{val}})$  is contained in a halfspace that has  $z$  on its boundary (in the limiting case  $\alpha \rightarrow 0$  the set  $B(z, \alpha, f(z))$  is exactly a halfspace). In particular if the next query point  $y_{k+1}$  is the center of gravity of  $R_k$  then we have that the volume of  $R_{k+1}$  is at most  $1 - 1/e$  times the volume of  $R_k$  (see [111]), thus leading to an exponential convergence rate. However the region  $R_k$  can be very large initially, and the center of gravity might have a large function value and gradient, which means that  $R_k$  would be intersected with a large sphere (possibly so large that it is close to a halfspace). On the other hand the first order method recommends to query  $x$ , which we can think of as a local improvement of  $y_k$ , which should lead to a much smaller sphere. The issue is that the position of this sphere might be such that the intersection with  $R_k$  is almost as large as the sphere itself. In order to balance between the geometric and function value/gradient considerations we propose for the new query to do a line search between the center of  $R_k$  and the recommended query  $x$ . The geometric politician follows this recipe with two important modifications: (i) there are many choices of centers that would guarantee an exponential convergence rate while being much easier to compute than the center of gravity, and we choose here to consider the *volumetric center*, see Section 11.4.3 for the definition and more details about this notion; (ii) we use a simple heuristic to adapt online the strong convexity parameter  $\alpha$ , namely we start with some large value for  $\alpha$  and if it happens that the feasible region  $R_k$  is empty then we know that  $\alpha$  was too large, in which case we reduce it. We can now describe formally the geometric politician, see Algorithm 29. Importantly one can verify that the geometric politician is affine invariant and thus can be implemented efficiently (see the proof of Proposition 11.4.3).

### ■ 11.4.3 Volumetric Center

The volumetric barrier for a polytope was introduced in [253] to construct an algorithm with both the oracle complexity of the center of gravity method and the computational complexity of the ellipsoid method (see [Section 2.3, [40]] for more details and Chapter 8 for recent advances on this construction). Recalling that the standard logarithmic barrier  $F_P$  for the polytope  $P = \{x \in \mathbb{R}^n : a_i^\top x < b_i, i \in [m]\}$  is defined by

$$F_P(x) = - \sum_{i=1}^m \log(b_i - a_i^\top x),$$

one defines the volumetric barrier  $v_P$  for  $P$  by

$$v_P(x) = \log \det(\nabla^2 F_P(x)).$$

The volumetric center  $c(P)$  is then defined as the minimizer of  $v_P$ . In the context of the geometric politician (see Algorithm 29) we are dealing with an intersection of balls rather than an intersection of halfspaces. More precisely the region of interest is of the form:

$$R = \bigcap_{i=1}^k \{x \in \mathbb{R}^n : \|x - c_i\| \leq r_i\}.$$

For such a domain the natural self-concordant barrier to consider is:

$$F_R(x) = -\frac{1}{2} \sum_{i=1}^k \log(r_i^2 - \|x - c_i\|^2).$$

The volumetric barrier is defined as before by

$$v_R(x) = \log \det(\nabla^2 F_R(x)),$$

and the volumetric center of  $R$  is the minimizer of  $v_R$ . It is shown in [12] that  $v_R$  is a self-concordant barrier which means that the center can be updated (when a new ball is added to  $R$ ) via few iterations of Newton's method. Often in practice, it takes less than 5 iterations to update the minimizer of a self-concordant barrier [103, 26] when we add a new constraint. Hence, the complexity merely depends on how fast we can compute the gradient and Hessian of  $F_R$  and  $v_R$ .

**Proposition 11.4.4.** *For the analytic barrier  $F_R$ , we have that*

$$\begin{aligned} \nabla F_R(x) &= A^\top \mathbf{1}_{k \times 1}, \\ \nabla^2 F_R(x) &= 2A^\top A + \lambda^{(1)} I \end{aligned}$$

where  $d$  is a vector defined by  $(r_i^2 - \|x - c_i\|^2)^{-1}$ ,  $A$  is a  $k \times n$  matrix with  $i^{\text{th}}$  row given by  $d_i(x)(x - c_i)$ ,  $\lambda^{(p)} = \sum_{i \in [k]} d_i^p(x)$  and  $\mathbf{1}_{k \times 1}$  is a  $k \times 1$  matrix with all entries being 1.

For the volumetric center, we have that

$$\begin{aligned} \nabla v_R(x) &= ((2\text{Tr}H^{-1})I + 4H^{-1})A^\top d + 8A^\top \sigma, \\ \nabla^2 v_R(x) &= 48A^\top \Sigma A - 64A^\top (AH^{-1}A^\top)^{(2)} A \\ &\quad + (8\text{Tr}(D\Sigma) + 2\lambda^{(2)}\text{Tr}(H^{-1}))I + 4\lambda^{(2)}H^{-1} \\ &\quad + 8\text{Tr}(H^{-1})A^\top DA + 16\text{sym}((A^\top DAH^{-1})) \\ &\quad - 4\text{Tr}(H^{-2})A^\top DJDA - 8H^{-1}A^\top DJDAH^{-1} \\ &\quad - 8\text{sym}((A^\top DJDAH^{-2})) - 8(d^\top AH^{-1}A^\top d)H^{-1} \\ &\quad - 16\text{sym}((A^\top \mathbf{Diag}(AH^{-2}A^\top)JDA)) \\ &\quad - 32\text{sym}((A^\top \mathbf{Diag}(AH^{-1}A^\top d)AH^{-1})) \end{aligned}$$

where  $H = \nabla^2 F_R(x)$ ,  $\sigma_i = e_i^\top AH^{-1}A^\top e_i$ ,  $e_i$  is the indicator vector with  $i^{\text{th}}$  coordinate,  $J$  is a  $k \times k$  matrix with all entries being 1,  $\text{sym}(\cdot)B = B + B^\top$ ,  $\mathbf{Diag}(v)$  is a diagonal matrix with  $\mathbf{Diag}(v)_{ii} = v_i$ ,  $D = \mathbf{Diag}(d)$ ,  $\Sigma = \mathbf{Diag}(\sigma)$ , and  $B^{(2)}$  is the Schur square of  $B$  defined by  $B_{ij}^{(2)} = B_{ij}^2$ .

The above proposition shows that one step of Newton method for analytic center requires 1 dense matrix multiplication and solving 1 linear system; and for volumetric center, it requires 5 dense matrix multiplications, 1 matrix inversion and solving 1 linear system if implemented correctly. Although the analytic center is a more popular choice for “geometrical” algorithms, we choose volumetric center here because it gives a better convergence rate [253, 20] and the extra cost  $\psi(k)$  is negligible to the cost of updating QR decomposition  $nk$ .

## ■ 11.5 Experiments

In this section, we compare the geometric politician against two libraries for first order methods, minFunc [230] and TFOCS [32]. Both are popular MATLAB libraries for minimizing general smooth convex functions. Since the focus of this chapter is all about how to find a good step direction, we use the exact line search (up to machine accuracy) whenever possible. This eliminates the effect of different line searches and reduces the number of algorithms we need to test. TFOCS is the only algorithm we use which does not use line search because they do not provide such option. To compensate on the unfairness to TFOCS, we note that the algorithm TFOCS uses is accelerated gradient descent and hence we implement the Gonzaga-Karas’s accelerated gradient descent [108], which is specifically designed to be used with exact line search. Another reason we pick this variant of accelerated gradient descent is because we found it to be the fastest variant of accelerated gradient descent (excluding the geometric descent in Section 11.3) for our tested data (Gonzaga and Karas also observed that on their own dataset).

The algorithms to be tested are the following:

- [SD] Steepest descent algorithm in minFunc.
- [Nes] Accelerated gradient descent, General Scheme 2.2.6 in [206].
- [TFOCS] Accelerated gradient descent in TFOCS.
- [GK] Gonzaga-Karas’s of Accelerated Gradient Descent (Section 11.5.1).
- [Geo] Geometric Descent (Section 11.3).
- [CG] Non-Linear Conjugate Gradient in minFunc.
- [BFGS] Broyden–Fletcher–Goldfarb–Shanno algorithm in minFunc.
- [PCG] Preconditioned Non-Linear Conjugate Gradient in minFunc.
- [ $\emptyset$ +] Geometric Politician itself (Section 11.5.1).
- [GK+] Using GK with Geometric Oracle (Section 11.5.1).
- [BFGS+] Using BFGS with Geometric Oracle (Section 11.5.1).

We only tested the geometric oracle on GK and BFGS because they are respectively the best algorithms in theory and practice on our tested data. The  $\emptyset$ + algorithm is used as the control group to test if the geometric politician by itself is sufficient to achieve good convergence rate. We note that all algorithms except Nes are parameter free; each step of SD, Nes, TFOCS, GK, Geo, CG takes  $O(n)$  time and each step of BFGS, PCG,  $\emptyset$ +, GK+ and BFGS+ takes roughly  $O(nk)$  time for  $k^{th}$  iteration.



**Algorithm 30:**  $\emptyset+$ 


---

**Input:**  $x_0$ .  
**for**  $k \leftarrow 1, 2, \dots$  **do**  
  | Set  $x_{k+1} \leftarrow \Phi_f(x_k, (x_i, f(x_i), \nabla f(x_i))_{i \in [k]})$ .  
**end**

---

**Algorithm 31:** Gonzaga-Karas's variant of Accelerated Gradient Descent

---

**Input:**  $x_1$ .  
 $\gamma = 2\alpha$ ,  $v_0 = x_0$  and  $y_0 = x_0$ .  
**for**  $k \leftarrow 1, 2, \dots$  **do**  
   $y_k \leftarrow \Phi_f(y_{k-1})$ .  
   $x_{k+1} = \text{line\_search}(y_k, -\nabla f(y_k))$ .  
  **if**  $\alpha \geq \gamma/1.02$  *and we are using first order oracle* **then**  $\alpha = \gamma/2$ . (\*);  
  **if**  $\alpha \geq \frac{\|\nabla f(y_k)\|^2}{2(f(y_k) - f(x_{k+1}))}$  **then**  $\alpha = \frac{\|\nabla f(y_k)\|^2}{20(f(y_k) - f(x_{k+1}))}$ ;  
   $G = \gamma \left( \frac{\alpha}{2} \|v_k - y_k\|^2 + \langle \nabla f(y_k), v_k - y_k \rangle \right)$ .  
   $A = G + \frac{1}{2} \|\nabla f(y_k)\|^2 + (\alpha - \gamma)(f(x_k) - f(y_k))$ .  
   $B = (\alpha - \gamma)(f(x_{k+1}) - f(x_k)) - \gamma(f(y_k) - f(x_k)) - G$ .  
   $C = \gamma(f(x_{k+1}) - f(x_k))$ .  
   $\beta = \frac{-B + \sqrt{B^2 - 4AC}}{2A}$ ,  $\gamma = (1 - \beta)\gamma + \beta\alpha$ .  
   $v_{k+1} = \frac{1}{\gamma}((1 - \beta)\gamma v_k + \beta(\alpha y_k - \nabla f(y_k)))$ .  
**end**

---

**■ 11.5.1 Details of Implementations**

The first algorithm we implement is the  $\emptyset+$  algorithm which simply repeatedly call the politician. As we will see, this algorithm is great for non smooth problems but not competitive for smooth problems.

The second algorithm we implement is the accelerated gradient descent proposed by Gonzaga and Karas [108]. This algorithm uses line search to learn the the smoothness parameter and strong convexity parameter, see Algorithm 31. We disable the line (\*) in the algorithm if  $\Phi_f$  is a politician instead of an oracle because  $\gamma \geq \alpha$  does not hold for the strong convexity parameter  $\alpha$  if  $\Phi_f$  is not an oracle.

The third algorithm we implemented is the Broyden–Fletcher–Goldfarb–Shanno (BFGS) algorithm. This algorithm uses the gradients to reconstruct the Hessian and use it to approximate Newton's method, see Algorithm 32. We note that another natural way to employ the politician with BFGS is to set  $x_{k+1} = \text{line\_search}(\Phi_f(x_k), p)$  and this runs faster in practice; however, this algorithm computes two gradients per iteration (namely  $\nabla f(x_k)$  and  $\nabla f(\Phi_f(x_k))$ ) while we restrict ourselves to algorithms which compute one gradient per iteration.

**■ 11.5.2 Quadratic Function**

We consider the function

$$f(x) = (x - c)^\top D(x - c), \quad (11.3)$$

where  $D$  is a diagonal matrix with entries uniformly sampled from  $[0, 1]$  and  $c$  is a random vector with entries uniformly sampled from the normal distribution  $N(0, 1)$ . Since this is a quadratic function, CG, BFGS and BFGS+ are equivalent and optimal, namely, they output the minimum point in the

**Algorithm 32:** BFGS

---

```

Input:  $x_1$ .
for  $k \leftarrow 1, 2, \dots$  do
   $p = -\nabla f(x_k)$ .
  for  $i \leftarrow k-1, \dots, 1$  do
     $\alpha_i \leftarrow \langle s_i, p \rangle / \langle s_i, y_i \rangle$ .
     $p = p - \alpha_i y_i$ .
  end
   $p = \langle s_{k-1}, y_{k-1} \rangle / \langle y_{k-1}, y_{k-1} \rangle p$ .
  for  $i \leftarrow 1, \dots, k-1$  do
     $\beta_i \leftarrow \langle y_i, p \rangle / \langle s_i, y_i \rangle$ .
     $p = p + (\alpha_i - \beta_i) y_i$ .
  end
   $x_{k+1} = \Phi_f(\text{line\_search}(x_k, p))$ .
   $s_k = x_{k+1} - x_k, y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$ .
end

```

---

span of all previous gradients.

### ■ 11.5.3 Variant of Nesterov's Worst Function

[206] introduced the function

$$f(x) = (1 - x[1])^2 + \sum_{k=1}^{n-1} (x[k] - x[k+1])^2$$

and used it to give a lower bound for all first-order methods. To distinguish the performance between CG, BFGS and BFGS+, we consider the following non-quadratic variant

$$f(x) = g(1 - x[1]) + \sum_{k=1}^{n-1} g(x[k] - x[k+1]) \quad (11.4)$$

for some function  $g$  to be defined. If we pick  $g(x) = |x|$  then all first order methods takes at least  $n$  iterations to minimize  $f$  exactly. On the other hand with  $g(x) = \max(|x| - 0.1, 0)$  one of the minimizer of  $f$  is  $(1, \frac{9}{10}, \frac{8}{10}, \dots, \frac{1}{10}, 0, 0, \dots, 0)$ , and thus it takes at least 11 iterations for first order methods to minimize  $f$  in this case. We “regularize” the situation a bit and consider the function

$$g(x) = \begin{cases} \sqrt{(x - 0.1)^2 + 0.001^2} - 0.001 & \text{if } x \geq 0.1 \\ \sqrt{(x + 0.1)^2 + 0.001^2} - 0.001 & \text{if } x \leq -0.1 \\ 0 & \text{otherwise} \end{cases}$$

Since this function is far from quadratic, our algorithms ( $\emptyset+$ , GK+, BFGS+) converge much faster. This is thus a nice example where the geometric politician helps a lot because the underlying dimension of the problem is small.

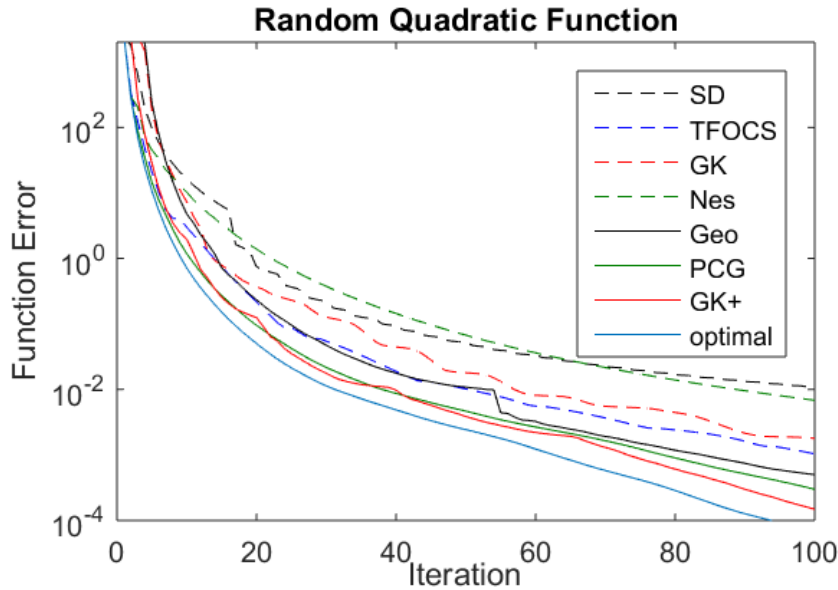


Figure 11-3: Comparison of first-order methods for the function (11.3) with  $n = 10000$ .

#### ■ 11.5.4 Binary Regression with Smoothed Hinge Loss

We consider the binary classification problem on the datasets from [48]. The problem is to minimize the regularized empirical risk:

$$f_t(x) = \frac{1}{n} \sum_{i=1}^n \varphi_t(b_i a_i^T x) + \frac{\lambda}{2} |x|^2 \quad (11.5)$$

where  $a_i \in \mathbb{R}^d$ ,  $b_i \in \mathbb{R}$  are given by the datasets,  $\lambda$  is the regularization coefficient,  $\varphi_t$  is the smoothed hinge loss defined by

$$\varphi_t(z) = \begin{cases} 0 & \text{if } z \leq -1 \\ z + 1 - \frac{t}{2} & \text{if } z \geq -1 + t \\ \frac{1}{2t}(z + 1)^2 & \text{otherwise} \end{cases}$$

and  $t$  is the smoothness parameter. The usual choice for  $t$  is 1, here we test both  $t = 1$  and  $t = 10^{-4}$ . The latter case is to test how well the algorithms perform when the function is non-smooth.

We note that for this problem it would be natural to compare ourselves with SGD (stochastic gradient descent) or more refined stochastic algorithms such as SAG [159] or SVRG [129]. However since the focus of this chapter is on general black-box optimization we stick to comparing only to general methods. It is an interesting open problem to extend our algorithms to the stochastic setting, see Section 11.6.

In figures 11-5 and 11-6, we show the performance profile for problems in the LIBSVM datasets (and with different values for the regularization parameter  $\lambda$ ). More precisely for a given algorithm we plot  $x \in [1, 10]$  versus the fraction of datasets that the algorithm can solve (up to a certain prespecified accuracy) in a number of iterations which is at most  $x$  times the number of iterations of the best algorithm for this dataset. Figure 11-5 shows the case  $t = 1$  with the targeted accuracy  $10^{-6}$ ; Figure 11-6 shows the case  $t = 10^{-4}$  with the targeted accuracy  $10^{-3}$ . We see that TFOCS is slower than SD for many problems, this is simply because SD uses the line search while TFOCS does not, and this makes a huge difference for simple problems. Among algorithms taking  $O(n)$  time per iteration, CG and Geo perform the best, while for the  $O(nk)$  algorithms we see that BFGS, BFGS+

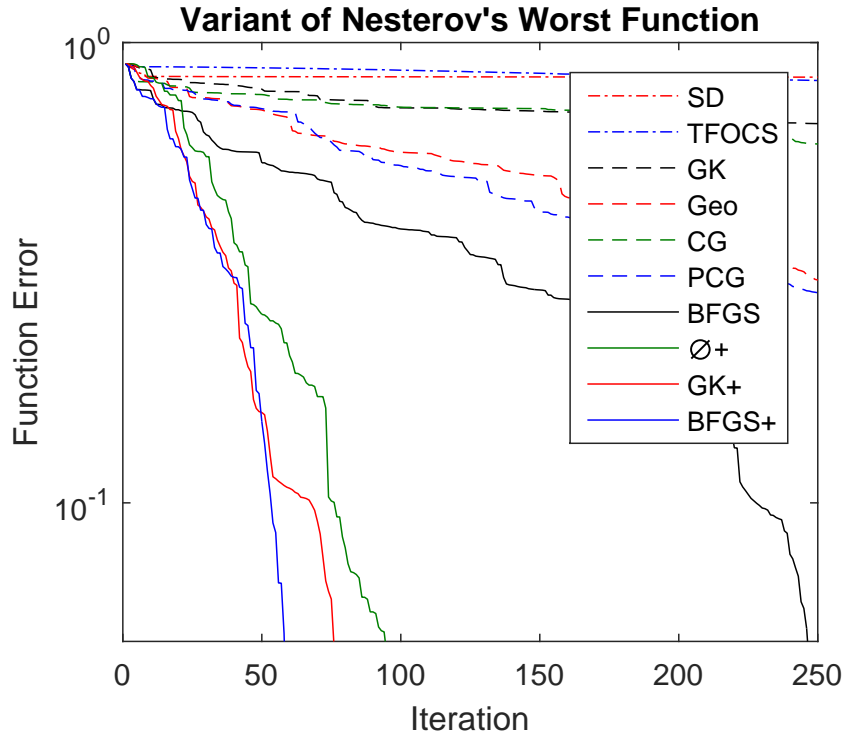


Figure 11-4: Comparison of first-order methods for the function (11.4) with  $n = 10000$ .

and GK+ perform the best. The gap in performance is particularly striking in the non-smooth case where BFGS+ is the fastest algorithm on almost all problems and all other methods (except GK+) are lagging far behind (for 20% of the problems all other methods take 10 times more iterations than BFGS+ and GK+).

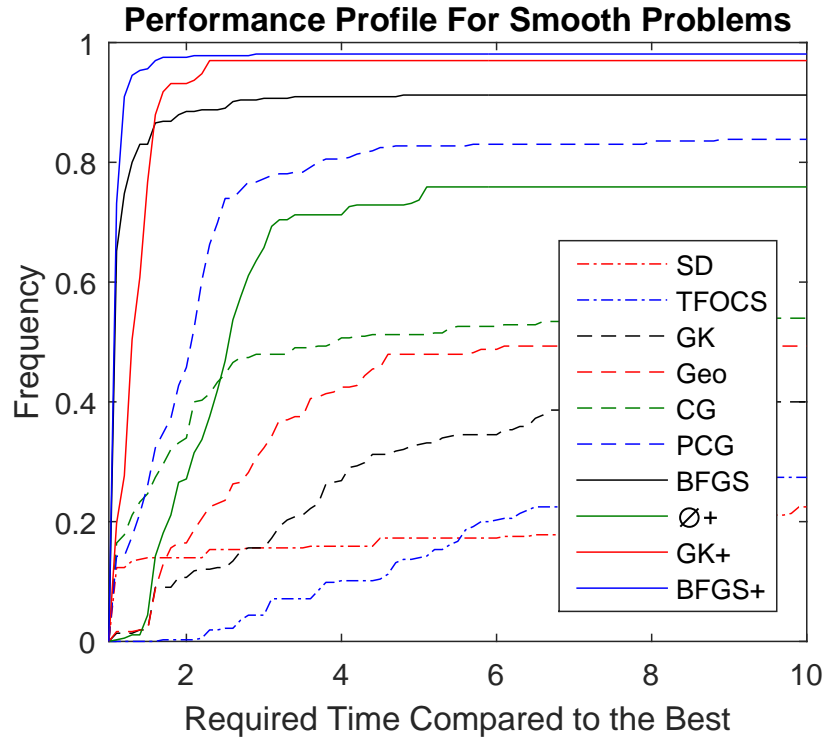
Finally in figures 11-7 and 11-8 we test five algorithms on three specific datasets (respectively in the smooth and non-smooth case). In both figures we see that BFGS+ performs the best for all three datasets. BFGS performs second for smooth problems while GK+ performs second for nonsmooth problems.

### ■ 11.5.5 Summary

The experiments show that BFGS+ and BFGS perform the best among all methods for smooth test problems while BFGS+ and GK+ perform the best for nonsmooth test problems. The first phenomenon is due to the optimality of these algorithm for quadratic problems. We leave the explanation for the second phenomenon as an open problem. At least, the experiments show that this is not due to the geometric oracle itself since  $\emptyset+$  is much slower, and this is not due to the original algorithm since GK performs much worse than GK+ for those problems. Overall these experiments are very promising for the geometric oracle as a replacement of quasi Newton method for non-smooth problems and as a general purpose solver due to its robustness.

## ■ 11.6 Discussion

First order methods generally involve only very basic operations at each step (addition, scalar multiplication). In last section, we formalize each step's operations (besides the gradient calculation) as



**Figure 11-5:** Performance profile on problem (11.5) with  $t = 1$  and  $\lambda = 10^{-4}, 10^{-5}, 10^{-6}, 10^{-7}, 10^{-8}$ .

the work of the politician. We showed that the cost per step of an affine invariant politician  $\psi(k)$  is negligible compared to the gradient calculation (which is  $\Omega(n)$ ). This opens up a lot of possibilities: instead of basic addition or scalar multiplication one can imagine computing a center of gravity, solving a linear program, or even searching over an exponential space (indeed, say  $k < 30$  and  $n > 10^{10}$ , then  $2^k < n$ ). Our experiments demonstrate the effectiveness of this strategy. On the other hand from a theoretical point of view a lot remains to be done. For example one can prove results of the following flavor:

**Theorem 11.6.1.** *Let  $f$  such that  $\alpha I \preceq \nabla^2 f(x) \preceq \beta I, \forall x \in \mathbb{R}^n$  and let  $\kappa = \beta/\alpha$ . Suppose that in the Geometric Politician we replace the volumetric center by the center of gravity or the center of the John ellipsoid. Let  $y_k$  be the output of the  $k^{\text{th}}$  step of SD+ with some initial point  $x_0$ . Then, we have that*

$$f(y_k) - f(x^*) \leq \kappa \left( 1 - \frac{1}{\Theta(\min(n \log(\kappa), \kappa))} \right)^k (f(x_0) - f(x^*)).$$

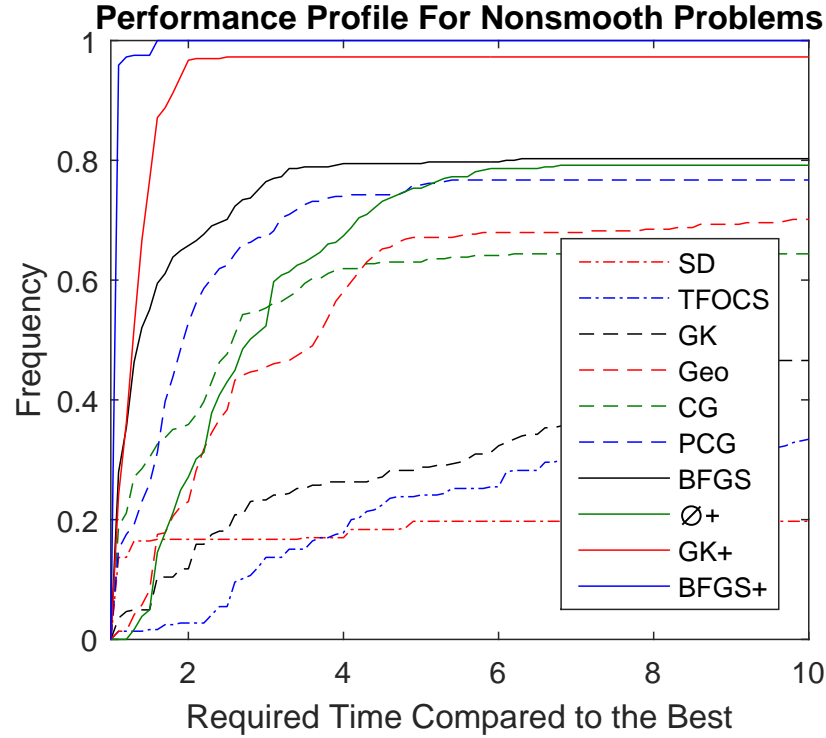
and

$$f(y_k) - f(x^*) \leq \frac{2\beta R^2}{k+4}$$

where  $R = \max_{f(x) \leq f(x_0)} \|x - x^*\|$ .

This claim says that, up to a logarithmic factor, SD+ enjoys simultaneously the incremental progress of gradient descent and the geometrical progress of cutting plane methods. There are three caveats in this claim:

- We use the center of gravity or the center of the John ellipsoid instead of the volumetric center. Note however that it is well-known that the volumetric center is usually more difficult to analyze, [253, 20].



**Figure 11-6:** Performance profile on problem (11.5) with  $t = 10^{-4}$  and  $\lambda = 10^{-4}, 10^{-5}, 10^{-6}, 10^{-7}, 10^{-8}$ .

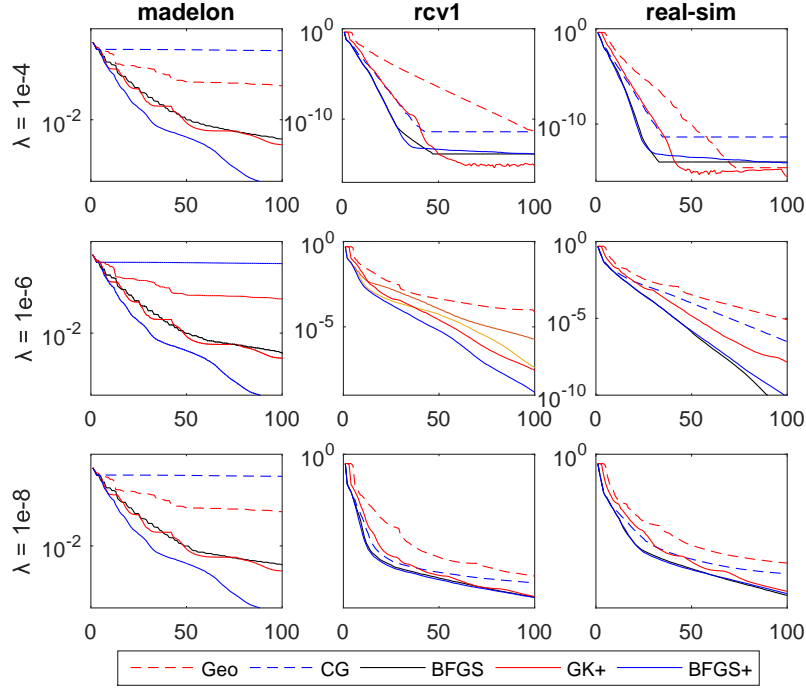
- The extraneous  $\log(\kappa)$  comes from the number of potential restart when we decrease  $\alpha$ . Is there a better way to learn  $\alpha$  that would not incur this additional logarithmic term?
- Theorem 11.3.1 shows essentially that one can combine the ellipsoid method with gradient descent to achieve the optimal  $1 - \sqrt{1/\kappa}$  rate. Can we prove such a result for SD+?

The geometric politician could be refined in many ways. Here are two simple questions that we leave for future work:

- One can think that gradient descent stores 1 gradient information, accelerated gradient descent stores 2 gradient information, and our method stores all past gradient information. We believe that neither 1, 2 nor all is the correct answer. Instead, the algorithm should dynamically decide the number of gradients to store based on the size of its memory, the cost of computing gradients, and the information each gradient reveals.
- Is there a stochastic version of our algorithm? How well would such a method compare with state of the art stochastic algorithms such as SAG [159] and SVRG [129]?

## ■ 11.7 Appendix: Convergence of SD+

Let  $f$  such that  $\alpha I \preceq \nabla^2 f(x) \preceq \beta I$  for all  $x \in \mathbb{R}^n$  and let  $\kappa = \beta/\alpha$ . Let  $y_k$  be the output of the  $k^{\text{th}}$  step of SD+ (where the volumetric center is replaced by the center of gravity or the center of the John ellipsoid) with some initial point  $x_0$ . We prove two rates of convergence for SD+, one with the condition number  $\kappa$ , and one with the ambient dimension  $n$ . We start by the former.



**Figure 11-7:** Comparison between Geo, CG, BFGS, GK+, BFGS+ on problem (11.5) with  $t = 1$  and  $\lambda = 10^{-4}, 10^{-6}, 10^{-8}$ .

**Theorem 11.7.1.** *One has*

$$f(y_k) - f(x^*) \leq \left(1 - \frac{1}{\kappa}\right)^k (f(x_0) - f(x^*))$$

and

$$f(y_k) - f(x^*) \leq \frac{2\beta r^2}{k + 4}$$

where  $r = \max_{f(x) \leq f(x_0)} \|x - x^*\|$ .

*Proof.* Let  $\delta_k = f(y_k) - f(x^*)$ . Since  $f$  is  $\alpha$ -strongly convex we have that

$$\delta_k \leq \frac{1}{2\alpha} \|\nabla f(y_k)\|^2.$$

Due to the decrease guarantee of politicians and the line search in steepest descent, we have that  $f(y_{k+1}) \leq f(x_{k+1}) \leq f(y_k) - \frac{1}{2\beta} \|\nabla f(y_k)\|^2$  and hence

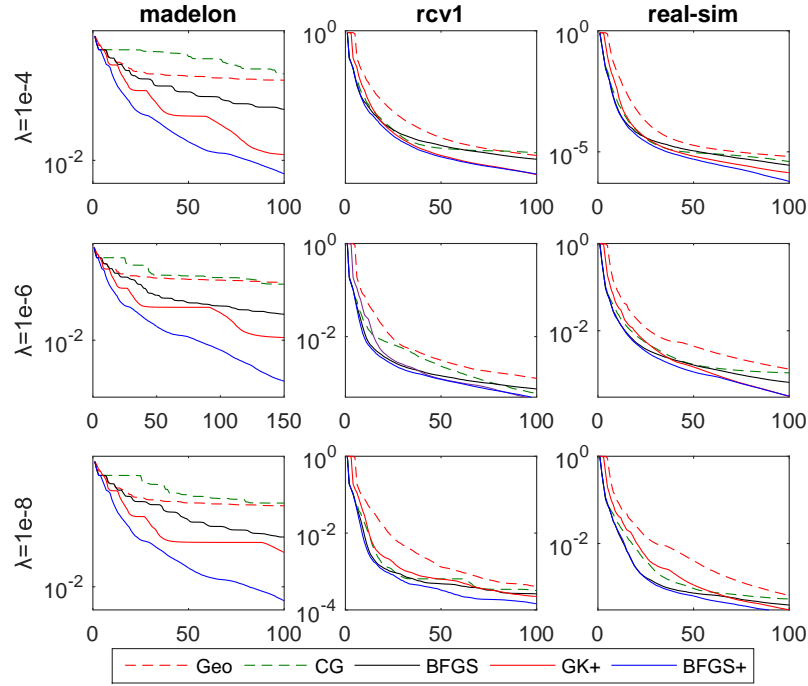
$$\delta_k - \delta_{k+1} \geq \frac{1}{2\beta} \|\nabla f(y_k)\|^2 \geq \frac{\delta_k}{\kappa}. \tag{11.6}$$

Hence, we have  $\delta_{k+1} \leq \left(1 - \frac{1}{\kappa}\right) \delta_k$  and this gives the first inequality.

To obtain a rate independent of  $\alpha$  we instead use the following estimate

$$\delta_k \leq \langle \nabla f(y_k), y_k - x^* \rangle \leq \|\nabla f(y_k)\| \cdot \|y_k - x^*\|.$$

Using the decrease guarantee of politicians and line search we have that  $f(y_k) \leq f(x_k) \leq f(y_{k-1}) \leq$



**Figure 11-8:** Comparison between Geo, CG, BFGS, GK+, BFGS+ on problem (11.5) with  $t = 10^{-4}$  and  $\lambda = 10^{-4}, 10^{-6}, 10^{-8}$ .

$\dots \leq f(x_0)$ , and thus by definition of  $R$ :

$$\|y_k - x^*\| \leq r.$$

Due to the line search in steepest descent again, we have that

$$\delta_k - \delta_{k+1} \geq \frac{1}{2\beta} \|\nabla f(y_k)\|^2 \geq \frac{1}{2\beta} \left(\frac{\delta_k}{r}\right)^2.$$

Since  $\delta_k \geq \delta_{k+1}$ , we have

$$\frac{1}{\delta_{k+1}} - \frac{1}{\delta_k} = \frac{\delta_k - \delta_{k+1}}{\delta_k \delta_{k+1}} \geq \frac{\delta_k - \delta_{k+1}}{\delta_k^2} \geq \frac{1}{2\beta r^2}.$$

So, by induction, we have that  $\frac{1}{\delta_k} \geq \frac{1}{\delta_0} + \frac{k}{2\beta r^2}$ . Now, we note that

$$\delta_0 \leq \langle \nabla f(x^*), x_0 - x^* \rangle + \frac{\beta}{2} \|x_0 - x^*\|^2 \leq \frac{\beta r^2}{2}.$$

Thus, we have that

$$\delta_k \leq \frac{2\beta r^2}{k+4}.$$

□

We now turn to the dimension dependent analysis of SD+. We first show a simple geometric result, namely that if an intersection of spheres has a “small” volume then the intersection must lie close to the boundary of one of the spheres.



**Lemma 11.7.2.** *Let  $R = \cap_{i=1}^k \{x \in \mathbb{R}^n : \|x - c_i\| \leq r_i\}$ ,  $D = \max_{i \in [k]} r_i$ , and  $\omega_n$  the volume of the unit ball in  $\mathbb{R}^n$ . Then, there exists  $i \in [k]$  such that for all  $x \in R$ ,*

$$\|x - c_i\|^2 \geq r_i^2 - 24k^2 \left( \frac{\text{vol}R}{D^n \omega_n} \right)^{1/n} D^2.$$

*Proof.* Since  $-\log(1 - \|x\|^2)$  is a 1-self concordant barrier function,  $2F_R$  is a  $k$ -self concordant function. Let  $y$  be the minimizer of  $F_R$ . Let  $E = \{x \in \mathbb{R}^n : x^\top \nabla^2(2F_R)(y)x \leq 1\}$ . Theorem 4.2.6 in [206] shows that

$$y + E \subset R \subset y + (k + 2\sqrt{k})E. \quad (11.7)$$

In particular, this shows that  $\text{vol}E \leq \text{vol}R$ . We have that

$$(\det \nabla^2 F_R(y))^{1/2} = \frac{\omega_n}{2^{n/2}} \frac{1}{\text{vol}E} \geq \frac{\omega_n}{2^{n/2}} \frac{1}{\text{vol}R}.$$

By the AM-GM inequality, we have that

$$\frac{\text{Tr} \nabla^2 F_R(y)}{n} \geq \frac{1}{2} \left( \frac{\omega_n}{\text{vol}R} \right)^{2/n}. \quad (11.8)$$

By Proposition 11.4.3, we have that  $\nabla^2 F_R(y) = 2A^\top A + \lambda^{(1)}I$  and hence,

$$\begin{aligned} \text{Tr} \nabla^2 F_R(y) &= 2\text{Tr} A^\top A + n\lambda^{(1)} \\ &= 2 \sum_{i=1}^k \frac{\|y - c_i\|^2}{(r_i^2 - \|y - c_i\|^2)^2} + n \sum_{i=1}^k \frac{1}{r_i^2 - \|y - c_i\|^2}. \end{aligned}$$

Applying (11.8), we have that

$$\frac{2}{n} \sum_{i=1}^k \frac{\|y - c_i\|^2}{(r_i^2 - \|y - c_i\|^2)^2} + \sum_{i=1}^k \frac{1}{r_i^2 - \|y - c_i\|^2} \geq \frac{1}{2} \left( \frac{\omega_n}{\text{vol}R} \right)^{2/n}$$

So, there exists  $i$  such that

$$\frac{\|y - c_i\|^2}{(r_i^2 - \|y - c_i\|^2)^2} \geq \frac{n}{8k} \left( \frac{\omega_n}{\text{vol}R} \right)^{2/n} \quad \text{or} \quad \frac{1}{r_i^2 - \|y - c_i\|^2} \geq \frac{1}{4k} \left( \frac{\omega_n}{\text{vol}R} \right)^{2/n}.$$

Using  $\|y - c_i\| \leq r_i \leq D$  and  $\text{vol}R \leq D^n \omega_n$ , we have that

$$\begin{aligned} r_i^2 - \|y - c_i\|^2 &\leq \max \left( \sqrt{\frac{8k}{n}} \left( \frac{\text{vol}R}{\omega_n} \right)^{1/n} \|y - c_i\|, 4k \left( \frac{\text{vol}R}{\omega_n} \right)^{2/n} \right) \\ &\leq 4k \left( \frac{\text{vol}R}{D^n \omega_n} \right)^{1/n} D^2. \end{aligned}$$

Therefore, the width of the ellipsoid  $E$  in the direction  $y - c_i$  is at most

$$r_i - \|y - c_i\| \leq r_i - \sqrt{r_i^2 - 4k \left( \frac{\text{vol}R}{D^n \omega_n} \right)^{1/n} D^2}.$$

The right hand side of (11.7) shows that, for all  $x \in R$ , we have

$$\begin{aligned} \|x - c_i\| &\geq r_i - (1 + k + 2\sqrt{k}) \left( r_i - \sqrt{r_i^2 - 4k \left( \frac{\text{vol}R}{D^n \omega_n} \right)^{1/n} D^2} \right) \\ &\geq r_i - (1 + k + 2\sqrt{k}) \left( r_i - r_i \left( 1 - 4k \left( \frac{\text{vol}R}{D^n \omega_n} \right)^{1/n} \frac{D^2}{r_i^2} \right) \right) \\ &\geq \left[ 1 - 12k^2 \left( \frac{\text{vol}R}{D^n \omega_n} \right)^{1/n} \frac{D^2}{r_i^2} \right] r_i. \end{aligned}$$

Hence, we have

$$\|x - c_i\|^2 \geq \left[ 1 - 24k^2 \left( \frac{\text{vol}R}{D^n \omega_n} \right)^{1/n} \frac{D^2}{r_i^2} \right] r_i^2.$$

□

Finally, equipped with the above geometrical result, we can bound the convergence of SD+ using the dimension  $n$ . We start with a lemma taking care of the adaptivity to the strong convexity in the geometric politician.

**Lemma 11.7.3.** *In the first  $k = \Theta(n \log(\frac{\kappa n}{\varepsilon}))$  iterations, either SD+ restarts the estimate of the strong convexity or*

$$f(y_k) - f(x^*) \leq \varepsilon (f(x_0) - f(x^*)).$$

*Proof.* The decrease guarantee and the smoothness imply that

$$\frac{\|\nabla f(y_k)\|^2}{2\beta} \leq f(y_k) - f(x^*) \leq f(x_0) - f(x^*).$$

Therefore, all the spheres found by the geometric politician have radius squared at most  $D^2$  where, denoting  $\bar{\alpha}$  for the convexity upper bound the algorithm is currently using,

$$D^2 = \max_{k \geq 1} \frac{\|\nabla f(y_k)\|^2}{\bar{\alpha}^2} \leq \frac{2\beta(f(x_0) - f(x^*))}{\bar{\alpha}^2}.$$

Lemma 11.7.2 shows that for any step  $k$ , there is  $i \in [k]$  such that for all  $x \in R_k$ ,

$$\|x - c_i\|^2 \geq r_i^2 - \frac{48\beta k^2}{\bar{\alpha}^2} \left( \frac{\text{vol}R_k}{D^n \omega_n} \right)^{1/n} (f(x_0) - f(x^*)).$$

Let  $k = \Theta(n \log(\frac{\kappa n}{\varepsilon}))$  and recall the discussion in Section 11.4.2 about the volume decrease of the geometric politician with the center of gravity (the same discussion applies to the John ellipsoid). We see that if the algorithm does not restart  $\bar{\alpha}$  within the first  $k$  iterations then we have

$$\frac{\text{vol}R_k}{D^n \omega_n} = \left( O\left(\frac{\varepsilon}{\kappa^2 k^2}\right) \right)^n,$$

and hence (for an appropriate numerical constant in  $k$ )

$$\|x - c_i\|^2 \geq r_i^2 - \frac{\varepsilon(f(x_0) - f(x^*))}{\bar{\alpha}\kappa}. \quad (11.9)$$

Recall from (11.6) that

$$f(y_{k+1}) \leq f(y_k) - \frac{f(y_k) - f(x^*)}{\kappa},$$

and therefore we have (by the improvement of the previous balls):

$$R_{k+1} \subset \left\{ \|x - c_i\|^2 \leq r_i^2 - \frac{2(f(y_k) - f(x^*))}{\bar{\alpha}\kappa} \right\} \cap R_k.$$

However, from (11.9), we know that either the above intersection is empty or  $f(y_k) - f(x^*) < \varepsilon(f(x_0) - f(x^*))$ . This proves the statement.  $\square$

**Theorem 11.7.4.** *We have that*

$$f(y_k) - f(x^*) \leq \kappa \left( 1 - \frac{1}{\Theta(n \log(\kappa))} \right)^k (f(x_0) - f(x^*)).$$

*Proof.* If  $\kappa < n$ , the statements follows from Theorem 11.7.1. Hence, we can assume  $\kappa \geq n$ .

Set  $T = \Theta(n \log(\frac{n\kappa}{\varepsilon}) \log(\kappa))$ , Lemma 11.7.3 shows that for every  $\Theta(n \log(\frac{n\kappa}{\varepsilon}))$  iteration, the algorithm either finds  $y$  such that

$$f(y) - f(x^*) \leq \varepsilon (f(x_0) - f(x^*))$$

or decreases  $\alpha_k$  by a constant where  $\alpha_k$  is the convexity upper bound the algorithm is using at  $k^{\text{th}}$  iteration. Note that  $\alpha_1 \leq \beta$  because of the line search, and thus the algorithm can restart  $\alpha_k$  at most  $\log(\kappa)$  many times. Hence, after  $T$  iterations, we must have

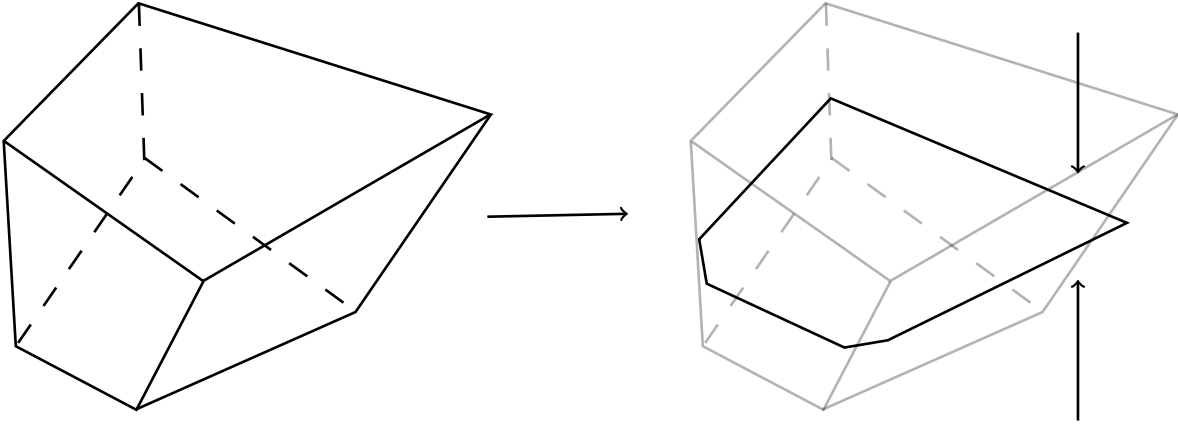
$$f(y_T) - f(x^*) \leq \varepsilon (f(x_0) - f(x^*)),$$

thus concluding the proof.  $\square$

# Part III

---

# Collapsing



---



---

## Undirected MaxFlow In Almost-Linear Time

### ■ 12.1 Introduction

In this chapter, we study the maximum flow problem. Given a graph  $G = (V, E)$  in which each edge  $e \in E$  is assigned a nonnegative capacity  $\mu_e$ , the **maximum  $s$ - $t$  flow problem** asks us to find a flow  $\mathbf{f}$  that routes as much flow as possible from a source vertex  $s$  to a sink vertex  $t$  while sending at most  $\mu_e$  units of flow over each edge  $e$ . Its generalization, the **maximum concurrent multicommodity flow problem**, supplies  $k$  source-sink pairs  $(s_i, t_i)$  and asks for the maximum  $\alpha$  such that we may simultaneously route  $\alpha$  units of flow between each source-sink pair. That is, it asks us to find flows  $\mathbf{f}_1, \dots, \mathbf{f}_k$  (which we think of as corresponding to  $k$  different commodities) such that  $\mathbf{f}_i$  sends  $\alpha$  units of flow from  $s_i$  to  $t_i$ , and  $\sum_i |\mathbf{f}_i(e)| \leq \mu_e$  for all  $e \in E$ .

These problems lie at the core of graph algorithms and combinatorial optimization, and they have been extensively studied over the past 60 years [232, 4]. They have found a wide range of theoretical and practical applications [5], and they are widely used as key subroutines in other algorithms (see [18, 233]).

In this chapter, we introduce a new framework for approximately solving flow problems in capacitated, undirected graphs and apply it to provide asymptotically faster algorithms for the maximum  $s$ - $t$  flow and maximum concurrent multicommodity flow problems. For graphs with  $n$  vertices and  $m$  edges, it allows us to find an  $\varepsilon$ -approximately maximum  $s$ - $t$  flows in time  $O(m^{1+o(1)}\varepsilon^{-2})$ , improving on the previous best bound of  $\tilde{O}(mn^{1/3}\text{poly}(1/\varepsilon))$ [51]. Applying the same framework in the multi-commodity setting solves a maximum concurrent multicommodity flow problem with  $k$  commodities in  $O(m^{1+o(1)}\varepsilon^{-2}k^2)$  time, improving on the existing bound of  $\tilde{O}(m^{4/3}\text{poly}(k, \varepsilon^{-1}))$ [140].

We believe that both our general framework and several of the pieces necessary for its present instantiation are of independent interest, and we hope that they will find other applications. These include:

- a non-Euclidean generalization of gradient descent, bounds on its performance, and a way to use this to reduce approximate maximum flow and maximum concurrent flow to oblivious routing;
- the definition and efficient construction of *flow sparsifiers*; and
- the construction of a new oblivious routing scheme that can be implemented extremely efficiently.

We have aimed to make our algorithm fairly modular and have thus occasionally worked in slightly more generality than is strictly necessary for the problem at hand. This has slightly increased the length of the exposition, but we believe that it clarifies the high-level structure of the argument, and it will hopefully facilitate the application of these tools in other settings.

### ■ 12.1.1 Related Work

For the first several decades of its study, the fastest algorithms for the maximum flow problem were essentially all deterministic algorithms based on combinatorial techniques, such as augmenting paths, blocking flows, preflows, and the push-relabel method. These culminated in the work of Goldberg and Rao [104], which computes exact maximum flows in time  $O(\min(n^{2/3}, m^{1/2}) \log(n^2/m) \log U)$  on graphs with edge weights in  $\{0, \dots, U\}$ . We refer the reader to [104] for a survey of these results.

More recently, a collection of new techniques based on randomization, spectral graph theory and numerical linear algebra, graph decompositions and embeddings, and iterative methods for convex optimization have emerged. These have allowed researchers to provide better provable algorithms for a wide range of flow and cut problems, particularly when one aims to obtain approximately optimal solutions on undirected graphs.

Our algorithm draws extensively on the intellectual heritage established by these works. In this section, we will briefly review some of the previous advances that inform our algorithm. We do not give a comprehensive review of the literature, but instead aim to provide a high-level view of the main tools that motivated the present work, along with the limitations of these tools that had to be overcome. For simplicity of exposition, we primarily focus on the maximum  $s$ - $t$  flow problem for the remainder of the introduction.

**Sparsification** In [33], Benczur and Karger showed how to efficiently approximate any graph  $G$  with a sparse graph  $G'$  on the same vertex set. To do this, they compute a carefully chosen probability  $p_e$  for each  $e \in E$ , sample each edge  $e$  with probability  $p_e$ , and include  $e$  in  $G'$  with its weight increased by a factor of  $1/p_e$  if it is sampled. Using this, they obtain, in nearly linear time, a graph  $G'$  with  $O(n \log n / \varepsilon^2)$  edges such that the total weight of the edges crossing any cut in  $G'$  is within a multiplicative factor of  $1 \pm \varepsilon$  of the weight crossing the corresponding cut in  $G$ . In particular, the Max-Flow Min-Cut Theorem implies that the value of the maximum flow on  $G'$  is within a factor of  $1 \pm \varepsilon$  of that of  $G$ .

This is an extremely effective tool for approximately solving cut problems on a dense graph  $G$ , since one can simply solve the corresponding problem on the sparsified graph  $G'$ . However, while this means that one can approximately compute the *value* of the maximum  $s$ - $t$  flow on  $G$  by solving the problem on  $G'$ , it is not known how to use the maximum  $s$ - $t$  flow on  $G'$  to obtain an actual approximately maximum flow on  $G$ . Intuitively, this is because the weights of edges included in  $G'$  are larger than they were in  $G$ , and the sampling argument does not provide any guidance about how to route flows over these edges in the original graph  $G$ .

**Iterative algorithms based on linear systems and electrical flows** In 2010, Christiano *et al.* [51] described a new linear algebraic approach to the problem that found  $\varepsilon$ -approximately maximum  $s$ - $t$  flows in time  $\tilde{O}(mn^{1/3} \text{poly}(1/\varepsilon))$ . They treated the edges of  $G$  as electrical resistors and then computed the *electrical flow* that would result from sending electrical current from  $s$  to  $t$  in the corresponding circuit. They showed that these flows can be computed in nearly-linear time using fast Laplacian linear system solvers [147, 146, 141, 163], which we further discuss below. The electrical flow obeys the flow conservation constraints, but it could violate the capacity constraints. They then adjusted the resistances of edges to penalize the edges that were flowing too much current and repeated the process. Kelner, Miller, and Peng [140] later showed how to use more general objects that they called *quadratically coupled flows* to use a similar approach to solve the maximum concurrent multicommodity flow problem in time  $\tilde{O}(m^{4/3} \text{poly}(k, 1/\varepsilon))$ .

Following this, Lee, Rao, and Srivastava [161] proposed another iterative algorithm that uses electrical flows, but in a way that was substantially different than in [51]. Instead of adjusting the

resistances of the edges in each iteration to correct overflowing edges, they keep the resistances the same but compute a new electrical flow to reroute the excess current. They explain how to interpret this as gradient descent in a certain space, from which a standard analysis would give an algorithm that runs in time  $\tilde{O}(m^{3/2}\text{poly}(1/\varepsilon))$ . By replacing the standard gradient descent step with Nesterov's accelerated gradient descent method [206] and using a regularizer to make the penalty function smoother, they obtain an algorithm that runs in time  $\tilde{O}(mn^{1/3}\text{poly}(1/\varepsilon))$  in unweighted graphs.

In all of these algorithms, the superlinear running times arise from an intrinsic  $\Theta(\sqrt{m})$  factor introduced by using electrical flows, which minimize an  $\ell_2$  objective function, to approximate the maximum congestion, which is an  $\ell_\infty$  quantity.

**Fast solvers for Laplacian linear systems** In their breakthrough paper [243], Spielman and Teng showed how to solve Laplacian systems in nearly-linear time. (This was later sped up and simplified by Koutis, Miller, and Peng [147, 146] and Kelner, Orecchia, Sidford, and Zhu [141].) Their algorithm worked by showing how to approximate the Laplacian  $\mathcal{L}_G$  of a graph  $G$  with the Laplacian  $\mathcal{L}_H$  of a much simpler graph  $H$  such that one could use the ability to solve linear systems in  $\mathcal{L}_H$  to accelerate the solution of a linear system in  $\mathcal{L}_G$ . They then applied this recursively to solve the linear systems in  $\mathcal{L}_H$ . In addition to providing the electrical flow primitive used by the algorithms described above, the structure of their recursive sequence of graph simplifications provides the motivating framework for much of the technical content of our oblivious routing construction.

**Oblivious routing** In an *oblivious routing scheme*, one specifies a linear operator taking any demand vector to a flow routing these demands over the edges of  $G$ . Given a collection of demand vectors, one can produce a multicommodity flow meeting these demands by routing each demand vector using this pre-specified operator, independently of the others. The competitive ratio of such an operator is the worst possible ratio between the congestion incurred by a set of demands in this scheme and the congestion of the best multicommodity flow routing these demands.

In [225], Räcke showed how to construct an oblivious routing scheme with a competitive ratio of  $O(\log n)$ . His construction worked by providing a probability distribution over trees  $T_i$  such that  $G$  embeds into each  $T_i$  with congestion at most 1, and such that the corresponding convex combination of trees embeds into  $G$  with congestion  $O(\log n)$ . In a sense, one can view this as showing how to approximate  $G$  by a probability distribution over trees. Using this, he was able to show how to obtain polylogarithmic approximations for a variety of cut and flow problems, given only the ability to solve these problems on trees.

We note that such an oblivious routing scheme clearly yields a logarithmic approximation to the maximum flow and maximum concurrent multicommodity flow problems. However, Räcke's construction took time substantially superlinear time, making it too slow to be useful for computing approximately maximum flows. Furthermore, it only gives a logarithmic approximation, and it is not clear how to use this a small number of times to reduce the error to a multiplicative  $\varepsilon$ .

In a later paper [182], Madry applied a recursive technique similar to the one employed by Spielman and Teng in their Laplacian solver to accelerate many of the applications of Räcke's construction at the cost of a worse approximation ratio. Using this, he obtained almost-linear-time polylogarithmic approximation algorithms for a wide variety of cut problems.

Unfortunately, his algorithm made extensive use of sparsification, which, for the previously mentioned reasons, made it unable to solve the corresponding flow problems. This meant that, while it could use flow-cut duality to find a polylogarithmic approximation of the value of a maximum flow, it could not construct a corresponding flow or repeatedly apply such a procedure a small number of times to decrease the error to a multiplicative  $\varepsilon$ .

In simultaneous, independent work [233], Jonah Sherman used somewhat different techniques to find another almost-linear-time algorithm for the (single-commodity) maximum flow problem. His approach is essentially dual to ours: Our algorithm maintains a flow that routes the given demands throughout its execution and iteratively works to improve its congestion. Our main technical tools thus consist of efficient methods for finding ways to route flow in the graph while maintaining flow conservation. Sherman, on the other hand, maintains a flow that does *not* route the given demands, along with a bound on the congestion required to route the excess flow at the vertices. He then uses this to iteratively work towards achieving flow conservation. (In a sense, our algorithm is more in the spirit of augmenting paths, whereas his is more like preflow-push.) As such, his main technical tools are efficient methods for producing dual objects that give congestion bounds. Objects meeting many of his requirements were given in the work of Madry [182] (whereas there were no previous constructions of flow-based analogues, requiring us to start from scratch); leveraging these allows him to avoid some of the technical complexity required by our approach. We believe that these papers nicely complement each other, and we enthusiastically refer the reader to Sherman’s paper.

### ■ 12.1.2 Our Approach

In this section, we give a high-level description of how we overcome the obstacles described in the previous section. For simplicity, we suppose for the remainder of this introduction that all edges have capacity 1.

The problem is thus to send as many units of flow as possible from  $s$  to  $t$  without sending more than one unit over any edge. It will be more convenient for us to work with an equivalent congestion minimization problem, where we try to find the unit  $s$ - $t$  flow  $\mathbf{f}$  (i.e., a flow sending one unit from  $s$  to  $t$ ) that minimizes  $\|\mathbf{f}\|_\infty = \max_e |\mathbf{f}_e|$ . If we begin with some initial unit  $s$ - $t$  flow  $\mathbf{f}_0$ , the goal will be thus be to find the circulation  $\mathbf{c}$  to add to  $\mathbf{f}_0$  that minimizes  $\|\mathbf{f}_0 + \mathbf{c}\|_\infty$ .

We give an iterative algorithm to approximately find such a  $\mathbf{c}$ . There are  $2^{O(\sqrt{\log n \log \log n})}/\varepsilon^2$  iterations, each of which adds a circulation to the present flow and runs in  $m \cdot 2^{O(\sqrt{\log n \log \log n})}$  time. Constructing this scheme consists of two main parts: an iterative scheme that reduces the problem to the construction of a projection matrix with certain properties; and the construction of such an operator.

#### The iterative scheme: Non-Euclidean gradient descent

The simplest way to improve the flow would be to just perform gradient descent on the maximum congestion of an edge. There are two problems with this:

The first problem is that gradient descent depends on having a smoothly varying gradient, but the infinity norm is very far from smooth. This is easily remedied by a standard technique: we replace the infinity norm with a smoother “soft max” function. Doing this would lead to an update that would be a linear projection onto the space of circulations. This could be computed using an electrical flow, and the resulting algorithm would be very similar to the unaccelerated gradient descent algorithm in [161].

The more serious problem is the one discussed in the previous section: the difference between  $\ell_2$  and  $\ell_\infty$ . Gradient steps choose a direction by optimizing a local approximation of the objective function over a sphere, whereas the  $\ell_\infty$  constraint asks us to optimize over a cube. The difference between the size of the largest sphere inside a cube and the smallest sphere containing it gives rise to an inherent  $O(\sqrt{m})$  in the number of iterations, unless one can somehow exploit additional structure in the problem.

To deal with this, we introduce and analyze a non-Euclidean variant of gradient descent that



operates with respect to an arbitrary norm.<sup>1</sup> Rather than choosing the direction by optimizing a local linearization of the objective function over the sphere, it performs an optimization over the unit ball in the given norm. By taking this norm to be  $\ell_\infty$  instead of  $\ell_2$ , we are able to obtain a much smaller bound on the number of iterations, albeit at the expense of having to solve a nonlinear minimization problem at every step. The number of iterations required by the gradient descent method depends on how quickly the gradient can change over balls in the norm we are using, which we express in terms of the Lipschitz constant of the gradient in the chosen norm.

To apply this to our problem, we write flows meeting our demands as  $\mathbf{f}_0 + \mathbf{c}$ , as described above. We then need a parametrization of the space of circulations so that the objective function (after being smoothed using soft max) has a good bound on its Lipschitz constant. Similarly to what occurs in [141], this comes down to finding a good linear representation of the space of circulations, which we show amounts in the present setting to finding a matrix that projects into the space of circulations while meeting certain norm bounds.

### Constructing a projection matrix

This reduces our problem to the construction of such a projection matrix. A simple calculation shows that any linear oblivious routing scheme  $A$  with a good competitive ratio gives rise to a projection matrix with the desired properties, and thus leads to an iterative algorithm that converges in a small number of iterations. Each of these iterations performs a matrix-vector multiplication with both  $A$  and  $A^T$ .

Intuitively, this is letting us replace the electrical flows used in previous algorithms with the flows given by an oblivious routing scheme. Since the oblivious routing scheme was constructed to meet  $\ell_\infty$  guarantees, while the electrical flow could only obtain such guarantees by relating  $\ell_2$  to  $\ell_\infty$ , it is quite reasonable that we should expect this to lead to a better iterative algorithm.

However, the computation involved in existing oblivious routing schemes is not fast enough to be used in this setting. Our task thus becomes constructing an oblivious routing scheme that we can compute and work with very efficiently. We do this with a recursive construction that reduces oblivious routing in a graph to oblivious routing in various successively simpler graphs.

To this end, we show that if  $G$  can be embedded with low congestion into  $H$  (existentially), and  $H$  can be embedded with low congestion into  $G$  *efficiently*, one can use an oblivious routing on  $H$  to obtain an oblivious routing on  $G$ . The crucial difference between the simplification operations we perform here and those in previous papers (e.g., in the work of Benczur-Karger [33] and Madry [182]) is that ours are accompanied by such embeddings, which enables us to transfer flows from the simpler graphs to the more complicated ones.

We construct our routing scheme by recursively composing two types of reductions, each of which we show how to implement without incurring a large increase in the competitive ratio:

- **Vertex elimination** This shows how to efficiently reduce oblivious routing on a graph  $G = (V, E)$  to routing on  $t$  graphs with roughly  $\tilde{O}(|E|/t)$  vertices. To do this, we show how to efficiently embed  $G$  into  $t$  simpler graphs, each consisting of a tree plus a subgraph supported on roughly  $\tilde{O}(|E|/t)$  vertices. This follows easily from a careful reading of Madry's paper [182]. We then show that routing on such a graph can be reduced to routing on a graph with at most  $\tilde{O}(|E|/t)$  vertices by collapsing paths and eliminating leaves.

---

<sup>1</sup>This idea and analysis seems to be implicit in other work, e.g., [205] However, we could not find a clean statement like the one we need in the literature, and we have not seen it previously applied in similar settings. We believe that it will find further applications, so we state it in fairly general terms before specializing to what we need for flow problems.

- **Flow sparsification** This allows us to efficiently reduce oblivious routing on an arbitrary graph to oblivious routing on a graph with  $\tilde{O}(|V|)$  edges, which we call a flow sparsifier. To construct flow sparsifiers, we use local partitioning to decompose the graph into well-connected clusters that contain many of the original edges. (These clusters are not quite expanders, but they are contained in graphs with good expansion in a manner that is sufficient for our purposes.) We then sparsify these clusters using standard techniques and then show that we can embed the sparse graph back into the original graph using electrical flows. If the graph was originally dense, this results in a sparser graph, and we can recurse on the result. While the implementation of these steps is somewhat different, the outline of this construction parallels Spielman and Teng's approach to the construction of spectral sparsifiers [243, 244].

Combining these two reductions recursively yields an efficient oblivious routing scheme, and thus an algorithm for the maximum flow problem.

Finally, we show that the same framework can be applied to the maximum concurrent multicommodity flow problem. While the norm and regularization change, the structure of the argument and the construction of the oblivious routing scheme go through without requiring substantial modification.

## ■ 12.2 Preliminaries

Here, we define some notations that is used throughout in this chapter.

**General Notation:** We let  $\mathbf{I}$  be the identity matrix and  $\mathbf{I}_{a \rightarrow b} \in \mathbb{R}^{b \times a}$  denote the matrix such that for all  $i \leq \min\{a, b\}$  we have  $\mathbf{I}_{ii} = 1$  and  $\mathbf{I}_{ij} = 0$  otherwise.

**Graphs:** Throughout this chapter we let  $G = (V, E, \boldsymbol{\mu})$  denote an undirected capacitated graph with  $n = |V|$  vertices,  $m = |E|$  edges, and non-negative capacities  $\boldsymbol{\mu} \in \mathbb{R}^E$ . We let  $w_e \geq 0$  denote the weight of an edge and let  $r_e \stackrel{\text{def}}{=} 1/w_e$  denote the resistance of an edge. Here we make no connection between  $\mu_e$  and  $r_e$ ; we fix their relationship later.

**Sizes:** For all  $a \in V$  we let  $d_a \stackrel{\text{def}}{=} \sum_{\{a,b\}} w_{a,b}$  denote the (*weighted*) *degree* of vertex  $a$  and we let  $\deg(a) \stackrel{\text{def}}{=} |\{e \in E \mid e = \{a, b\} \text{ for some } b \in V\}|$  denote its (*combinatorial*) *degree*. We let  $\mathbf{D} \in \mathbb{R}^{V \times V}$  be the diagonal matrix where  $\mathbf{D}_{a,a} = d_a$ . Furthermore, for any vertex subset  $S \subseteq V$  we define its *volume* by  $\text{vol}(S) \stackrel{\text{def}}{=} \sum_{a \in V} d_a$ .

**Cuts:** For any vertex subset  $S \subseteq V$  we denote the cut induced by  $S$  by the edge subset

$$\partial(S) \stackrel{\text{def}}{=} \{e \in E \mid e \not\subseteq S \text{ and } e \not\subseteq E \setminus S\}$$

and we denote the cost of  $F \subseteq E$  by  $w(F) \stackrel{\text{def}}{=} \sum_{e \in F} w_e$ . We denote the *conductance* of  $S \subseteq V$  by

$$\Phi(S) \stackrel{\text{def}}{=} \frac{w(\partial(S))}{\min\{\text{vol}(S), \text{vol}(V - S)\}}$$

and we denote the conductance of a graph by

$$\Phi(G) \stackrel{\text{def}}{=} \min_{S \subseteq V : S \neq \emptyset, V} \phi(S)$$

**Subgraphs:** For a graph  $G = (V, E)$  and a vertex subset  $S \subseteq V$  let  $G(S)$  denote the subgraph of  $G$  consisting of vertex set  $S$  and all the edges of  $E$  with both endpoints in  $S$ , i.e.  $\{(a, b) \in E \mid a, b \in S\}$ . When we refer to a graph property, such as  $\text{vol}$  or  $\Phi$ , we use subscripts to specify which graph we are considering. For example  $\text{vol}_{G(S)}(A)$  denotes the volume of vertex set  $A$  in the subgraph of  $G$  induced by  $S$ .

**Congestion:** Thinking of edge vectors,  $\mathbf{f} \in \mathbb{R}^E$ , as flows we let the *congestion*<sup>2</sup> of  $\mathbf{f}$  be given by  $\text{cong}(\mathbf{f}) \stackrel{\text{def}}{=} \|\mathbf{U}^{-1}\mathbf{f}\|_\infty$ . For any collection of flows  $\{\mathbf{f}_i\} = \{\mathbf{f}_1, \dots, \mathbf{f}_k\}$  we overload notation and let their *total congestion* be given by

$$\text{cong}(\{\mathbf{f}_i\}) \stackrel{\text{def}}{=} \left\| \mathbf{U}^{-1} \sum_i \mathbf{f}_i \right\|_\infty$$

**Demands and Multicommodity Flow:** We call a vector  $\boldsymbol{\chi} \in \mathbb{R}^V$  a *demand vector* if it is the case that  $\sum_{a \in V} \boldsymbol{\chi}(a) = 0$  and we say  $\mathbf{f} \in \mathbb{R}^E$  meets the demands if  $\mathbf{B}^T \mathbf{f} = \boldsymbol{\chi}$ . Given a set of demands  $D = \{\boldsymbol{\chi}_1, \dots, \boldsymbol{\chi}_k\}$ , i.e.  $\forall i \in [k], \sum_{a \in V} \boldsymbol{\chi}_i(a) = 0$ , we denote the optimal low congestion routing of these demands as follows

$$\text{OPT}(D) \stackrel{\text{def}}{=} \min_{\{\mathbf{f}_i\} \in \mathbb{R}^E : \{\mathbf{B}^T \mathbf{f}_i\} = \{\boldsymbol{\chi}_i\}} \text{cong}(\{\mathbf{f}_i\})$$

We call a set of flows  $\{\mathbf{f}_i\}$  that meet demands  $\{\boldsymbol{\chi}_i\}$ , i.e.  $\forall i, \mathbf{B}^T \mathbf{f}_i = \boldsymbol{\chi}_i$ , a *multicommodity flow* meeting the demands.

**Running Time:** For matrix  $\mathbf{A}$ , we let  $\mathcal{T}(\mathbf{A})$  denote the maximum amount of time needed to apply  $\mathbf{A}$  or  $\mathbf{A}^T$  to a vector.

## ■ 12.3 Solving Max-Flow Using a Circulation Projection

### ■ 12.3.1 Gradient Descent

In this section, we discuss the gradient descent method for general norms. Let  $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$  be an arbitrary norm on  $\mathbb{R}^n$  and recall that the gradient of  $f$  at  $\mathbf{x}$  is defined to be the vector  $\nabla f(\mathbf{x}) \in \mathbb{R}^n$  such that

$$f(\mathbf{y}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + o(\|\mathbf{y} - \mathbf{x}\|). \quad (12.1)$$

The gradient descent method is a greedy minimization method that updates the current vector,  $\mathbf{x}$ , using the direction which minimizes  $\langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$ . To analyze this method's performance, we need a tool to compare the improvement  $\langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$  with the step size,  $\|\mathbf{y} - \mathbf{x}\|$ , and the quantity,  $\|\nabla f(\mathbf{x})\|$ . For  $\ell_2$  norm, this can be done by Cauchy Schwarz inequality and in general, we can define a new norm for  $\nabla f(\mathbf{x})$  to make this happens. We call this the *dual norm*  $\|\cdot\|^*$  defined as follows

$$\|\mathbf{x}\|^* \stackrel{\text{def}}{=} \max_{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y}\| \leq 1} \langle \mathbf{y}, \mathbf{x} \rangle.$$

Fact 12.8.4 shows that this definition indeed yields that  $\langle \mathbf{y}, \mathbf{x} \rangle \leq \|\mathbf{y}\|^* \|\mathbf{x}\|$ . Next, we define the fastest increasing direction  $\mathbf{x}^\#$ , which is an arbitrary point satisfying the following

$$\mathbf{x}^\# \stackrel{\text{def}}{=} \arg \max_{\mathbf{s} \in \mathbb{R}} \langle \mathbf{x}, \mathbf{s} \rangle - \frac{1}{2} \|\mathbf{s}\|^2.$$

In the appendix, we provide some facts about  $\|\cdot\|^*$  and  $\mathbf{x}^\#$  that we will use in this section. Using the notations defined, the gradient descent method simply produces a sequence of  $\mathbf{x}_k$  such that

$$\mathbf{x}_{k+1} := \mathbf{x}_k - t_k (\nabla f(\mathbf{x}_k))^\#$$

<sup>2</sup>Note that here and in the rest of the chapter we will focus our analysis with congestion with respect to the norm  $\|\cdot\|_\infty$  and we will look at oblivious routing strategies that are competitive with respect to this norm. However, many of the results present are easily generalizable to other norms. These generalizations are outside the scope of this thesis.

where  $t_k$  is some chosen step size for iteration  $k$ . To determine what these step sizes should be we need some information about the smoothness of the function, in particular, the magnitude of the second order term in (12.1). The natural notion of smoothness for gradient descent is the Lipschitz constant of the gradient of  $f$ , that is the smallest constant  $L$  such that

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n \quad : \quad \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^* \leq L \cdot \|\mathbf{x} - \mathbf{y}\|.$$

In the appendix we provide an equivalent definition and a way to compute  $L$ , which is useful later.

Let  $X^* \subseteq \mathbb{R}^n$  denote the set of optimal solutions to the unconstrained minimization problem  $\min_{\mathbf{x} \in \mathbb{R}^n} f$  and let  $f^*$  denote the optimal value of this minimization problem, i.e.

$$\forall \mathbf{x} \in X^* \quad : \quad f(\mathbf{x}) = f^* = \min_{\mathbf{y} \in \mathbb{R}^n} f(\mathbf{y}) \quad \text{and} \quad \forall \mathbf{x} \notin X^* \quad : \quad f(\mathbf{x}) > f^*$$

We assume that  $X^*$  is non-empty. Now, we are ready to estimate the convergence rate of the gradient descent method.

**Theorem 12.3.1.** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a convex continuously differentiable function and let  $L$  be the Lipschitz constant of  $\nabla f$ . For initial point  $\mathbf{x}_0 \in \mathbb{R}^n$  we define a sequence of  $\mathbf{x}_k$  by the update rule*

$$\mathbf{x}_{k+1} := \mathbf{x}_k - \frac{1}{L}(\nabla f(\mathbf{x}_k))^\#$$

For all  $k \geq 0$ , we have

$$f(\mathbf{x}_k) - f^* \leq \frac{2 \cdot L \cdot R^2}{k + 4} \quad \text{where} \quad R \stackrel{\text{def}}{=} \max_{\mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) \leq f(\mathbf{x}_0)} \min_{\mathbf{x}^* \in X^*} \|\mathbf{x} - \mathbf{x}^*\|.$$

*Proof.* <sup>3</sup> By the Lipschitz continuity of the gradient of  $f$  and Lemma 12.8.5 we have

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{1}{2L} (\|\nabla f(\mathbf{x}_k)\|^*)^2.$$

Furthermore, by the convexity of  $f$ , we know that

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n \quad : \quad f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle.$$

Using this and the fact that  $f(\mathbf{x}_k)$  decreases monotonically with  $k$ , we get

$$f(\mathbf{x}_k) - f^* \leq \min_{\mathbf{x}^* \in X^*} \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle \leq \min_{\mathbf{x}^* \in X^*} \|\nabla f(\mathbf{x}_k)\|^* \|\mathbf{x}_k - \mathbf{x}^*\| \leq R \|\nabla f(\mathbf{x}_k)\|^*.$$

Therefore, letting  $\phi_k \stackrel{\text{def}}{=} f(\mathbf{x}_k) - f^*$ , we have

$$\phi_k - \phi_{k+1} \geq \frac{1}{2L} (\|\nabla f(\mathbf{x}_k)\|^*)^2 \geq \frac{\phi_k^2}{2 \cdot L \cdot R^2}.$$

Furthermore, since  $\phi_k \geq \phi_{k+1}$ , we have

$$\frac{1}{\phi_{k+1}} - \frac{1}{\phi_k} = \frac{\phi_k - \phi_{k+1}}{\phi_k \phi_{k+1}} \geq \frac{\phi_k - \phi_{k+1}}{\phi_k^2} \geq \frac{1}{2 \cdot L \cdot R^2}.$$

So, by induction, we have that

$$\frac{1}{\phi_k} - \frac{1}{\phi_0} \geq \frac{k}{2 \cdot L \cdot R^2}.$$

<sup>3</sup>The structure of this specific proof was modeled after a proof in [205] for a slightly different problem.

Now, note that since  $\nabla f(\mathbf{x}^*) = 0$ , we have that

$$f(\mathbf{x}_0) \leq f(\mathbf{x}^*) + \langle \nabla f(\mathbf{x}^*), \mathbf{x}_0 - \mathbf{x}^* \rangle + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \leq f(\mathbf{x}^*) + \frac{L}{2} R^2.$$

So, we have that  $\phi_0 \leq \frac{L}{2} R^2$  and putting this all together yields that

$$\frac{1}{\phi_k} \geq \frac{1}{\phi_0} + \frac{k}{2 \cdot L \cdot R^2} \geq \frac{4}{2 \cdot L \cdot R^2} + \frac{k}{2 \cdot L \cdot R^2}.$$

□

### ■ 12.3.2 Maximum Flow Formulation

For an arbitrary set of demands  $\chi \in \mathbb{R}^V$  we wish to solve the following *maximum flow* problem

$$\max_{\alpha \in \mathbb{R}, \mathbf{f} \in \mathbb{R}^E} \alpha \quad \text{subject to} \quad \mathbf{B}^T \mathbf{f} = \alpha \chi \quad \text{and} \quad \|\mathbf{U}^{-1} \mathbf{f}\|_\infty \leq 1.$$

Equivalently, we want to compute a *minimum congestion flow*

$$\min_{\mathbf{f} \in \mathbb{R}^E : \mathbf{B}^T \mathbf{f} = \chi} \|\mathbf{U}^{-1} \mathbf{f}\|_\infty.$$

where we call  $\|\mathbf{U}^{-1} \mathbf{f}\|_\infty$  the *congestion* of  $\mathbf{f}$ .

Letting  $\mathbf{f}_0 \in \mathbb{R}^E$  be some initial *feasible flow*, i.e.  $\mathbf{B}^T \mathbf{f}_0 \stackrel{\text{def}}{=} \chi$ , we write the problem equivalently as

$$\min_{\mathbf{c} \in \mathbb{R}^E : \mathbf{B}^T \mathbf{c} = 0} \|\mathbf{U}^{-1}(\mathbf{f}_0 + \mathbf{c})\|_\infty$$

where the output flow is  $\mathbf{f} = \mathbf{f}_0 + \mathbf{c}$ . Although the gradient descent method is applicable to constrained optimization problems and has a similar convergence guarantee, the sub-problem involved in each iteration is a constrained optimization problem, which is quite complicated in this case. Since the domain is a linear subspace, the constraints can be avoided by projecting the variables onto this subspace.

Formally, we define a circulation projection matrix as follows.

**Definition 12.3.2.** A matrix  $\tilde{\mathbf{P}} \in \mathbb{R}^{E \times E}$  is a *circulation projection matrix* if it is a projection matrix onto the circulation space, i.e. it satisfies the following

- $\forall \mathbf{x} \in \mathbb{R}^E$  we have  $\mathbf{B}^T \tilde{\mathbf{P}} \mathbf{x} = \mathbf{0}$ .
- $\forall \mathbf{x} \in \mathbb{R}^E$  with  $\mathbf{B}^T \mathbf{x} = \mathbf{0}$  we have  $\tilde{\mathbf{P}} \mathbf{x} = \mathbf{x}$ .

Then, the problem becomes

$$\min_{\mathbf{c} \in \mathbb{R}^E} \|\mathbf{U}^{-1}(\mathbf{f}_0 + \tilde{\mathbf{P}} \mathbf{c})\|_\infty.$$

Applying gradient descent on this problem is similar to applying projected gradient method on the original problem. But, instead of using the orthogonal projection that is not suitable for  $\|\cdot\|_\infty$ , we will pick a better projection matrix.

Applying the change of basis  $\mathbf{x} = \mathbf{U}^{-1} \mathbf{c}$  and letting  $\alpha_0 = \mathbf{U}^{-1} \mathbf{f}_0$  and  $\mathbf{P} = \mathbf{U}^{-1} \tilde{\mathbf{P}} \mathbf{U}$ , we write the problem equivalently as

$$\min_{\mathbf{x} \in \mathbb{R}^E} \|\alpha_0 + \mathbf{P} \mathbf{x}\|_\infty$$

where the output maximum flow is

$$f(\mathbf{x}) = U(\alpha_0 + \mathbf{P}\mathbf{x}) / \|U(\alpha_0 + \mathbf{P}\mathbf{x})\|_\infty.$$

### ■ 12.3.3 An Approximate Maximum Flow Algorithm

Since the gradient descent method requires the objective function to be differentiable, we introduce a smooth version of  $\|\cdot\|_\infty$  which we call  $\text{smax}_t$ . In next section, we prove that there is a convex differentiable function  $\text{smax}_t$  such that  $\nabla \text{smax}_t$  is Lipschitz continuous with Lipschitz constant  $\frac{1}{t}$  and such that

$$\forall \mathbf{x} \in \mathbb{R}^E \quad : \quad \|\mathbf{x}\|_\infty - t \ln(2m) \leq \text{smax}_t(\mathbf{x}) \leq \|\mathbf{x}\|_\infty \quad .$$

Now we consider the following regularized optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^E} g_t(\mathbf{x}) \quad \text{where} \quad g_t(\mathbf{x}) = \text{smax}_t(\alpha_0 + \mathbf{P}\mathbf{x}).$$

For the rest of this section, we consider solving this optimization problem using gradient descent under  $\|\cdot\|_\infty$ .

First, we bound the Lipschitz constant of the gradient of  $g_t$ .

**Lemma 12.3.3.** *The gradient of  $g_t$  is Lipschitz continuous with Lipschitz constant  $L = \frac{\|\mathbf{P}\|_\infty^2}{t}$ .*

*Proof.* By Lemma 12.8.5 and the Lipschitz continuity of  $\nabla \text{smax}_t$ , we have

$$\text{smax}_t(\mathbf{y}) \leq \text{smax}_t(\mathbf{x}) + \langle \nabla \text{smax}_t(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2t} \|\mathbf{y} - \mathbf{x}\|_\infty.$$

Setting  $\mathbf{x} \leftarrow \alpha_0 + \mathbf{P}\mathbf{x}$  and  $\mathbf{y} \leftarrow \alpha_0 + \mathbf{P}\mathbf{y}$ , we have

$$\begin{aligned} g_t(\mathbf{y}) &\leq g_t(\mathbf{y}) + \langle \nabla \text{smax}_t(\alpha_0 + \mathbf{P}\mathbf{x}), \mathbf{P}\mathbf{y} - \mathbf{P}\mathbf{x} \rangle + \frac{1}{2t} \|\mathbf{P}\mathbf{y} - \mathbf{P}\mathbf{x}\|_\infty^2 \\ &\leq g_t(\mathbf{y}) + \langle \mathbf{P}^T \nabla \text{smax}_t(\alpha_0 + \mathbf{P}\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2t} \|\mathbf{P}\|_\infty^2 \|\mathbf{y} - \mathbf{x}\|_\infty^2 \\ &= g_t(\mathbf{y}) + \langle \nabla g_t(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2t} \|\mathbf{P}\|_\infty^2 \|\mathbf{y} - \mathbf{x}\|_\infty^2. \end{aligned}$$

Hence, the result follows from Lemma 12.8.5. □

Now, we apply gradient descent to find an approximate max flow as follows.

---

#### Algorithm 33: MaxFlow

---

**Input:** any initial feasible flow  $\mathbf{f}_0$  and  $\text{OPT} = \min_{\mathbf{x}} \|\mathbf{U}^{-1} \mathbf{f}_0 + \mathbf{P}\mathbf{x}\|_\infty$ .

Let  $\alpha_0 = (\mathbf{I} - \mathbf{P}) \mathbf{U}^{-1} \mathbf{f}_0$  and  $\mathbf{x}_0 = 0$ .

Let  $t = \varepsilon \text{OPT} / 2 \ln(2m)$  and  $k = 300 \|\mathbf{P}\|_\infty^4 \ln(2m) / \varepsilon^2$ .

Let  $g_t = \text{smax}_t(\alpha_0 + \mathbf{P}\mathbf{x})$ .

**for**  $i = 1, 2, \dots, k$  **do**

$\mathbf{x}_{i+1} = \mathbf{x}_i - \frac{t}{\|\mathbf{P}\|_\infty^2} (\nabla g_t(\mathbf{x}_i))^\#$ . (See Lemma 12.3.5)

**end**

**Output:**  $U(\vec{\alpha}_0 + \mathbf{P}\mathbf{x}_k) / \|\vec{\alpha}_0 + \mathbf{P}\mathbf{x}_k\|_\infty$ .

---

We remark that the initial flow can be obtained by BFS and the OPT value can be approximated using binary search. In Section 12.7, we will give an algorithm with better dependence on  $\|\mathbf{P}\|$ .

**Theorem 12.3.4.** Let  $\tilde{\mathbf{P}}$  be a cycle projection matrix, let  $\mathbf{P} = \mathbf{U}^{-1}\tilde{\mathbf{P}}\mathbf{U}$ , and let  $\varepsilon < 1$ . *MaxFlow* outputs an  $(1 - \varepsilon)$ -approximate maximum flow in time

$$O\left(\frac{\|\mathbf{P}\|_\infty^4 \ln(m) (\mathcal{T}(\mathbf{P}) + m)}{\varepsilon^2}\right).$$

*Proof.* First, we bound  $\|\alpha_0\|_\infty$ . Let  $\mathbf{x}^*$  be a minimizer of  $\min_{\mathbf{x}} \|\mathbf{U}^{-1}\mathbf{f}_0 + \mathbf{P}\mathbf{x}\|_\infty$  such that  $\mathbf{P}\mathbf{x}^* = \mathbf{x}^*$ . Then, we have

$$\begin{aligned} \|\alpha_0\|_\infty &= \|\mathbf{U}^{-1}\mathbf{f}_0 - \mathbf{P}\mathbf{U}^{-1}\mathbf{f}_0\|_\infty \\ &\leq \|\mathbf{U}^{-1}\mathbf{f}_0 + \mathbf{x}^*\|_\infty + \|\mathbf{x}^* + \mathbf{P}\mathbf{U}^{-1}\mathbf{f}_0\|_\infty \\ &= \|\mathbf{U}^{-1}\mathbf{f}_0 + \mathbf{x}^*\|_\infty + \|\mathbf{P}\mathbf{x}^* + \mathbf{P}\mathbf{U}^{-1}\mathbf{f}_0\|_\infty \\ &\leq (1 + \|\mathbf{P}\|_\infty) \|\mathbf{U}^{-1}\mathbf{f}_0 + \mathbf{x}^*\|_\infty \\ &= (1 + \|\mathbf{P}\|_\infty) \text{OPT}. \end{aligned}$$

Second, we bound  $R$  in Theorem 12.3.1. Note that

$$g_t(\mathbf{x}_0) = \text{smax}_t(\alpha_0) \leq \|\alpha_0\|_\infty \leq (1 + \|\mathbf{P}\|_\infty) \text{OPT}.$$

Hence, the condition  $g_t(\mathbf{x}) \leq g_t(\mathbf{x}_0)$  implies that

$$\|\alpha_0 + \mathbf{P}\mathbf{x}\|_\infty \leq (1 + \|\mathbf{P}\|_\infty) \text{OPT} + t \ln(2m).$$

For any  $\mathbf{y} \in X^*$  let  $\mathbf{c} = \mathbf{x} - \mathbf{P}\mathbf{x} + \mathbf{y}$  and note that  $\mathbf{P}\mathbf{c} = \mathbf{P}\mathbf{y}$  and therefore  $\mathbf{c} \in X^*$ . Using these facts, we can bound  $R$  as follows

$$\begin{aligned} R &= \max_{\mathbf{x} \in \mathbb{R}^E : g_t(\mathbf{x}) \leq g_t(\mathbf{x}_0)} \left\{ \min_{\mathbf{x}^* \in X^*} \|\mathbf{x} - \mathbf{x}^*\|_\infty \right\} \\ &\leq \max_{\mathbf{x} \in \mathbb{R}^E : g_t(\mathbf{x}) \leq g_t(\mathbf{x}_0)} \|\mathbf{x} - \mathbf{c}\|_\infty \\ &\leq \max_{\mathbf{x} \in \mathbb{R}^E : g_t(\mathbf{x}) \leq g_t(\mathbf{x}_0)} \|\mathbf{P}\mathbf{x} - \mathbf{P}\mathbf{y}\|_\infty \\ &\leq \max_{\mathbf{x} \in \mathbb{R}^E : g_t(\mathbf{x}) \leq g_t(\mathbf{x}_0)} \|\mathbf{P}\mathbf{x}\|_\infty + \|\mathbf{P}\mathbf{y}\|_\infty \\ &\leq 2\|\alpha_0\|_\infty + \|\alpha_0 + \mathbf{P}\mathbf{x}\|_\infty + \|\alpha_0 + \mathbf{P}\mathbf{y}\|_\infty \\ &\leq 2\|\alpha_0\|_\infty + 2\|\alpha_0 + \mathbf{P}\mathbf{x}\|_\infty \\ &\leq 4(1 + \|\mathbf{P}\|_\infty) \text{OPT} + 2t \ln(2m). \end{aligned}$$

From Lemma 12.3.3, we know that the Lipschitz constant of  $\nabla g_t$  is  $\|\mathbf{P}\|_\infty^2/t$ . Hence, Theorem 12.3.1 shows that

$$\begin{aligned} g_t(\mathbf{x}_k) &\leq \min_{\mathbf{x}} g_t(\mathbf{x}) + \frac{2 \cdot L \cdot R^2}{k+4} \\ &\leq \text{OPT} + \frac{2 \cdot L \cdot R^2}{k+4}. \end{aligned}$$

So, we have

$$\begin{aligned} \|\alpha_0 + \mathbf{P}\mathbf{x}_k\|_\infty &\leq g_t(\mathbf{x}_k) + t \ln(2m) \\ &\leq \text{OPT} + t \ln(2m) + \frac{2\|\mathbf{P}\|_\infty^2}{t(k+4)} (4(1 + \|\mathbf{P}\|_\infty) \text{OPT} + 2t \ln(2m))^2. \end{aligned}$$

Using  $t = \varepsilon \text{OPT} / 2 \ln(2m)$  and  $k = 300 \|\mathbf{P}\|_\infty^4 \ln(2m) / \varepsilon^2$ , we have

$$\|\boldsymbol{\alpha}_0 + \mathbf{P}\mathbf{x}_k\|_\infty \leq (1 + \varepsilon) \text{OPT}.$$

Therefore,  $\boldsymbol{\alpha}_0 + \mathbf{P}\mathbf{x}_k$  is an  $(1 - \varepsilon)$  approximate maximum flow.

Now, we estimate the running time. In each step 5, we are required to compute  $(\nabla g(\mathbf{x}_k))^\#$ . The gradient

$$\nabla g(\mathbf{x}) = \mathbf{P}^T \nabla \text{smax}_t(\boldsymbol{\alpha}_0 + \mathbf{P}\mathbf{x})$$

can be computed in  $O(\mathcal{T}(\mathbf{P}) + m)$  using the formula of the gradient of  $\text{smax}_t$ , applications of  $\mathbf{P}$  and  $\mathbf{P}^T$ . Lemma 12.3.5 shows that the  $\#$  operator can be computed in  $O(m)$ .  $\square$

**Lemma 12.3.5.** *In  $\|\cdot\|_\infty$ , the  $\#$  operator is given by the explicit formula*

$$\left(\mathbf{x}^\#\right)_e = \text{sign}(x_e) \|\mathbf{x}\|_1 \quad \text{for } e \in E.$$

*Proof.* Recall that

$$\mathbf{x}^\# = \arg \max_{\mathbf{s} \in \mathbb{R}} \langle \mathbf{x}, \mathbf{s} \rangle - \frac{1}{2} \|\mathbf{s}\|_\infty^2.$$

It is easy to see that for all  $e \in E$ ,  $\|\mathbf{x}^\#\|_\infty = |(\mathbf{x}^\#)_e|$ . In particular, we have

$$\left(\mathbf{x}^\#\right)_e = \text{sign}(x_e) \|\mathbf{x}^\#\|_\infty.$$

Fact 12.8.3 shows that  $\|\mathbf{x}^\#\|_\infty = \|\mathbf{x}\|_1$  and the result follows.  $\square$

### ■ 12.3.4 Properties of soft max

In this section, we define  $\text{smax}_t$  and discuss its properties. Formally, the regularized convex function can be found by smoothing technique using convex conjugate [208] [35, Sec 5.4]. For simplicity and completeness, we define it explicitly and prove its properties directly. Formally, we define

$$\forall \mathbf{x} \in \mathbb{R}^E, \forall t \in \mathbb{R}_{>0} \quad : \quad \text{smax}_t(\mathbf{x}) \stackrel{\text{def}}{=} t \ln \left( \frac{\sum_{e \in E} \exp\left(\frac{x_e}{t}\right) + \exp\left(-\frac{x_e}{t}\right)}{2m} \right).$$

For notational simplicity, for all  $\mathbf{x}$  where this vector is clear from context, we define  $\mathbf{c}$  and  $\mathbf{s}$  as follows

$$\forall e \in E \quad : \quad \mathbf{c}_e \stackrel{\text{def}}{=} \exp\left(\frac{x_e}{t}\right) + \exp\left(-\frac{x_e}{t}\right) \quad \text{and} \quad \mathbf{s}_e \stackrel{\text{def}}{=} \exp\left(\frac{x_e}{t}\right) - \exp\left(-\frac{x_e}{t}\right),$$

where the letters are chosen due to the very close resemblance to hyperbolic sine and hyperbolic cosine.

**Lemma 12.3.6.**

$$\forall \mathbf{x} \in \mathbb{R}^n \quad : \quad \nabla \text{smax}_t(\mathbf{x}) = \frac{1}{\mathbf{1}^T \mathbf{c}} \mathbf{s}$$

$$\forall \mathbf{x} \in \mathbb{R}^n \quad : \quad \nabla^2 \text{smax}_t(\mathbf{x}) = \frac{1}{t(\mathbf{1}^T \mathbf{c})} \left[ \text{diag}(\mathbf{c}) - \frac{\mathbf{s}\mathbf{s}^T}{\mathbf{1}^T \mathbf{c}} \right]$$



*Proof.* For all  $i \in E$  and  $\mathbf{x} \in \mathbb{R}^E$ , we have

$$\begin{aligned} \frac{\partial}{\partial x_i} \text{smax}_t(\mathbf{x}) &= \frac{\partial}{\partial x_i} \left( t \ln \left( \frac{\sum_{e \in E} \exp\left(\frac{x_e}{t}\right) + \exp\left(-\frac{x_e}{t}\right)}{2m} \right) \right) \\ &= \frac{\exp\left(\frac{x_i}{t}\right) - \exp\left(-\frac{x_i}{t}\right)}{\sum_{e \in E} \exp\left(\frac{x_e}{t}\right) + \exp\left(-\frac{x_e}{t}\right)}. \end{aligned}$$

For all  $i, j \in E$  and  $\mathbf{x} \in \mathbb{R}^E$ , we have

$$\begin{aligned} \frac{\partial^2}{\partial x_i \partial x_j} \text{smax}_t(\mathbf{x}) &= \frac{\partial^2}{\partial x_i \partial x_j} \left( t \ln \left( \frac{\sum_{e \in E} \exp\left(\frac{x_e}{t}\right) + \exp\left(-\frac{x_e}{t}\right)}{2m} \right) \right) \\ &= \frac{\partial}{\partial j} \left[ \frac{\exp\left(\frac{x_i}{t}\right) - \exp\left(-\frac{x_i}{t}\right)}{\sum_{e \in E} \exp\left(\frac{x_e}{t}\right) + \exp\left(-\frac{x_e}{t}\right)} \right] \\ &= \frac{1}{t} \frac{(\mathbf{1}^T \mathbf{c}) \mathbf{1}_{i=j}(\mathbf{c}_i) - \mathbf{s}_i \mathbf{s}_j}{(\mathbf{1}^T \mathbf{c})^2}. \end{aligned}$$

□

**Lemma 12.3.7.** *The function  $\text{smax}_t$  is a convex continuously differentiable function and it has Lipschitz continuous gradient with Lipschitz constant  $1/t$  and*

$$\|\mathbf{x}\|_\infty - t \ln(2m) \leq \text{smax}_t(\mathbf{x}) \leq \|\mathbf{x}\|_\infty$$

for  $\mathbf{x} \in \mathbb{R}^E$ .

*Proof.* By the formulation of the Hessian, for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^E$ , we have

$$\mathbf{y}^T (\nabla^2 \text{smax}_t(\mathbf{x})) \mathbf{y} \leq \frac{\sum_i c_i \mathbf{y}_i^2}{t(\mathbf{1}^T \mathbf{c})} \leq \frac{\sum_i c_i (\max_j \mathbf{y}_j^2)}{t(\mathbf{1}^T \mathbf{c})} \leq \frac{1}{t} \|\mathbf{y}\|_\infty^2.$$

On the other side, for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^E$ , we have by  $s_i \leq |s_i| \leq c_i$  and Cauchy Schwarz shows that

$$\mathbf{y}^T \mathbf{s} \mathbf{s}^T \mathbf{y} \leq (\mathbf{1}^T |s|)(\mathbf{y}^T \text{diag}(|s|)\mathbf{y}) \leq (\mathbf{1}^T \mathbf{c})(\mathbf{y}^T \text{diag}(\mathbf{c})\mathbf{y}).$$

and hence

$$0 \leq \mathbf{y}^T (\nabla^2 \text{smax}_t(\mathbf{x})) \mathbf{y}.$$

Thus, the first part follows from Lemma 12.8.6. For the later part, we have

$$\|\mathbf{x}\|_\infty \geq t \ln \left( \frac{\sum_{e \in E} \exp\left(\frac{x_e}{t}\right) + \exp\left(-\frac{x_e}{t}\right)}{2m} \right) \geq t \ln \left( \frac{\exp\left(\frac{\|\mathbf{x}\|_\infty}{t}\right)}{2m} \right) = \|\mathbf{x}\|_\infty - \ln(2m).$$

□

## ■ 12.4 Oblivious Routing

In the previous sections, we saw how a circulation projection matrix can be used to solve max flow. In the next few sections, we show how to efficiently construct a circulation projection matrix to obtain an almost linear time algorithm for solving max flow.

Our proof focuses on the notion of (linear) oblivious routings. Rather than constructing the circulation projection matrix directly, we show how the efficient construction of an oblivious routing

algorithm with a good competitive ratio immediately allows us to produce a circulation projection matrix.

In the remainder of this section, we formally define oblivious routings and prove the relationship between oblivious routing and circulation projection matrices (Section 12.4.1), provide a high level overview of our recursive approach and state the main theorems we will prove in later sections (Section 12.4.2). Finally, we prove the main theorem about our almost-linear-time construction of circulation projection with norm  $2^{O(\sqrt{\log(n)} \log \log(n))}$  assuming the proofs in the later sections (Section 12.4.3).

### ■ 12.4.1 From Oblivious Routing to Circulation Projection

Here we provide definitions and prove basic properties of *oblivious routings*, that is, fixed mappings from demands to flows that meet the input demands. While non-linear algorithms could be considered, we restrict our attention to linear oblivious routing strategies and use the term oblivious routing to refer to the linear subclass for the remainder of this chapter.<sup>4</sup>

**Definition 12.4.1** (Oblivious Routing). An *oblivious routing* on graph  $G = (V, E)$  is a linear operator  $\mathbf{A} \in \mathbb{R}^{E \times V}$  such that for all demands  $\chi \in \mathbb{R}^V$ ,  $\mathbf{B}^T \mathbf{A} \chi = \chi$ . We call  $\mathbf{A} \chi$  the *routing* of  $\chi$  by  $\mathbf{A}$ .

Given an oblivious routing strategy  $\mathbf{A}$  and a set of demands  $D = \{\chi_1, \dots, \chi_k\}$ , one can construct a multicommodity flow satisfying all the demands in  $D$  by using  $\mathbf{A}$  to route each demand individually, obliviously to the existence of the other demands. We measure the *competitive ratio* of such an oblivious routing strategy to be the ratio of the worst relative congestion of such a routing to the minimal-congestion routing of the demands.

**Definition 12.4.2** (Competitive Ratio). The *competitive ratio of oblivious routing*  $\mathbf{A} \in \mathbb{R}^{E \times V}$ , denoted  $\rho(\mathbf{A})$ , is given by

$$\rho(\mathbf{A}) \stackrel{\text{def}}{=} \max_{\{\chi_i\} : \forall i \chi_i \perp \mathbf{1}} \frac{\text{cong}(\{\mathbf{A} \chi_i\})}{\text{OPT}(\{\chi_i\})}$$

At times, it will be more convenient to analyze an oblivious routing as a linear algebraic object rather a combinatorial algorithm; towards this end, we note that the competitive ratio of a linear oblivious routing strategy can be gleaned from the operator norm of a related matrix (see also [158] and [139]). Below, we state and prove a generalization of this result to weighted graphs that will be vital to relating  $\mathbf{A}$  to  $\bar{\mathbf{P}}$ .

**Lemma 12.4.3.** For any oblivious routing  $\mathbf{A}$ , we have  $\rho(\mathbf{A}) = \|\mathbf{U}^{-1} \mathbf{A} \mathbf{B}^T \mathbf{U}\|_\infty$

*Proof.* For a set of demands  $D$ , let  $D_\infty$  be the set of demands that results by taking the routing of every demand in  $D$  by  $\text{OPT}(D)$  and splitting it up into demands on every edge corresponding to the flow sent by  $\text{OPT}(D)$ . Now, clearly  $\text{OPT}(D) = \text{OPT}(D_\infty)$  since routing  $D$  can be used to route  $D_\infty$  and vice versa, and clearly  $\text{cong}(\mathbf{A} D) \leq \text{cong}(\mathbf{A} D_\infty)$  by the linearity of  $\mathbf{A}$  (routing  $D_\infty$  simply doesn't reward  $\mathbf{A}$  routing for cancellations). Using this and the fact that  $D_\infty \subseteq D$ , we have:

$$\begin{aligned} \rho_p(\mathbf{A}) &= \max_D \frac{\text{cong}(\{\mathbf{A} D\})}{\text{OPT}(D)} = \max_{D_\infty} \frac{\text{cong}(\mathbf{A} D_\infty)}{\text{OPT}(D_\infty)} = \max_{x \in \mathbb{R}^E} \frac{\|\sum_{e \in E} x_e |\mathbf{U}^{-1} \mathbf{A} \chi_e|\|_\infty}{\|\mathbf{U}^{-1} x\|_\infty} \\ &= \max_{x \in \mathbb{R}^E} \frac{\|\mathbf{U}^{-1} \mathbf{A} \mathbf{B}^T x\|_\infty}{\|\mathbf{U}^{-1} x\|_\infty} = \max_{x \in \mathbb{R}^E} \frac{\|\mathbf{U}^{-1} \mathbf{A} \mathbf{B}^T \mathbf{U} x\|_\infty}{\|x\|_\infty}. \end{aligned}$$

□

<sup>4</sup>Note that the oblivious routing strategies considered in [139] [158] [225] are all linear oblivious routing strategies.

To make this lemma easily applicable in a variety of settings, we make use of the following easy-to-prove lemma.

**Lemma 12.4.4** (Operator Norm Bounds). *For all  $\mathbf{A} \in \mathbb{R}^{n \times m}$ , we have that*

$$\|\mathbf{A}\|_\infty = \|\mathbf{A}\|_\infty = \|\mathbf{A} \mathbf{1}\|_\infty = \max_{i \in n} \|\mathbf{A}^T \mathbf{1}_i\|_1.$$

The previous two lemmas make the connection between oblivious routings and circulation projection matrices clear. Below, we prove it formally.

**Lemma 12.4.5** (Oblivious Routing to Circulation Projection). *For oblivious routing  $\mathbf{A} \in \mathbb{R}^{E \times V}$  the matrix  $\tilde{\mathbf{P}} \stackrel{\text{def}}{=} \mathbf{I} - \mathbf{A}\mathbf{B}^T$  is a circulation projection matrix such that  $\|\mathbf{U}^{-1}\tilde{\mathbf{P}}\mathbf{U}\|_\infty \leq 1 + \rho(\mathbf{A})$ .*

*Proof.* First, we verify that  $\text{im}(\tilde{\mathbf{P}})$  is contained in cycle space:

$$\forall \mathbf{x} \in \mathbb{R}^E \quad : \quad \mathbf{B}^T \tilde{\mathbf{P}}\mathbf{x} = \mathbf{B}^T \mathbf{x} - \mathbf{B}^T \mathbf{A}\mathbf{B}^T \mathbf{x} = \mathbf{0}.$$

Next, we check that  $\tilde{\mathbf{P}}$  is the identity on cycle space

$$\forall \mathbf{x} \in \mathbb{R}^E \text{ s.t. } \mathbf{B}^T \mathbf{x} = \mathbf{0} \quad : \quad \tilde{\mathbf{P}}\mathbf{x} = \mathbf{x} - \mathbf{A}\mathbf{B}^T \mathbf{x} = \mathbf{x}.$$

Finally, we bound the  $\ell_\infty$ -norm of the scaled projection matrix:

$$\|\mathbf{U}^{-1}\tilde{\mathbf{P}}\mathbf{U}\|_\infty = \|\mathbf{I} - \mathbf{U}^{-1}\mathbf{A}\mathbf{B}^T\mathbf{U}\|_\infty \leq 1 + \rho(\mathbf{A}).$$

□

## ■ 12.4.2 A Recursive Approach by Embeddings

We construct an oblivious routing for a graph recursively. Given a generic, possibly complicated, graph, we show how to reduce computing an oblivious routing on this graph to computing an oblivious routing on a simpler graph on the same vertex set. A crucial concept in these constructions will be the notion of an embedding, which will allow us to relate the competitive ratios of an oblivious routing algorithms over graphs on the same vertex sets but different edge sets.

**Definition 12.4.6** (Embedding). Let  $G = (V, E, \boldsymbol{\mu})$  and  $G' = (V, E', \boldsymbol{\mu}')$  denote two undirected capacitated graphs on the same vertex set with incidence matrices  $\mathbf{B} \in \mathbb{R}^{E \times V}$  and  $\mathbf{B}' \in \mathbb{R}^{E' \times V}$  respectively. An *embedding* from  $G$  to  $G'$  is a matrix  $\mathbf{M} \in \mathbb{R}^{E' \times E}$  such that  $\mathbf{B}'^T \mathbf{M} = \mathbf{B}^T$ .

In other words, an embedding is a map from flows in one graph  $G$  to flows in another graph  $G'$  that preserves the demands met by the flow. We can think of an embedding as a way of routing any flow in graph  $G$  into graph  $G'$  that has the same vertex set, but different edges. We will be particularly interested in embeddings that increase the congestion of the flow by a small amount going from  $G$  to  $G'$ .

**Definition 12.4.7** (Embedding Congestion). Let  $\mathbf{M} \in \mathbb{R}^{E' \times E}$  be an embedding from  $G = (V, E, \boldsymbol{\mu})$  to  $G' = (V, E', \boldsymbol{\mu}')$  and let  $\mathbf{U} \in \mathbb{R}^{E \times E}$  and  $\mathbf{U}' \in \mathbb{R}^{E' \times E'}$  denote the capacity matrices of  $G$  and  $G'$  respectively. The *congestion* of embedding  $\mathbf{M}$  is given by

$$\text{cong}(\mathbf{M}) \stackrel{\text{def}}{=} \max_{\mathbf{x} \in \mathbb{R}^E} \frac{\|\mathbf{U}'^{-1} \mathbf{M} \mathbf{x}\|_\infty}{\|\mathbf{U}^{-1} \mathbf{x}\|_\infty} = \|\mathbf{U}'^{-1} |\mathbf{M}| \mathbf{U} \mathbf{1}\|_\infty .$$

We say  $G$  *embeds into*  $G'$  with congestion  $\alpha$  if there exists an embedding  $\mathbf{M}$  from  $G$  to  $G'$  such that  $\text{cong}(\mathbf{M}) \leq \alpha$ .

Embeddings potentially allow us to reduce computing an oblivious routing in a complicated graph to computing an oblivious routing in a simpler graph. Specifically, if we can embed a complicated graph in a simpler graph and we can efficiently embed the simple graph in the original graph, both with low congestion, then we can just focus on constructing oblivious routings in the simpler graph. We prove this formally as follows.

**Lemma 12.4.8** (Embedding Lemma). *Let  $G = (V, E, \mu)$  and  $G' = (V, E', \mu')$  denote two undirected capacitated graphs on the same vertex sets, let  $M \in \mathbb{R}^{E' \times E}$  denote an embedding from  $G$  into  $G'$ , let  $M' \in \mathbb{R}^{E \times E'}$  denote an embedding from  $G'$  into  $G$ , and let  $A' \in \mathbb{R}^{E' \times V}$  denote an oblivious routing algorithm on  $G'$ . Then  $A \stackrel{\text{def}}{=} M' A'$  is an oblivious routing algorithm on  $G$  and*

$$\rho(A) \leq \text{cong}(M) \cdot \text{cong}(M') \cdot \rho(A').$$

*Proof.* For all  $x \in \mathbb{R}^V$  we have by definition of embeddings and oblivious routings that

$$B^T A x = B^T M' A' x = B^T x.$$

To bound  $\rho(A)$ , we let  $U$  denote the capacity matrix of  $G$  and  $U'$  denote the capacity matrix of  $G'$ . Using Lemma 12.4.3, we get

$$\rho(A) = \|U^{-1} A B^T U\|_\infty = \|U^{-1} M' A' B^T U\|_\infty$$

Using that  $M$  is an embedding and therefore  $B^T M = B^T$ , we get

$$\rho(A) = \|U^{-1} M' A' B^T M U\|_\infty \leq \|U^{-1} M' U'\|_\infty \cdot \|U'^{-1} A' B^T U'\|_\infty \cdot \|U'^{-1} M U\|_\infty$$

By the definition of competitive ratio and congestion, we obtain the result.  $\square$

Note how in this lemma we only use the embedding from  $G$  to  $G'$  to certify the quality of flows in  $G'$ , we do not actually need to apply this embedding in the reduction.

Using this concept, we construct oblivious routings via recursive application of two techniques. First, in Section 12.5, we show how to take an arbitrary graph  $G = (V, E)$  and approximate it by a *sparse graph*  $G' = (V, E')$  (i.e. one in which  $|E'| = \tilde{O}(|V|)$ ) such that flows in  $G$  can be routed in  $G'$  with low congestion and that there is an  $\tilde{O}(1)$  embedding from  $G'$  to  $G$  that can be applied in  $\tilde{O}(|E|)$  time. We call such a construction a *flow sparsifiers* and prove the following theorem.

**Theorem 12.4.9.** *Let  $G = (V, E, \mu)$  be an undirected capacitated graph. In  $\tilde{O}(|E| \log U)$  time we can construct a graph  $G'$  on the same vertex set with at most  $\tilde{O}(|V| \log U)$  edges and capacity ratio at most  $U \cdot \text{poly}(|V|)$ . Moreover, given an oblivious routing  $A'$  on  $G'$ , in  $\tilde{O}(|E| \log U)$  time we can construct an oblivious routing  $A$  on  $G$  such that*

$$\mathcal{T}(A) = \tilde{O}(|E| \log U + \mathcal{T}(A')) \quad \text{and} \quad \rho(A) = \tilde{O}(\rho(A')).$$

Next, in Section 12.6 we show how to embed a graph into a collection of graphs consisting of trees plus extra edges. Then, we will show how to embed these graphs into better structured graphs consisting of trees plus edges so that by simply removing degree 1 and degree 2 vertices we are left with graphs with fewer vertices. Formally, we prove the following.

**Theorem 12.4.10.** *Let  $G = (V, E, \mu)$  be an undirected capacitated graph with capacity ratio  $U$ . For all  $t > 0$  in  $\tilde{O}(t \cdot |E|)$  time we can compute graphs  $G_1, \dots, G_t$  each with at most  $\tilde{O}(\frac{|E| \log(U)}{t})$  vertices, at most  $|E|$  edges, and capacity ratio at most  $|V| \cdot U$ . Moreover, given oblivious routings  $A_i$  for each*

$G_i$ , in  $\tilde{O}(t \cdot |E|)$  time we can compute an oblivious routing  $\mathbf{A}$  on  $G$  such that

$$\mathcal{T}(\mathbf{A}) = \tilde{O}(t \cdot |E|) + \sum_{i=1}^t \mathcal{T}(\mathbf{A}_i) \quad \text{and} \quad \rho(\mathbf{A}) = \tilde{O}(\max_i \rho(\mathbf{A}_i)).$$

In the next section we show that the careful application of these two ideas along with a powerful primitive for routing on constant sized graphs suffices to produce an oblivious routing with the desired properties.

### ■ 12.4.3 Efficient Oblivious Routing Construction Proof

First, we provide the lemma that will serve as the base case of our recursion. In particular, we show that electric routing can be used to obtain a routing algorithm with constant competitive ratio for constant-size graphs.

**Lemma 12.4.11** (Base Case). *Let  $G = (V, E, \mu)$  be an undirected capacitated graph and let us assign weights to edges so that  $\mathbf{W} = \mathbf{U}^2$ . For  $\mathcal{L} \stackrel{\text{def}}{=} \mathbf{B}^T \mathbf{W} \mathbf{B}$  we have that  $\mathbf{A} \stackrel{\text{def}}{=} \mathbf{W} \mathbf{B} \mathcal{L}^+$  is an oblivious routing on  $G$  with  $\rho(\mathbf{A}) \leq \sqrt{|E|}$  and  $\mathcal{T}(\mathcal{L}^+) = \tilde{O}(|E|)$ .*

*Proof.* To see that  $\mathbf{A}$  is an oblivious routing strategy we note that for any demands  $\chi \in \mathbb{R}^V$  we have  $\mathbf{B}^T \mathbf{A} = \mathcal{L} \chi = \chi$ . To see bound  $\rho(\mathbf{A})$  we note that by Lemma 12.4.3 and standard norm inequalities we have

$$\rho(\mathbf{A}) = \max_{\mathbf{x} \in \mathbb{R}^E} \frac{\|\mathbf{U}^{-1} \mathbf{W} \mathbf{B} \mathcal{L}^+ \mathbf{B}^T \mathbf{U} \mathbf{x}\|_\infty}{\|\mathbf{x}\|_\infty} \leq \max_{\mathbf{x} \in \mathbb{R}^E} \frac{\|\mathbf{U} \mathbf{B} \mathcal{L}^+ \mathbf{B}^T \mathbf{U} \mathbf{x}\|_2}{\frac{1}{\sqrt{|E|}} \|\mathbf{x}\|_2} = \sqrt{|E|} \cdot \|\mathbf{U} \mathbf{B} \mathcal{L}^+ \mathbf{B}^T \mathbf{U}\|_2$$

The result follows from the fact in [242] that  $\mathbf{\Pi} \stackrel{\text{def}}{=} \mathbf{U} \mathbf{B} \mathcal{L}^+ \mathbf{B}^T \mathbf{U}$  is an orthogonal projection, and therefore  $\|\mathbf{\Pi}\|_2 \leq 1$ , and the fact in [243, 147, 146, 141] that  $\mathcal{T}(\mathcal{L}^+) = \tilde{O}(|E|)$ .  $\square$

Assuming Theorem 12.4.9 and Theorem 12.4.10, which we prove in the next two sections, we prove that low-congestion oblivious routings can be constructed efficiently.

**Theorem 12.4.12.** *Given an undirected capacitated graph  $G = (V, E, \mu)$  with capacity ratio  $U$ . Assume  $U = \text{poly}(|V|)$ . We can construct an oblivious routing algorithm  $\mathbf{A}$  on  $G$  in time*

$$O(|E| 2^{O(\sqrt{\log |V| \log \log |V|})})$$

such that

$$\mathcal{T}(\mathbf{A}) = |E| 2^{O(\sqrt{\log |V| \log \log |V|})} \quad \text{and} \quad \rho(\mathbf{A}) = 2^{O(\sqrt{\log |V| \log \log |V|})}.$$

*Proof.* Let  $c$  be the constant hidden in the exponent terms, including  $\tilde{O}(\cdot)$  and  $\text{poly}(\cdot)$  in Theorem 12.4.9 and Theorem 12.4.10. Apply Theorem 12.4.9 to construct a sparse graph  $G^{(1)}$ , then apply Theorem 12.4.10 with  $t = \lceil 2\sqrt{\log |V| \log \log |V|} \rceil$  to get  $t$  graphs  $G_1^{(1)}, \dots, G_t^{(1)}$  such that each graphs have at most  $O(\frac{1}{t}|E| \log^{2c} |V| \log^{2c} U)$  vertices and at most  $U \cdot |V|^{2c}$  capacity ratio.

Repeat this process on each  $G_i^{(1)}$ , it produces  $t^2$  graphs  $G_1^{(2)}, \dots, G_{t^2}^{(2)}$ . Keep doing this until all graphs  $G_i$  produced have  $O(1)$  vertices. Let  $k$  be the highest level we go through in this process. Since at the  $k$ -th level the number of vertices of each graph is at most  $O(\frac{c^k}{t^k} |E| \log^{2kc} |V| \log^{2k}(U|V|^{2ck}))$

vertices, we have  $k = O\left(\sqrt{\frac{\log |V|}{\log \log |V|}}\right)$ .

On each graph  $G_i$ , we use Theorem 12.4.11 to get an oblivious routing algorithm  $\mathbf{A}_i$  for each  $G_i$  with

$$\mathcal{T}(\mathbf{A}_i) = \tilde{O}(1) \quad \text{and} \quad \rho(\mathbf{A}_i) = \tilde{O}(1).$$

Then, Theorem 12.4.10 and 12.4.9 shows that we have an oblivious routing algorithm  $\mathbf{A}$  for  $G$  with

$$\mathcal{T}(\mathbf{A}) = O(tk|E| \log^{ck}(|V|) \log^{2k}(U|V|^{2ck})) \quad \text{and} \quad \rho(\mathbf{A}) = O(\log^{2kc} |V| \log^k(U|V|^{2ck})).$$

The result follows from  $k = O\left(\sqrt{\frac{\log |V|}{\log \log |V|}}\right)$  and  $t = \left\lceil 2\sqrt{\log |V| \log \log |V|} \right\rceil$ .  $\square$

Using Theorem 12.4.12, Lemma 12.4.5 and Theorem 12.3.4, we have the following almost linear time max flow algorithm on undirected graph.

**Theorem 12.4.13.** *Given an undirected capacitated graph  $G = (V, E, \mu)$  with capacity ratio  $U$ . Assume  $U = \text{poly}(|V|)$ . There is an algorithm finds an  $(1 - \varepsilon)$  approximate maximum flow in time*

$$O\left(\frac{|E|2^{O\left(\sqrt{\log |V| \log \log |V|}\right)}}{\varepsilon^2}\right).$$

## ■ 12.5 Flow Sparsifiers

In order to prove Theorem 12.4.9, i.e. reduce the problem of efficiently computing a competitive oblivious routing on a dense graph to the same problem on a sparse graph, we introduce a new algorithmic tool called *flow sparsifiers*.<sup>5</sup> A flow sparsifier is an efficient cut-sparsification algorithm that also produces an efficiently-computable low-congestion embedding mapping the sparsified graph back to the original graph.

**Definition 12.5.1** (Flow Sparsifier). An algorithm is a  $(h, \varepsilon, \alpha)$ -flow sparsifier if on input graph  $G = (V, E, \mu)$  with capacity ratio  $U$  it outputs a graph  $G' = (V, E', \mu')$  with capacity ratio  $U' \leq U \cdot \text{poly}(|V|)$  and an embedding  $\mathbf{M} : \mathbb{R}^{E'} \rightarrow \mathbb{R}^E$  of  $G'$  into  $G$  with the following properties:

- **Sparsity:**  $G'$  is  $h$ -sparse, i.e.

$$|E'| \leq h$$

- **Cut Approximation:**  $G'$  is an  $\varepsilon$ -cut approximation of  $G$ , i.e.

$$\forall S \subseteq V \quad : \quad (1 - \varepsilon)\mu(\partial_G(S)) \leq \mu'(\partial_{G'}(S)) \leq (1 + \varepsilon)\mu(\partial_G(S))$$

- **Flow Approximation:**  $\mathbf{M}$  has congestion at most  $\alpha$ , i.e.

$$\text{cong}(\mathbf{M}) \leq \alpha.$$

- **Efficiency:** The algorithm runs in  $\tilde{O}(m \log U)$  time and  $\mathcal{T}(\mathbf{M})$  is also  $\tilde{O}(m \log U)$ .

Flow sparsifiers allow us to solve a multi-commodity flow problem on a possibly dense graph  $G$  by converting  $G$  into a sparse graph  $G'$  and solving the flow problem on  $G'$ , while suffering a loss of a factor of at most  $\alpha$  in the congestion when mapping the solution back to  $G$  using  $\mathbf{M}$ .

<sup>5</sup>Note that our flow sparsifiers aim to reduce the number of edges, and are different from the flow sparsifiers of Leighton and Moitra [167], which work in a different setting and reduce the number of vertices.

**Theorem 12.5.2.** *Consider a graph  $G = (V, E, \mu)$  and let  $G' = (V, E', \mu')$  be given by an  $(h, \varepsilon, \alpha)$ -flow sparsifier of  $G$ . Then, for any set of  $k$  demands  $D = \{\chi_1, \chi_2, \dots, \chi_k\}$  between vertex pairs of  $V$ , we have:*

$$\text{OPT}_{G'}(D) \leq \frac{O(\log k)}{1 - \varepsilon} \cdot \text{OPT}_G(D). \quad (12.2)$$

Given the optimum flow  $\{f_i^*\}$  over  $G'$ , we have

$$\text{cong}_G(\{\mathbf{M}f_i^*\}) \leq \alpha \cdot \text{OPT}_{G'}(D) \leq \frac{O(\alpha \log k)}{1 - \varepsilon} \cdot \text{OPT}_G(D).$$

*Proof.* By the flow-cut gap theorem of Aumann and Rabani [21], we have that, for any set of  $k$  demands  $D$  on  $V$  we have:

$$\text{OPT}_G(D) \geq O\left(\frac{1}{\log k}\right) \cdot \max_{S \subseteq V} \frac{D(\partial(S))}{\mu(\partial_G(S))}.$$

where  $D(\partial(S))$  denotes the total amount of demand separated by the cut between  $S$  and  $\bar{S}$ . As any cut  $S \subseteq V$  in  $G'$  has capacity  $\mu'(\partial_{G'}(S)) \geq (1 - \varepsilon)\mu(\partial_G(S))$ , we have:

$$\text{OPT}_{G'}(D) \leq \max_{S \subseteq V} \frac{D(\partial(S))}{\mu'(\partial_{G'}(S))} \leq \frac{1}{1 - \varepsilon} \cdot \max_{S \subseteq V} \frac{D(\partial(S))}{\mu(\partial_G(S))} \leq \frac{O(\log k)}{1 - \varepsilon} \cdot \text{OPT}_G(D).$$

The second part of the theorem follows as a consequence of the definition of the congestion of the embedding  $\mathbf{M}$ .  $\square$

Our flow sparsifiers should be compared with the cut-based decompositions of Räcke [225]. Räcke constructs a probability distribution over trees and gives explicit embeddings from  $G$  to this distribution and backwards, achieving a congestion of  $O(\log n)$ . However, this distribution over tree can include up to  $O(m \log n)$  trees and it is not clear how to use it to obtain an almost linear time algorithm. Flow sparsifiers answer this problem by embedding  $G$  into a single graph  $G'$ , which is larger than a tree, but still sparse. Moreover, they provide an explicit efficient embedding of  $G'$  into  $G$ . Interestingly, the embedding from  $G$  to  $G'$  is not necessary for our notion of flow sparsifier, and is replaced by the cut-approximation guarantee. This requirement, together with the application of the flow-cut gap [21], lets us argue that the optimal congestion of a  $k$ -commodity flow problem can change at most by a factor of  $O(\log k)$  between  $G$  and  $G'$ .

### Main Theorem on Flow Sparsifiers and Proof of Theorem 12.4.9

The main goal of this section will be to prove the following theorem:

**Theorem 12.5.3.** *For any constant  $\varepsilon \in (0, 1)$ , there is an  $(\tilde{O}(n \log U), \varepsilon, \tilde{O}(1))$ -flow sparsifier.*

Assuming Theorem 12.5.3, we can now prove Theorem 12.4.9, the main theorem necessary for edge reduction in our construction of low-congestion projections.

*Proof of Theorem 12.4.9.* We apply the flow sparsifier of Theorem 12.5.3 to  $G = (V, E, \mu)$  and obtain output  $G' = (V, E', \mu')$  with embedding  $\mathbf{M}$ . By the definition of flow sparsifier, we know that the capacity ratio  $U'$  of  $G'$  is at most  $U \cdot \text{poly}(|V|)$ , as required. Moreover, again by Theorem 12.5.3,  $G'$  has at most  $\tilde{O}(|V| \log U)$  edges. Given an oblivious routing  $\mathbf{A}'$  on  $G'$  consider the oblivious routing  $\mathbf{A} \stackrel{\text{def}}{=} \mathbf{M}\mathbf{A}'$ . By the definition of flow sparsifier, we have that  $\mathcal{T}(\mathbf{M}) = \tilde{O}(|E| \log U)$ . Hence  $\mathcal{T}(\mathbf{A}) = \mathcal{T}(\mathbf{M}) + \mathcal{T}(\mathbf{A}') = \tilde{O}(|E| \log U) + \mathcal{T}(\mathbf{A}')$ . To complete the proof, we bound the competitive ratio  $\rho(\mathbf{A})$ . Using the same argument as in Lemma 12.4.3, we can write  $\rho(\mathbf{A})$  as

$$\rho(\mathbf{A}) = \max_D \frac{\text{cong}_G(\{\mathbf{A}D\})}{\text{OPT}_G(D)} \leq \max_{D_\infty} \frac{\text{cong}_G(\mathbf{A}D_\infty)}{\text{OPT}_G(D_\infty)},$$

where  $D_\infty$  is the set of demands that result by taking the routing of every demand in  $D$  by  $\text{OPT}(D)$  and splitting it up into demands on every edge corresponding to the flow sent by  $\text{OPT}(D)$ . Notice that  $D_\infty$  has at most  $|E|$  demands that are routed between pairs of vertices in  $V$ . Then, because  $G'$  is an  $\varepsilon$ -cut approximation to  $G$ , the flow-cut gap of Aumann and Rabani [21] guarantees that

$$\text{OPT}_G(D_\infty) \geq \frac{1}{O(\log n)} \text{OPT}_{G'}(D_\infty).$$

As a result, we obtain:

$$\begin{aligned} \rho(\mathbf{A}) &\leq O(\log n) \cdot \max_{D_\infty} \frac{\text{cong}_G(\mathbf{A}D_\infty)}{\text{OPT}_{G'}(D_\infty)} = O(\log n) \cdot \max_{D_\infty} \frac{\text{cong}_G(\mathbf{M}\mathbf{A}'D_\infty)}{\text{OPT}_{G'}(D_\infty)} \\ &\leq O(\log n) \cdot \text{cong}(\mathbf{M}) \cdot \max_{D_\infty} \frac{\text{cong}_{G'}(\mathbf{A}'D_\infty)}{\text{OPT}_{G'}(D_\infty)} \leq \tilde{O}(\rho(\mathbf{A}')). \end{aligned}$$

□

## Techniques

We will construct flow sparsifiers by taking as a starting point the construction of spectral sparsifiers of Spielman and Teng [244]. Their construction achieves a sparsity of  $\tilde{O}\left(\frac{n}{\varepsilon^2}\right)$  edges, while guaranteeing an  $\varepsilon$ -spectral approximation. As the spectral approximation implies the cut approximation, the construction in [244] suffices to meet the first two conditions in Definition 12.5.1. Moreover, their algorithm also runs in time  $\tilde{O}(m)$ , meeting the fourth condition. Hence, to complete the proof of Theorem 12.5.3, we will modify the construction of Spielman and Teng to endow their sparsifier  $G'$  with an embedding  $\mathbf{M}$  onto  $G$  of low congestion that can be both computed and invoked efficiently. The main tool we use in constructing  $\mathbf{M}$  is the notion of electrical-flow routing and the fact that electrical-flow routing schemes achieve a low competitive ratio on near-expanders and subsets thereof [139, 158].

To exploit this fact and construct a flow sparsifier, we follow Spielman and Teng [244] and partition the input graph into vertex sets, where each sets induces a near-expanders and most edges of the graph do not cross set boundaries. We then sparsify these induced subgraphs using standard sparsification techniques and iterate on the edges not in the subgraphs. As each iteration removes a constant fraction of the edges, by using standard sparsification techniques, we immediately obtain the sparsity and cut approximation properties. To obtain the embedding  $\mathbf{M}$  with  $\text{cong}(\mathbf{M}) = \tilde{O}(1)$ , we prove a generalization of results in [139, 158] and show that the electrical-flow routing achieves a low competitive ratio on near-expanders and subsets thereof.

In the next two subsections, we introduce the necessary concept about electrical-flow routing and prove that it achieves low competitive ratio over near-expanders (and subsets of near-expanders).

### ■ 12.5.1 Subgraph Routing

Given an oblivious routing strategy  $\mathbf{A}$ , we may be interested only in routing demands coming from a subset of edge  $F \subseteq E$ . In this setting, given a set of demands  $D$  routable in  $F$ , we let  $\text{OPT}^F(D)$  denote the minimal congestion achieved by any routing restricted to only sending flow on edges in  $F$  and we measure the  $F$ -competitive ratio of  $\mathbf{A}$  by

$$\rho^F(\mathbf{A}) \stackrel{\text{def}}{=} \max_{D \text{ routable in } F} \frac{\text{cong}(\mathbf{A}D)}{\text{OPT}^F(D)}$$

Note that  $\mathbf{A}$  may use all the edges in  $G$  but  $\rho^F(\mathbf{A})$  compares it only against routings that are restricted to use only edges in  $F$ . As before, we can upper bound the  $F$ -competitive ratio  $\rho^F(\mathbf{A})$  by operator



norms.

**Lemma 12.5.4.** *Let  $\mathbf{1}_F \in \mathbb{R}^E$  denote the indicator vector for set  $F$  (i.e.  $\mathbf{1}_F(e) = 1$  if  $e \in F$  and  $\mathbf{1}_F(e) = 0$ ) and let  $\mathbf{I}_F \stackrel{\text{def}}{=} \text{diag}(\mathbf{1}_F)$ . For any  $F \subseteq E$  we have*

$$\rho^F(\mathbf{A}) = \|\mathbf{U}^{-1}\mathbf{A}\mathbf{B}^T\mathbf{U}\mathbf{I}_F\|_\infty$$

*Proof.* We use the same reasoning as in the non-subgraph case. For a set of demands  $D = \{\chi_i\}$ , we consider  $D_\infty^F$ , the demands on the edges in  $F$  used by  $\text{OPT}^F(D)$ . Then, it is the case that  $\text{OPT}^F(D) = \text{OPT}^F(D_\infty)$  and we know that cost of obviously routing  $D_P$  is greater than the cost of obviously routing  $D$ . Therefore, we have:

$$\begin{aligned} \rho^F &= \max_{\mathbf{x} \in \mathbb{R}^E : \mathbf{I}_{E \setminus F} \mathbf{x} = 0} \frac{\|\sum_{e \in E} |\mathbf{U}^{-1}\mathbf{A}\mathbf{B}^T \mathbf{1}_e \mathbf{x}_e|\|_\infty}{\|\mathbf{U}^{-1}\mathbf{x}\|_\infty} \\ &= \max_{\mathbf{y} \in \mathbb{R}^E : \mathbf{I}_{E \setminus F} \mathbf{y} = 0} \frac{\|\sum_{e \in E} |\mathbf{U}^{-1}\mathbf{A}\mathbf{B}^T \mathbf{U} \mathbf{1}_e \mathbf{y}_e|\|_\infty}{\|\mathbf{y}\|_\infty} \\ &= \max_{\mathbf{y} \in \mathbb{R}^E} \frac{\|\sum_{e \in E} |\mathbf{U}^{-1}\mathbf{A}\mathbf{B}^T \mathbf{U} \mathbf{I}_F \mathbf{1}_e \mathbf{y}_e|\|_\infty}{\|\mathbf{y}\|_\infty} = \|\mathbf{U}^{-1}\mathbf{A}\mathbf{B}^T \mathbf{U} \mathbf{I}_F\|_\infty = \|\mathbf{U}^{-1}\mathbf{A}\mathbf{B}^T \mathbf{U} \mathbf{I}_F\|_\infty \end{aligned}$$

(Having  $y_e \neq 0$  for  $e \in E \setminus F$  decreases the ratio.)

□

## ■ 12.5.2 Electrical-Flow Routings

In this section, we define the notion of electrical-flow routing and prove the results necessary to construct flow sparsifiers. Recall that  $\mathbf{R}$  is the diagonal matrix of resistances and the Laplacian  $\mathcal{L}$  is defined as  $\mathbf{B}^T \mathbf{R}^{-1} \mathbf{B}$ . For the rest of this section, we assume that resistances are set as  $\mathbf{R} = \mathbf{U}^{-1}$ .

**Definition 12.5.5.** Consider a graph  $G = (V, E, \mu)$  and set the edge resistances as  $r_e = \frac{1}{\mu_e}$  for all  $e \in E$ . The oblivious electrical-flow routing strategy is the linear operator  $\mathbf{A}_\mathcal{E}$  defined as

$$\mathbf{A}_\mathcal{E} \stackrel{\text{def}}{=} \mathbf{R}^{-1} \mathbf{B} \mathcal{L}^+,$$

In words, the electrical-flow routing strategy is the routing scheme that, for each demand  $\chi$  sends the electrical flow with boundary condition  $\chi$  on the graph  $G$  with resistances  $\mathbf{R} = \mathbf{U}^{-1}$ .

For the electrical-flow routing strategy  $\mathbf{A}_\mathcal{E}$ , the upper bound on the competitive ratio  $\rho(\mathbf{A}_\mathcal{E})$  in Lemma 12.4.3 can be rephrased in terms of the voltages induced on  $G$  by electrically routing an edge  $e \in E$ . This interpretation appears in [139, 158].

**Lemma 12.5.6.** *Let  $\mathbf{A}_\mathcal{E}$  be the electrical-flow routing strategy. For an edge  $e \in E$ , we let the voltage vector  $\mathbf{v}_e \in \mathbb{R}^V$  be given by  $\mathbf{v}_e \stackrel{\text{def}}{=} \mathcal{L}^+ \chi_e$ . We then have*

$$\rho(\mathbf{A}_\mathcal{E}) = \max_{e \in E} \sum_{(a,b) \in E} \frac{|v_e(a) - v_e(b)|}{r_{ab}}.$$

*Proof.* We have:

$$\rho(\mathbf{A}_\mathcal{E}) = \|\mathbf{B} \mathcal{L}^+ \mathbf{B}^T \mathbf{R}^{-1}\|_\infty = \max_{e \in E} \|\mathbf{R}^{-1} \mathbf{B} \mathcal{L}^+ \mathbf{B}^T \mathbf{1}_e\|_1 = \max_{e \in E} \sum_{(a,b) \in E} \frac{|v_e(a) - v_e(b)|}{r_{ab}}.$$

□

The same reasoning can be extended to the subgraph-routing case to obtain the following lemma.

**Lemma 12.5.7.** *For  $F \subseteq E$  and  $\mathbf{R} = \mathbf{U}^{-1}$  we have*

$$\rho^F(\mathbf{A}_{\mathcal{E}}) = \max_{e \in E} \sum_{(a,b) \in F} \frac{|v_e(a) - v_e(b)|}{r_{ab}}.$$

*Proof.* As before, we have:

$$\begin{aligned} \rho^F(\mathbf{A}_{\mathcal{E}}) &= \|\mathbf{B}\mathcal{L}^+ \mathbf{B}^T \mathbf{R}^{-1} \mathbf{I}_F\|_{\infty} && \text{(By Lemma 12.5.4)} \\ &= \max_{e \in E} \|\mathbf{I}_F \mathbf{R}^{-1} \mathbf{B}\mathcal{L}^+ \mathbf{B}^T \mathbf{1}_e\|_1 = \max_{e \in E} \sum_{(a,b) \in F} \frac{|v_e(a) - v_e(b)|}{r_{ab}} \end{aligned}$$

□

### 12.5.2.1 Bounding the Congestion

In this section, we prove that we can bound the  $F$ -competitive ratio of the oblivious electrical-routing strategy as long as the edges  $F$  that the optimum flow is allowed to route over are contained within an induced expander  $G(U) = (U, E(U))$  for some  $U \subseteq V$ . Towards this we provide and prove the following lemma. This is a generalization of a similar lemma proved in [139].

**Lemma 12.5.8.** *For weighted graph  $G = (V, E, w)$  with integer weights and vertex subset  $U \subseteq V$  the following holds:*

$$\rho^F(\mathbf{A}_{\mathcal{E}}) \leq \frac{8 \log(\text{vol}(G(U)))}{\Phi(G(U))^2}$$

*Proof.* By Lemma 12.5.7, for every edge  $e \in E$ ,

$$\rho^F(\mathbf{A}_{\mathcal{E}}) \leq \|\mathbf{I}_{E(U)} \mathbf{R}^{-1} \mathbf{B}\mathcal{L}^+ \chi_e\|_1$$

Fix any edge  $e \in E$  and let  $v \stackrel{\text{def}}{=} \mathcal{L}^+ \chi_e$ . Recall that with this definition

$$\|\mathbf{I}_{E(U)} \mathbf{R}^{-1} \mathbf{B}\mathcal{L}^+ \chi_e\|_1 = \sum_{(a,b) \in E(U)} \frac{|v(a) - v(b)|}{r_{ab}} = \sum_{(a,b) \in E(U)} w_{ab} \cdot |v(a) - v(b)| \quad (12.3)$$

We define the following vertex subsets:

$$\forall x \in \mathbb{R} \quad : \quad S_x^{\leq} \stackrel{\text{def}}{=} \{a \in U \mid v(a) \leq x\} \quad \text{and} \quad S_x^{\geq} \stackrel{\text{def}}{=} \{a \in U \mid v(a) \geq x\}$$

Since adding a multiple of the all-ones vector to  $v$  does not change the quantity of interest in Equation 12.3, we can assume without loss of generality that

$$\text{vol}_{G(U)}(S_0^{\geq}) \geq \frac{1}{2} (\text{vol}(G(U))) \quad \text{and} \quad \text{vol}_{G(U)}(S_0^{\leq}) \geq \frac{1}{2} (\text{vol}(G(U))).$$

For any vertex subset  $S \subseteq U$ , we denote the flow out of  $S$  and the weight out of  $S$  by

$$f(S) \stackrel{\text{def}}{=} \sum_{e=(a,b) \in E(U) \cap \partial(S)} w_e |v(a) - v(b)|, \quad \text{and} \quad w(S) \stackrel{\text{def}}{=} \sum_{e \in E(U) \cap \partial(S)} w_e.$$

At this point, we define a collections of subsets  $\{C_i \in S_0^{\geq}\}$ . For an increasing sequence of real numbers

$\{c_i\}$ , we let  $C_i \stackrel{\text{def}}{=} S_{c_i}^{\geq}$  and we define the sequence  $\{c_i\}$  inductively as follows:

$$c_0 \stackrel{\text{def}}{=} 0 \quad , \quad c_i \stackrel{\text{def}}{=} c_{i-1} + \Delta_{i-1} \quad , \quad \text{and} \quad \Delta_i \stackrel{\text{def}}{=} 2 \frac{f(C_i)}{w(C_i)} \quad .$$

In words, the  $c_{i+1}$  equals the sum of  $c_i$  and an increase  $\Delta_i$  which depends on how much the cut  $\delta(C_i) \cap E(U)$  was congested by the electrical flow.

Now,  $l_i \stackrel{\text{def}}{=} w(\partial_{E(U)}(C_{i-1}) - \partial_{E(U)}(C_i))$ , i.e. the weight of the edges in  $E(U)$  cut by  $C_{i-1}$  but not by  $C_i$ . We get

$$\begin{aligned} \text{vol}(C_{i+1}) &\leq \text{vol}(C_i) - l_i \\ &\leq \text{vol}(C_i) - \frac{w(C_i)}{2} && \text{(By choice of } l_i \text{ and } \Delta_i) \\ &\leq \text{vol}(C_i) - \frac{1}{2} \text{vol}(C_i) \Phi(G(U)) && \text{(Definition of conductance)} \end{aligned}$$

Applying this inductively and using our assumption on  $\text{vol}(S_0^{\geq})$  we have that

$$\text{vol}(C_i) \leq \left(1 - \frac{1}{2} \Phi(G(U))\right)^i \text{vol}(C_0) \leq \frac{1}{2} \left(1 - \frac{1}{2} \Phi(G(U))\right)^i \text{vol}(G(U))$$

Since  $\phi(G(U)) \in (0, 1)$ , for  $i + 1 = \frac{2 \log(\text{vol}(G(U)))}{\Phi(G(U))}$  we have that  $\text{vol}(S_i) \leq \frac{1}{2}$ . Since  $\text{vol}(S_i)$  decreases monotonically with  $i$ , if we let  $r$  be the smallest value such that  $C_{r+1} = \emptyset$ , we must have

$$r \leq \frac{2 \cdot \log(\text{vol}(G(U)))}{\Phi(G(U))}$$

Since  $v$  corresponds to a unit flow, we know that  $f(C_i) \leq 1$  for all  $i$ . Moreover, by the definition of conductance we know that  $w(C_i) \geq \Phi(G(U)) \cdot \text{vol}(C_i)$ . Therefore,

$$\Delta_i \leq \frac{2}{\Phi(G(U)) \cdot \text{vol}(C_i)}.$$

We can now bound the contribution of  $C_0^{\geq}$  to the volume of the linear embedding  $v$ . In the following, for a vertex  $a \in V$ , we let  $d(a) \stackrel{\text{def}}{=} \sum_{e=\{a,b\} \in E(U)} w_e$  be the degree of  $a$  in  $E(U)$ .

$$\begin{aligned} \sum_{a \in C_0^{\geq}} d(a)v(a) &= \sum_{i=0}^r \left[ \sum_{a \in C_i - C_{i+1}} d(a)v(a) \right] \\ &\leq \sum_{i=0}^r \left[ \sum_{a \in C_i - C_{i+1}} d(a) \left( \sum_{j=0}^i \Delta_j \right) \right] && \text{(By definition of } C_i) \\ &\leq \sum_{i=0}^r \left[ (\text{vol}(C_i) - \text{vol}(C_{i+1})) \cdot \left( \sum_{j=0}^i \Delta_j \right) \right] \\ &= \sum_{i=0}^r \text{vol}(C_i) \Delta_i \leq \frac{2r}{\Phi(G(U))} && \text{(Rearrangement and fact that } \text{vol}(C_{r+1}) = 0) \end{aligned}$$

By repeating the same argument on  $S_0^{\leq}$ , we get that  $\sum_{a \in S_0^{\leq}} d(a)v(a) \leq \frac{2r}{\Phi(G(U))}$ . Putting this all

together yields

$$\|\mathbf{I}_{E(U)}\mathbf{R}^{-1}\mathbf{B}\mathcal{L}^+\chi_e\| = \sum_{(a,b)\in G(U)} w_{ab} \cdot |v(a) - v(b)| \leq \sum_{a\in G(U)} d(a)v(a) \leq \frac{4r}{\Phi(G(U))}$$

□

From this lemma and Lemma 12.5.7, the following is immediate:

**Lemma 12.5.9.** *Let  $F \subseteq E$  be contained within some vertex induced subgraph  $G(U)$ , then for  $\mathbf{R} = \mathbf{U}^{-1}$  we have*

$$\rho^F(\mathbf{R}^{-1}\mathbf{B}\mathcal{L}^+) \leq \rho^{E(U)}(\mathbf{R}^{-1}\mathbf{B}\mathcal{L}^+) \leq \frac{8 \log(\text{vol}(G(U)))}{\Phi(G(U))^2}.$$

### ■ 12.5.3 Construction and Analysis of Flow Sparsifiers

In the remainder of this section we show how to produce an efficient  $O(\log^c)$ -flow sparsifier for some fixed constant  $c$ , proving Theorem 12.5.3. In this chapter, we make no attempt to optimize the value of  $c$ . For the rest of this section, we again assume that we choose the resistance of an edge to be the inverse of its capacity, i.e.  $\mathbf{U} = \mathbf{W} = \mathbf{R}^{-1}$ .

As discussed before, our approach follows closely that of Spielman and Teng [244] to the construction of spectral sparsifiers. The first step of this line of attack is to reduce the problem to the unweighted case.

**Lemma 12.5.10.** *Given an  $(h, \varepsilon, \alpha)$ -flow-sparsifier algorithm for unweighted graphs, it is possible to construct an  $(h \cdot \log U, \varepsilon, \alpha)$ -flow-sparsifier algorithm for weighted graphs  $G = (V, E, \mu)$  with capacity ratio  $U$  obeying*

$$U = \frac{\max_{e\in E} \mu_e}{\min_{e\in E} \mu_e} = \text{poly}(|V|).$$

*Proof.* We write each edge in binary so that  $G = \sum_{i=0}^{\log U} 2^i G_i$  for some unweighted graphs  $\{G_i = (V, E_i)_{i\in[\log U]}\}$ , where  $|E_i| \leq m$  for all  $i$ . We now apply the unweighted flow-sparsifier to each  $G_i$  in turn to obtain graphs  $\{G'_i\}$ . We let  $G' \stackrel{\text{def}}{=} \sum_{i=0}^{\log U} 2^i G'_i$  be the weighted flow-sparsified graph. By the assumption on the unweighted flow-sparsifier, each  $G'_i$  is  $h$ -sparse, so that  $G'$  must have at most  $h \cdot \log U$  edges. Similarly,  $G'$  is an  $\varepsilon$ -cut approximation of  $G$ , as each  $G'_i$  is an  $\varepsilon$ -cut approximation of the corresponding  $G_i$ . Letting  $\mathbf{M}_i$  be the embedding of  $G'_i$  into  $G_i$ , we can consider the embedding  $\mathbf{M} = \sum_{i=0}^{\log U} 2^i \mathbf{M}_i$  of  $G'$  into  $G$ . As each  $\mathbf{M}_i$  has congestion bounded by  $\alpha$ , it must be the case that  $\mathbf{M}$  also has congestion bounded by  $\alpha$ . The time to run the weighted flow sparsifier and to invoke  $\mathbf{M}$  is now  $\tilde{O}(m) \cdot \log U = \tilde{O}(m)$  by our assumption on  $U$ . □

The next step is to construct a routine which flow-sparsifies a constant fraction of the edges of  $E$ . This routine will then be applied iteratively to produce the final flow-sparsifier.

**Lemma 12.5.11.** *On input an unweighted graph  $G = (V, E)$ , there is an algorithm that runs in  $\tilde{O}(m)$  and computes a partition of  $E$  into  $(F, \bar{F})$ , an edge set  $F' \subseteq F$  with weight vector  $w_{F'} \in \mathbb{R}^E$ ,  $\text{support}(w_{F'}) = F'$ , and an embedding  $\mathbf{H} : \mathbb{R}^{F'} \rightarrow \mathbb{R}^E$  with the following properties:*

1.  $F$  contains most of the volume of  $G$ , i.e.

$$|F| \geq \frac{|E|}{2};$$

2.  $F'$  contains only  $\tilde{O}(n)$  edges, i.e.  $|F'| \leq \tilde{O}(n)$ .

3. The weights  $w_{F'}$  are bounded

$$\forall e \in F' \quad , \quad \frac{1}{\text{poly}(n)} \leq w_{F'}(e) \leq n.$$

4. The graph  $H' = (V, F', w_{F'})$  is an  $\varepsilon$ -cut approximation to  $H = (V, F)$ , i.e. for all  $S \subseteq V$ :

$$(1 - \varepsilon)|\partial_H(S)| \leq w_{F'}(\partial_{H'}(S)) \leq (1 + \varepsilon)|\partial_H(S)|.$$

5. The embedding  $\mathbf{H}$  from  $H = (V, F', w_{F'})$  to  $G$  has bounded congestion

$$\text{cong}(\mathbf{H}) = \tilde{O}(1).$$

and can be applied in time  $\tilde{O}(m)$ .

Given Lemma 12.5.10 and Lemma 12.5.11, it is straightforward to complete the proof of Theorem 12.5.3.

*Proof.* Using Lemma 12.5.10, we reduce the objective of Theorem 12.5.3 to running a  $(\tilde{O}(n), \varepsilon, \tilde{O}(1))$ -flow sparsifier on  $\log U$  unweighted graphs. To construct this unweighted flow sparsifier, we apply Lemma 12.5.11 iteratively as follows. Starting with the instance unweighted graph  $G_1 = (V, E_1)$ , we run the algorithm of Lemma 12.5.11 on the current graph  $G_t = (V, E_t)$  to produce the sets  $F_t$  and  $F'_t$ , the weight vector  $w_{F'_t}$  and the embedding  $\mathbf{H}_t : \mathbb{R}^{F'_t} \rightarrow \mathbb{R}^E$ . To proceed to the next iteration, we then define  $E_{t+1} \stackrel{\text{def}}{=} E_t \setminus F_t$  and move on to  $G_{t+1}$ .

By Lemma 12.5.11, at every iteration  $t$ ,  $|F_t| \geq \frac{1}{2} \cdot |E_t|$ , so that  $|E_{t+1}| \leq \frac{1}{2} \cdot |E_t|$ . This shows that there can be at most  $T \leq \log(|E_1|) = O(\log n)$  iterations.

After the last iteration  $T$ , we have effectively partitioned  $E_1$  into disjoint subsets  $\{F_t\}_{t \in [T]}$ , where each  $F_t$  is well-approximated but the weighted edgeset  $F'_t$ . We then output the weighted graph  $G' = (V, E' \stackrel{\text{def}}{=} \bigcap_{t=1}^T F'_t, w' \stackrel{\text{def}}{=} \sum_{t=1}^T w_{F'_t})$ , which is the sum of the disjoint weighted edges sets  $\{F'_t\}_{t \in [T]}$ . We also output the embedding  $\mathbf{M} : \mathbb{R}^{E'} \rightarrow \mathbb{R}^E$  from  $G'$  to  $G$ , defined as the direct sum

$$\mathbf{M} = \bigoplus_{t=1}^T \mathbf{H}_t.$$

In words,  $\mathbf{M}$  maps an edge  $e' \in E'$  by finding  $t$  for which  $e' \in F'_t$  and applying the corresponding  $\mathbf{H}_t$ .

We are now ready to prove that this algorithm with output  $G'$  and  $\mathbf{M}$  is an efficient  $(\tilde{O}(n), \varepsilon, \tilde{O}(n))$ -flow sparsifier. To bound the capacity ratio  $U'$  of  $G'$ , we notice that

$$U' \leq \max_t \frac{\max_{e \in F'_t} w_{F'_t}(e)}{\min_{e \in F'_t} w_{F'_t}(e)} \leq \text{poly}(n),$$

where we used the fact that the sets  $F'_t$  are disjoint and the guarantee on the range of  $w_{F'_t}$ .

Next, we bound the sparsity of  $G'$ . By Lemma 12.5.11,  $F'_t$  contains at most  $\tilde{O}(n)$  edges. As a result, we get the required bound

$$|E'| = \sum_{t=1}^T |F'_t| \leq \tilde{O}(Tn) = \tilde{O}(n).$$

For the cut approximation, we consider any  $S \subseteq V$ . By the cut guarantee of Lemma 12.5.11, we have that, for all  $t \in [T]$ ,

$$(1 - \varepsilon)|\partial_G(S) \cap F_t| \leq w_{F'_t}(\partial_G(S) \cap F'_t) \leq (1 + \varepsilon)|\partial_G(S) \cap F_t|.$$

Summing over all  $t$ , as  $E' = \dot{\bigcup} F'_t$  and  $E = \dot{\bigcup} F_t$ , we obtain the required approximation

$$(1 - \varepsilon)|\partial_G(S)| \leq w'(\partial_{G'}(S)) \leq (1 + \varepsilon)|\partial_G(S)|.$$

The congestion of  $\mathbf{M}$  can be bounded as follows

$$\text{cong}(\mathbf{M}) \leq \sum_{t=1}^T \text{cong}(\mathbf{H}_t) = \tilde{O}(T) = \tilde{O}(1).$$

To conclude the proof, we address the efficiency of the flow sparsifier. The algorithm applies the routine of Lemma 12.5.11 for  $T = \tilde{O}(1)$  times and hence runs in time  $\tilde{O}(m)$ , as required. Invoking the embedding  $\mathbf{M}$  requires invoking each of the  $T$  embeddings  $\mathbf{H}_t$ . This takes time  $\tilde{O}(Tm) = \tilde{O}(m)$ .  $\square$

### 12.5.3.1 Flow Sparsification of Unweighted Graphs: Proof of Lemma 12.5.11

In this subsection, we prove Lemma 12.5.11. Our starting point is the following decomposition statement, which shows that we can form a partition of an unweighted graph where most edges do not cross the boundaries and the subgraphs induced within each set of this partition are near-expanders. The following lemma is implicit in Spielman and Teng's local clustering approach to spectral sparsification [244].

**Lemma 12.5.12** (Decomposition Lemma). *For an unweighted graph  $G = (V, E)$ , in  $\tilde{O}(m)$ -time we can produce a partition  $V_1, \dots, V_k$  of  $V$  and a collection of sets  $S_1, \dots, S_k \subseteq V$  with the following properties:*

- For all  $i$ ,  $S_i$  is contained in  $V_i$ .
- For all  $i$ , there exists a set  $T_i$  with  $S_i \subseteq T_i \subseteq V_i$ , such that

$$\Phi(G(T_i)) \geq \Omega\left(\frac{1}{\log^2 n}\right).$$

- At least half of the edges are found within the sets  $\{S_i\}$ , i.e.

$$\sum_{i=1}^k |E(S_i)| = \sum_{i=1}^k |\{e = \{a, b\} : a \in S_i, b \in S_i\}| \geq \frac{1}{2}|E|.$$

To design an algorithm satisfying the requirements of Lemma 12.5.11, we start by applying the Decomposition Lemma to our unweighted input graph  $G = (V, E)$  to obtain the partition  $\{V_i\}_{i \in [k]}$  and the sets  $\{S_i\}_{i \in [k]}$ . We let  $G_i \stackrel{\text{def}}{=} (S_i, E(S_i))$ . To reduce the number of edges, while preserving cuts, we apply a spectral sparsification algorithm to each  $G_i$ . Concretely, by applying the spectral sparsification by effective resistances of Spielman and Srivastava [242] to each  $G_i$ , we obtain weighted graphs  $G'_i = (S_i, E'_i \subseteq E(S_i), w'_i)$  in time  $\sum_{i=1}^k \tilde{O}(|E(S_i)|) \leq \tilde{O}(|E|)$  with  $|E'_i| \leq \tilde{O}(|S_i|)$  and the property that cuts are preserved<sup>6</sup> for all  $i$ :

$$\forall S \subseteq S_i \quad (1 - \varepsilon) \cdot |\partial_{G_i}(S)| \leq w'_i(\partial_{G'_i}(S)) \leq (1 + \varepsilon) \cdot |\partial_{G_i}(S)|.$$

Moreover, the spectral sparsification of [242] constructs the weights  $\{w'_i(e)\}_{e \in E'_i}$  such that

$$\forall e \in E'_i \quad \frac{1}{\text{poly}(n)} \leq \frac{1}{\text{poly}(|S_i|)} \leq w'_i(e) \leq |S_i| \leq n.$$

<sup>6</sup>The spectral sparsification result actually yields the stronger spectral approximation guarantee, but for our purposes the cut guarantee suffices.

To complete the description of the algorithm, we output the partition  $(F, \bar{F})$  of  $E$ , where

$$F \stackrel{\text{def}}{=} \bigcup_{i=1}^k E(S_i).$$

We also output the set of weighted sparsified edges  $F'$ .

$$F' \stackrel{\text{def}}{=} \bigcup_{i=1}^k E'_i.$$

The weight  $w_{F'}(e)$  of edge  $e \in F'$  is given by finding  $i$  such that  $e \in E'_i$  and setting  $w_{F'}(e) = w'_i(e)$ .

We now depart from Spielman and Teng's construction by endowing our  $F'$  with an embedding onto  $G$ . The embedding  $\mathbf{H} : \mathbb{R}^{F'} \rightarrow \mathbb{R}^E$  of the graph  $H = (V, F', w_{F'})$  to  $G$  is constructed by using the oblivious electrical-flow routing of  $E(S_i)$  into  $G(V_i)$ . More specifically, as the sets  $\{V_i\}$  partition  $V$ , the embedding  $\mathbf{H}$  can be expressed as the following direct sum over the orthogonal subspaces  $\{\mathbb{R}^{E(V_i) \times E'_i}\}$ .

$$\mathbf{H} \stackrel{\text{def}}{=} \left( \bigoplus_{i=1}^k \mathbf{B}_{E(V_i)} \mathcal{L}_{G(V_i)}^+ \mathbf{B}_{E(V_i)}^T \mathbf{I}_{(E(V_i), E'_i)} \right),$$

where  $\mathbf{I}_{(E(V_i), E'_i)}$  is the identity mapping of the edges  $E'_i \subseteq E(V_i)$  of  $F'$  over  $V_i$  to the edges  $E(V_i)$  of  $V_i$  in  $G$ . Notice that there is no dependence on the resistances over  $G$  as  $G$  is unweighted.

This complete the description of the algorithm. We are now ready to give the proof of Lemma 12.5.11.

[Proof of Lemma 12.5.11]. The algorithm described above performs a decomposition of the input graph  $G = (V, E)$  in time  $\tilde{O}(m)$  by the Decomposition Lemma. By the result of Spielman and Srivastava [242], each  $G_i$  is sparsified in time  $\tilde{O}(|E(S_i)|)$ . Hence, the sparsification step requires time  $\tilde{O}(m)$  as well. This shows that the algorithm runs in  $\tilde{O}(m)$ -time, as required.

By the Decomposition Lemma, we know that  $|F| = \sum_{i=1}^k |E(S_i)| \geq \frac{|E|}{2}$ , which satisfies the requirement of the Lemma. Moreover, by the spectral sparsification result, we know that  $|F'| = \sum_{i=1}^k |E'_i| \leq \sum_{i=1}^k \tilde{O}(|S_i|) \leq \tilde{O}(n)$ , as required. We also saw that by construction the weights  $w_{F'}$  are bounded:

$$\forall e \in F' \quad , \quad \frac{1}{\text{poly}(n)} \leq w_{F'}(e) \leq n.$$

To obtain the cut-approximation guarantee, we use the fact that for every  $i$ , by spectral sparsification,

$$\forall S \subseteq S_i \quad , \quad (1 - \varepsilon) \cdot |\partial_{G_i}(S)| \leq w'_i(\partial_{G'_i}(S)) \leq (1 + \varepsilon) \cdot |\partial_{G_i}(S)|.$$

We have  $H' = (V, F', w_{F'})$  and  $H = (V, F)$ . Consider now  $T \subseteq V$  and apply the previous bound to  $T \cap S_i$  for all  $i$ . Because  $F' \subseteq F = \cup_{i=1}^k E(S_i)$ , we have that summing over the  $k$  bounds yields

$$\forall T \subseteq V \quad , \quad (1 - \varepsilon) |\partial_H(T)| \leq w_{F'}(\partial_{H'}(T)) \leq (1 + \varepsilon) |\partial_H(T)|,$$

which is the desired cut-approximation guarantee.

Finally, we are left to prove that the embedding  $\mathbf{H}$  from  $H' = (V, F', w_{F'})$  to  $G = (V, E)$  has low congestion and can be applied efficiently. By definition of congestion,

$$\text{cong}(\mathbf{H}) = \max_{\mathbf{x} \in \mathbb{R}^{F'}} \frac{\|\mathbf{H}\mathbf{x}\|_\infty}{\|\mathbf{U}_{F'}^{-1}\mathbf{x}\|_\infty} = \|\|\mathbf{H}\mathbf{U}_{F'}\mathbf{1}_{F'}\|_\infty = \left\| \left\| \bigoplus_{i=1}^k \mathbf{B}_{E(V_i)} \mathcal{L}_{G(V_i)}^+ \mathbf{B}_{E(V_i)}^T \mathbf{I}_{(E(V_i), E'_i)} \right\| \mathbf{U}_{F'} \mathbf{1}_{F'} \right\|_\infty.$$

Decomposing  $\mathbb{R}^E$  into the subspaces  $\{\mathbb{R}^{E(V_i)}\}$  and  $\mathbb{R}^{F'}$  into the subspaces  $\{\mathbb{R}^{E'_i}\}$  we have:

$$\text{cong}(\mathbf{H}) \leq \max_{i \in [k]} \left\| \left\| \mathbf{B}_{E(V_i)} \mathcal{L}_{G(V_i)}^+ \mathbf{B}_{E(V_i)}^T \mathbf{I}_{(E(V_i), E'_i)} \right\|_{\mathbf{U}_{E'_i} \mathbf{1}_{E'_i}} \right\|_{\infty}.$$

For each  $i \in [k]$ , consider now the set of demands  $D_i$  over  $V_i$ ,  $D_i \stackrel{\text{def}}{=} \{\chi_e\}_{e \in E'_i}$ , given by the edges of  $E'_i$  with their capacities  $w'_i$ . That is,  $\chi_e \in \mathbb{R}^{V_i}$  is the demand corresponding to edge  $e \in E'_i$  with weight  $w'_i(e)$ . Consider also the electrical routing  $\mathbf{A}_{\mathcal{E}, i} = \mathbf{B}_{E(V_i)} \mathcal{L}_{G(V_i)}^+$  over  $G(V_i)$ . Then:

$$\text{cong}(\mathbf{H}) \leq \max_{i \in [k]} \text{cong}(\mathbf{A}_{\mathcal{E}, i} D_i)$$

Notice that, by construction,  $D_i$  is routable in  $G'_i = (S_i, E'_i, w'_i)$  and  $\text{OPT}_{G'_i}(D_i) = 1$ . But, by our use of spectral sparsifiers in the construction,  $G'_i$  is an  $\varepsilon$ -cut approximation of  $G_i$ . Hence, by the flow-cut gap of Aumann and Rabani [21], we have:

$$\text{OPT}_{G_i}(D_i) \leq O(\log(|D_i|)) \cdot \text{OPT}_{G'_i}(D_i) \leq \tilde{O}(1).$$

When we route  $D_i$  oblivious in  $G(V_i)$ , we can consider the  $E(S_i)$ -competitive ratio  $\rho^{E(S_i)}(\mathbf{A}_{\mathcal{E}, i})$  of the electrical routing  $\mathbf{A}_{\mathcal{E}, i} = \mathbf{B}_{E(V_i)} \mathcal{L}_{G(V_i)}^+$ , as  $D_i$  is routable in  $E(S_i)$ , because  $E'_i \subseteq E(S_i)$ . We have

$$\text{cong}(\mathbf{H}) \leq \max_{i \in [k]} \rho_{G(V_i)}^{E(S_i)}(\mathbf{A}_{\mathcal{E}, i}) \cdot \text{OPT}_{G(V_i)}^{E(S_i)}(D_i) = \max_{i \in [k]} \rho_{G(V_i)}^{E(S_i)}(\mathbf{A}_{\mathcal{E}, i}) \cdot \text{OPT}_{G_i}(D_i),$$

Finally, putting these bounds together, we have:

$$\text{cong}(\mathbf{H}) \leq \max_{i \in [k]} \rho_{G(V_i)}^{E(S_i)}(\mathbf{A}_{\mathcal{E}, i}) \cdot \text{OPT}_{G_i}(D_i) \leq \tilde{O}(1) \cdot \max_{i \in [k]} \rho_{G(V_i)}^{E(S_i)}(\mathbf{A}_{\mathcal{E}, i}).$$

But, by the Decomposition Lemma, there exists  $T_i$  with  $S_i \subseteq T_i \subseteq V_i$  such that

$$\Phi(G(T_i)) \geq \Omega\left(\frac{1}{\log^2 n}\right).$$

Then, by Lemma 12.5.9, we have that:

$$\rho_{G(V_i)}^{E(S_i)}(\mathbf{A}_{\mathcal{E}, i}) \leq O\left(\frac{\log \text{vol}(G(T_i))}{\Phi(G(T_i))^2}\right) \leq \tilde{O}(1).$$

This concludes the proof that  $\text{cong}(\mathbf{H}) \leq \tilde{O}(1)$ . To complete the proof of the Lemma, we just notice that  $\mathbf{H}$  can be invoked in time  $\tilde{O}(m)$ . A call of  $\mathbf{H}$  involves solving  $k$ -electrical-problems, one for each  $G(V_i)$ . This can be done in time  $\sum_{i=1}^k \tilde{O}(|E(V_i)|) \leq \tilde{O}(m)$ , using any of the nearly-linear Laplacian system solvers available, such as [141].  $\square$

## ■ 12.6 Removing Vertices in Oblivious Routing Construction

In this section we show how to reduce computing an efficient oblivious routing on a graph  $G = (V, E)$  to computing an oblivious routing for  $t$  graphs with  $\tilde{O}(\frac{|V|}{t})$  vertices and at most  $|E|$  edges. Formally we show

**Theorem 12.6.1.** *Let  $G = (V, E, \mu)$  be an undirected capacitated graph with capacity ratio  $U$ . For all  $t > 0$  in  $\tilde{O}(t \cdot |E|)$  time we can compute graphs  $G_1, \dots, G_t$  each with at most  $\tilde{O}(\frac{|E| \log(U)}{t})$  vertices, at most  $|E|$  edges, and capacity ratio at most  $|V| \cdot U$ , such that given oblivious routings  $\mathbf{A}_i$  for each*



$G_i$ , in  $\tilde{O}(t \cdot |E|)$  time we can compute an oblivious routing  $\mathbf{A} \in \mathbb{R}^{E \times V}$  on  $G$  such that

$$\mathcal{T}(\mathbf{A}) = \tilde{O}\left(t \cdot |E| + \sum_{i=1}^t \mathcal{T}(\mathbf{A}_i)\right) \quad \text{and} \quad \rho(\mathbf{A}) = \tilde{O}\left(\max_i \rho(\mathbf{A}_i)\right)$$

We break this proof into several parts. First we show how to embed  $G$  into a collection of  $t$  graphs consisting of trees minus some edges which we call *patrial tree embeddings* (Section 12.6.1). Then we show how to embed a partial tree embedding in an “almost  $j$ -tree” [182], that is a graph consisting of a tree and a subgraph on at most  $j$  vertices, for  $j = 2t$  (Section 12.6.2). Finally, we show how to reduce oblivious routing on an almost  $j$ -tree to oblivious routing on a graph with at most  $O(j)$  vertices by removing degree-1 and degree-2 vertices (Section 12.6.3). Finally, in Section 12.6.4 we put this all together to prove Theorem 12.4.10.

We remark that much of the ideas in the section were either highly influenced from [182] or are direct restatements of theorems from [182] adapted to our setting. We encourage the reader to look over that paper for further details regarding the techniques used in this section.

### ■ 12.6.1 From Graphs to Partial Tree Embeddings

To prove Theorem 12.4.10, we make heavy use of spanning trees and various properties of them. In particular, we use the facts that for every pair of vertices there is a unique *tree path* connecting them, that every edge in the tree *induces a cut* in the graph, and that we can embed a graph in a tree by simply routing every edge over its tree path and that the congestion of this embedding will be determined by the *load* the edges place on tree edges. We define these quantities formally below.

**Definition 12.6.2** (Tree Path). For undirected graph  $G = (V, E)$ , spanning tree  $T$ , and all  $a, b \in V$  we let  $P_{a,b} \subseteq E$  denote the unique path from  $a$  to  $b$  using only edges in  $T$  and we let  $\mathbf{p}_{a,b} \in \mathbb{R}^E$  denote the vector representation of this path corresponding to the unique vector sending one unit from  $a$  to  $b$  that is nonzero only on  $T$  (i.e.  $\mathbf{B}^T \mathbf{p}_{a,b} = \chi_{a,b}$  and  $\forall e \in E \setminus T$  we have  $\mathbf{p}_{a,b}(e) = 0$ ).

**Definition 12.6.3** (Tree Cuts). For undirected  $G = (V, E)$  and spanning tree  $T \subseteq E$  the *edges cut by*  $e$ ,  $\partial_T(F)$ , and the edges cut by  $F$ ,  $\partial_T(e)$ , are given by

$$\partial_T(e) \stackrel{\text{def}}{=} \{e' \in E \mid e' \in P_e\} \quad \text{and} \quad \partial_T(F) \stackrel{\text{def}}{=} \cup_{e \in F} \partial(e)$$

**Definition 12.6.4** (Tree Load). For undirected capacitated  $G = (V, E, \boldsymbol{\mu})$  and spanning tree  $T \subseteq E$  the load on edge  $e \in E$  by  $T$ ,  $\text{cong}_T(e)$  is given by  $\text{load}_T(e) = \sum_{e' \in E \mid e \in P_{e'}} \boldsymbol{\mu}_{e'}$

While these properties do highlight the fact that we could just embed our graph into a collection of trees to simplify the structure of our graph, this approach suffers from a high computational cost [225]. Instead we show that we can embed parts of the graph onto collections of trees at a lower computational cost but higher complexity. In particular we will consider what we call partial tree embeddings.

**Definition 12.6.5** ([Partial Tree Embedding<sup>7</sup>]). For undirected capacitated graph  $G = (V, E, \boldsymbol{\mu})$ , spanning tree  $T$  and spanning tree subset  $F \subseteq T$  we define the *partial tree embedding graph*  $H = H(G, T, F) = (V, E', \boldsymbol{\mu}')$  to be a graph on the same vertex set where  $E' = T \cup \partial_T(F)$  and

$$\forall e \in E' \quad : \quad \boldsymbol{\mu}'(e) = \begin{cases} \text{load}_T(e) & \text{if } e \in T \setminus F. \\ \boldsymbol{\mu}(e) & \text{otherwise} \end{cases}$$

<sup>7</sup>This is a restatement of the  $H(T, F)$  graphs in [182].

Furthermore, we let  $\mathbf{M}_H \in \mathbb{R}^{E' \times E}$  denote the embedding from  $G$  to  $H(G, T, F)$  where edges not cut by  $F$  are routed over the tree and other edges are mapped to themselves.

$$\forall e \in E \quad : \quad \mathbf{M}_H(e) = \begin{cases} \mathbf{p}_e & e \notin \partial_T(F) \\ \mathbf{1}_e & \text{otherwise} \end{cases}$$

and we let  $\mathbf{M}'_H \in \mathbb{R}^{E \times E'}$  denote the embedding from  $H$  to  $G$  that simply maps edges in  $H$  to their corresponding edges in  $G$ , i.e.  $\forall e \in E'$ ,  $\mathbf{M}'_H(e) = \mathbf{1}_e$ .

Note that by definition  $\text{cong}(\mathbf{M}_H) \leq 1$ , i.e. a graph embeds into its partial tree embedding with no congestion. However, to get embedding guarantees in the other direction more work is required. For this purpose we use a lemma from Madry [182] saying that we can construct a convex combination or a distribution of partial tree embeddings we can get such a guarantee.

**Lemma 12.6.6** ([Probabilistic Partial Tree Embedding<sup>8</sup>]). *For any undirected capacitated graph  $G = (V, E, \boldsymbol{\mu})$  and any  $t > 0$  in  $\tilde{O}(t \cdot m)$  time we can find a collection of partial tree embeddings  $H_1 = H(G, T_1, F_1), \dots, H_t = H(G, T_t, F_t)$  and coefficients  $\lambda_i \geq 0$  with  $\sum_i \lambda_i = 1$  such that  $\forall i \in [t]$  we have  $|F_i| = \tilde{O}(\frac{m \log U}{t})$  and such that  $\sum_i \lambda_i \mathbf{M}'_{H_i}$  embeds  $G' = \sum_i \lambda_i G_i$  into  $G$  with congestion  $\tilde{O}(1)$ .*

Using this lemma, we can prove that we can reduce constructing an oblivious routing for a graph to constructing oblivious routings on several partial tree embeddings.

**Lemma 12.6.7.** *Let the  $H_i$  be graphs produced by Lemma 12.6.6 and for all  $i$  let  $\mathbf{A}_i$  be an oblivious routing algorithm for  $H_i$ . It follows that  $\mathbf{A} = \sum_i \lambda_i \mathbf{M}'_{H_i} \mathbf{A}_i$  is an oblivious routing on  $G$  with  $\rho(\mathbf{A}) \leq \tilde{O}(\max_i \rho(\mathbf{A}_i) \log n)$  and  $\mathcal{T}(\mathbf{A}) = O(\sum_i \mathcal{T}(\mathbf{A}_i))$ .*

*Proof.* The proof is similar to the proof of Lemma 12.4.8. For all  $i$  let  $\mathbf{U}_i$  denote the capacity matrix of graph  $G_i$ . Then using Lemma 12.4.3 we get

$$\rho(\mathbf{A}) = \|\mathbf{U}^{-1} \mathbf{A} \mathbf{B}^T \mathbf{U}\|_\infty = \left\| \sum_{i=1}^t \lambda_i \mathbf{U}^{-1} \mathbf{M}'_{H_i} \mathbf{A}_i \mathbf{B}^T \mathbf{U} \right\|_\infty$$

Using that  $\mathbf{M}_{H_i}$  is an embedding and therefore  $\mathbf{B}_{H_i}^T \mathbf{M}_{H_i} = \mathbf{B}^T$  we get

$$\rho(\mathbf{A}) = \left\| \sum_{i=1}^t \lambda_i \mathbf{U}^{-1} \mathbf{M}'_{H_i} \mathbf{A}_i \mathbf{B}_{H_i}^T \mathbf{M}_{H_i} \mathbf{U} \right\|_\infty \leq \max_{j,k} \left\| \sum_{i=1}^t \lambda_i \mathbf{U}^{-1} \mathbf{M}'_i \mathbf{U}_i \right\|_\infty \cdot \rho(\mathbf{A}_j) \cdot \text{cong}(\mathbf{M}_{H_k})$$

Since  $\sum_i \lambda_i \mathbf{M}'_{H_i}$  is an embedding of congestion of at most  $\tilde{O}(1)$  and  $\text{cong}(\mathbf{M}_{H_k}) \leq 1$  we have the desired result.  $\square$

## ■ 12.6.2 From Partial Tree Embeddings To Almost- $j$ -trees

Here we show how to reduce constructing an oblivious routing for a partial tree embedding to constructing an oblivious routing for what Madry [182] calls an “almost  $j$ -tree,” the union of a tree plus a subgraph on at most  $j$  vertices. First we define such objects and then we prove the reduction.

**Definition 12.6.8** (Almost -tree). We call a graph  $G = (V, E)$  an *almost  $j$ -tree* if there is a spanning tree  $T \subseteq E$  such that the endpoints of  $E \setminus T$  include at most  $j$  vertices.

<sup>8</sup>This is an adaptation of Corollary 5.6 in [182]

**Lemma 12.6.9.** *For undirected capacitated  $G = (V, E, \boldsymbol{\mu})$  and partial tree embedding  $H = H(G, T, F)$  in  $\tilde{O}(|E|)$  time we can construct an almost  $2 \cdot |F|$ -tree  $G' = (V, E', \boldsymbol{\mu}')$  with  $|E'| \leq |E|$  and an embedding  $\mathbf{M}'$  from  $G'$  to  $H$  such that  $H$  is embeddable into  $G'$  with congestion 2,  $\text{cong}(\mathbf{M}') = 2$ , and  $\mathcal{T}(\mathbf{M}') = \tilde{O}(|E|)$ .*

*Proof.* For every  $e = (a, b) \in E$ , we let  $v^1(e) \in V$  denote the first vertex on tree path  $P_{(a,b)}$  incident to  $F$  and we let  $v^2(e) \in V$  denote the last vertex incident to  $F$  on tree path  $P_{(a,b)}$ . Note that for every  $e = (a, b) \in T$  we have that  $(v^1(e), v^2(e)) = e$ .

We define  $G' = (V, E', \boldsymbol{\mu}')$  to simply be the graph that consists of all these  $(v^1(e), v^2(e))$  pairs

$$E' = \{(a, b) \mid \exists e \in E \text{ such that } (a, b) = (v^1(e), v^2(e))\}$$

and we define the weights to simply be the sums

$$\forall e' \in E' \quad : \quad \boldsymbol{\mu}'(e') \stackrel{\text{def}}{=} \sum_{e \in E \mid e=(v^1(e'), v^2(e'))} \boldsymbol{\mu}(e)$$

Now to embed  $H$  in  $G'$  we define  $\mathbf{M}$  by

$$\forall e = (a, b) \in E \quad : \quad \mathbf{M}\mathbf{1}_e = \mathbf{p}_{a, v^1(e)} + \mathbf{1}_{(v^1(e), v^2(e))} + \mathbf{p}_{v^2(e), b}$$

and to embed  $G'$  in  $H$  we define  $\mathbf{M}'$  by

$$\forall e' \in E' \quad : \quad \mathbf{M}'\mathbf{1}_{e'} = \sum_{e=(a,b) \in E \mid e'=(v^1(e), v^2(e))} \frac{\boldsymbol{\mu}(e)}{\boldsymbol{\mu}'(e')} \left[ \mathbf{p}_{v^1(e), a} + \mathbf{1}_{(a, b)} + \mathbf{p}_{b, v^2(e)} \right]$$

In other words we route edges in  $H$  along the tree until we encounter nodes in  $F$  and then we route them along added edges and we simply route the other way for the reverse embedding. By construction clearly the congestion of the embedding in either direction is 2.

To bound the running time, we note that by having every edge  $e$  in  $H$  maintain its  $v^1(e)$  and  $v^2(e)$  information, having every edge  $e'$  in  $E'$  maintain the set  $\{e \in E \mid e' = (v^1(e), v^2(e))\}$  in a list, and using link cut trees [240] or the static tree structure in [141] to update information along tree paths we can obtain the desired value of  $\mathcal{T}(\mathbf{M}')$ .  $\square$

### ■ 12.6.3 From Almost- $J$ Trees to Less Vertices

Here we show that by “greedy elimination” [243] [147] [146], i.e. removing all degree 1 and degree 2 vertices in  $O(m)$  time we can reduce oblivious routing in almost- $j$ -trees to oblivious routing in graphs with  $O(j)$  vertices while only losing  $O(1)$  in the competitive ratio. Again, we remark that the lemmas in this section are derived heavily from [182] but repeated for completeness and to prove additional properties that we will need for our purposes.

We start by showing that an almost- $j$ -tree with no degree 1 or degree 2 vertices has at most  $O(j)$  vertices.

**Lemma 12.6.10.** *For any almost  $j$ -tree  $G = (V, E)$  with no degree 1 or degree 2 vertices, we have  $|V| \leq 3j - 2$ .*

*Proof.* Since  $G$  is an almost  $j$ -tree, there is some  $J \subseteq V$  with  $|J| \leq j$  such that the removal of all edges with both endpoints in  $J$  creates a forest. Now, since  $K = V - J$  is incident only to forest edges clearly the sum of the degrees of the vertices in  $K$  is at most  $2(|V| - 1)$  (otherwise there would be a cycle). However, since the minimum degree in  $G$  is 3, clearly this sum is at least  $3(|V| - j)$ . Combining yields that  $3|V| - 3j \leq 2|V| - 2$ .  $\square$

Next, we show how to remove degree one vertices efficiently.

**Lemma 12.6.11** (Removing Degree One Vertices). *Let  $G = (V, E, \boldsymbol{\mu})$  be an unweighted capacitated graph, let  $a \in V$  be a degree 1 vertex, let  $e = (a, b) \in E$  be the single edge incident to  $a$ , and let  $G' = (V', E', \boldsymbol{\mu}')$  be the graph that results from simply removing  $e$  and  $a$ , i.e.  $V' = V \setminus \{a\}$  and  $E' = E \setminus \{e\}$ . Given  $a \in V$  and an oblivious routing algorithm  $\mathbf{A}'$  in  $G'$  in  $O(1)$  time we can construct an oblivious routing algorithm  $\mathbf{A}$  in  $G$  such that*

$$\mathcal{T}(\mathbf{A}) = O(\mathcal{T}(\mathbf{A}') + 1) \quad , \quad \text{and} \quad \rho(\mathbf{A}) = \rho(\mathbf{A}')$$

*Proof.* For any demand vector  $\boldsymbol{\chi}$ , the only way to route demand at  $a$  in  $G$  is over  $e$ . Therefore, if  $\mathbf{B}\mathbf{f} = \boldsymbol{\chi}$  then  $\mathbf{f}(e) = \boldsymbol{\chi}$ . Therefore, to get an oblivious routing algorithm on  $G$ , we can simply send demand at  $a$  over edge  $e$ , modify the demand at  $b$  accordingly, and then run the oblivious routing algorithm on  $G'$  on the remaining vertices. The routing algorithm we get is the following

$$\mathbf{A} \stackrel{\text{def}}{=} \mathbf{I}_{E' \rightarrow E} \mathbf{A}' (\mathbf{I} + \mathbf{1}_b \mathbf{1}_a^T) + \mathbf{1}_e \mathbf{1}_a^T$$

Since all routing algorithms send this flow on  $e$  we get that  $\rho(\mathbf{A}) = \rho(\mathbf{A}')$  and since the above operators not counting  $\mathbf{A}$  have only  $O(1)$  entries that are not the identity we can clearly implement the operations in the desired running time.  $\square$

Using the above lemma we show how to remove all degree 1 and 2 vertices in  $O(m)$  time while only increasing the congestion by  $O(1)$ .

**Lemma 12.6.12** (Greedy Elimination). *Let  $G = (V, E, \boldsymbol{\mu})$  be an unweighted capacitated graph and let  $G' = (V', E', \boldsymbol{\mu}')$  be the graph the results from iteratively removing vertices of degree 1 and replacing degree 2 vertices with an edge connecting its neighbors of the minimum capacity of its adjacent edges. We can construct  $G'$  in  $O(m)$  time and given an oblivious routing algorithm  $\mathbf{A}'$  in  $G'$  in  $O(1)$  time we can construct an oblivious routing algorithm  $\mathbf{A}$  in  $G$  such that <sup>9</sup>*

$$\mathcal{T}(\mathbf{A}) = O(\mathcal{T}(\mathbf{A}') + |E|) \quad , \quad \text{and} \quad \rho(\mathbf{A}) \leq 4 \cdot \rho(\mathbf{A}')$$

*Proof.* First we repeatedly apply Lemma 12.6.11 repeatedly to in reduce to the case that there are no degree 1 vertices. By simply array of the degrees of every vertex and a list of degree 1 vertices this can be done in  $O(m)$  time. We denote the result of these operations by graph  $K$ .

Next, we repeatedly find degree two vertices that have not been explored and explore this vertices neighbors to get a path of vertices,  $a_1, a_2, \dots, a_k \in V$  for  $k > 3$  such that each vertex  $a_2, \dots, a_{k-1}$  is of degree two. We then compute  $j = \arg \min_{i \in [k-1]} \boldsymbol{\mu}(a_i, a_{i+1})$ , remove edge  $(a_j, a_{j+1})$  and add an edge  $(a_1, a_k)$  of capacity  $\boldsymbol{\mu}(a_j, a_{j+1})$ . We denote the result of doing this for all degree two vertices by  $K'$  and note that again by careful implementation this can be performed in  $O(m)$  time.

Note that clearly  $K$  is embeddable in  $K'$  with congestion 2 just by routing every edge over itself except the removed edges which we route by the path plus the added edges. Furthermore,  $K'$  is embeddable in  $K$  with congestion 2 again by routing every edge on itself except for the edges which we added which we route back over the paths they came from. Furthermore, we note that clearly this embedding and the transpose of this operator is computable in  $O(m)$  time.

Finally, by again repeatedly applying Lemma 12.6.11 to  $K'$  until there are no degree 1 vertices we get a graph  $G'$  that has no degree one or degree two vertices (since nothing decreased the degree of vertices with degree more than two). Furthermore, by Lemma 12.6.11 and by Lemma 12.4.8 we see that we can compose these operators to compute  $\mathbf{A}$  with the desired properties.  $\square$

<sup>9</sup>Note that the constant of 4 below is improved to 3 in [182].

## ■ 12.6.4 Putting It All Together

Here we put together the previous components to prove the main theorem of this section.

*Node Reduction Theorem 12.4.10.* Using Lemma 12.6.6, we can construct  $G' = \sum_{i=1}^t \lambda_i G_i$  and embeddings  $M_1, \dots, M_t$  from  $G_i$  to  $G$ . Next we can apply Lemma 12.6.9 to each  $G_i$  to get almost- $j$ -trees  $G'_1, \dots, G'_t$  and embeddings  $M'_1, \dots, M'_t$  from  $G'_i$  to  $G_i$ . Furthermore, using Lemma 12.6.12 we can construct graphs  $G''_1, \dots, G''_t$  with the desired properties (the congestion ratio property follows from the fact that we only add capacities during these reductions)

Now given oblivious routing algorithms  $A''_1, \dots, A''_t$  on the  $G''_i$  and again by Lemma 12.6.12 we could get oblivious routing algorithms  $A'_1, \dots, A'_t$  on the  $G'_i$  with constant times more congestion. Finally, by the guarantees of Lemma 12.4.8 we have that  $A \stackrel{\text{def}}{=} \sum_{i=1}^t \lambda M_i M'_i A'_i$  is an oblivious routing algorithm that satisfies the requirements.  $\square$

## ■ 12.7 Nonlinear Projection and Maximum Concurrent Flow

### ■ 12.7.1 Gradient Descent Method for Nonlinear Projection Problem

In this section, we strengthen and generalize the **MaxFlow** algorithm to a more general setting. We believe this algorithm may be of independent interest as it includes maximum concurrent flow problem, the compressive sensing problem, etc. For some norms, e.g.  $\|\cdot\|_1$  as typically of interest compressive sensing, the Nesterov algorithm [208] can be used to replace gradient descent. However, this kind of accelerated method is not known in the general norm settings as good proxy function may not exist at all. Even worse, in the non-smooth regime, the minimization problem on the  $\|\cdot\|_p$  space with  $p > 2$  is difficult under some oracle assumption [201]. For these reasons we focus here on the gradient descent method which is always applicable.

Given a norm  $\|\cdot\|$ , we wish to solve the what we call the *non-linear projection* problem

$$\min_{\mathbf{x} \in L} \|\mathbf{x} - \mathbf{y}\|$$

where  $\mathbf{y}$  is an given point and  $L$  is a linear subspace. We assume the following:

#### Assumption 12.7.1.

1. There are a family of convex differentiable functions  $f_t$  such that for all  $\mathbf{x} \in L$ , we have

$$\|\mathbf{x}\| \leq f_t(\mathbf{x}) \leq \|\mathbf{x}\| + Kt$$

and the Lipschitz constant of  $\nabla f_t$  is  $\frac{1}{t}$ .

2. There is a projection matrix  $\mathbf{P}$  onto the subspace  $L$ .

In other words we assume that there is a family of regularized objective functions  $f_t$  and a projection matrix  $\mathbf{P}$ , which we can think of as an approximation algorithm of this projection problem.

Now, let  $\mathbf{x}^*$  be a minimizer of  $\min_{\mathbf{x} \in L} \|\mathbf{x} - \mathbf{y}\|$ . Since  $\mathbf{x}^* \in L$ , we have  $\mathbf{P}\mathbf{x}^* = \mathbf{x}^*$  and hence

$$\begin{aligned} \|\mathbf{P}\mathbf{y} - \mathbf{y}\| &\leq \|\mathbf{y} - \mathbf{x}^*\| + \|\mathbf{x}^* - \mathbf{P}\mathbf{y}\| \\ &\leq \|\mathbf{y} - \mathbf{x}^*\| + \|\mathbf{P}\mathbf{x}^* - \mathbf{P}\mathbf{y}\| \\ &\leq (1 + \|\mathbf{P}\|) \min_{\mathbf{x} \in L} \|\mathbf{x} - \mathbf{y}\|. \end{aligned} \tag{12.4}$$

Therefore, the approximation ratio of  $\mathbf{P}$  is  $1 + \|\mathbf{P}\|$  and we see that our problem is to show that we can solve nonlinear projection using a decent linear projection matrix. Our algorithm for solving this problem is shown in **NonlinearProjection**.

---

**Algorithm 34: NonlinearProjection**


---

**Input:** a point  $\mathbf{y}$  and  $\text{OPT} = \min_{\mathbf{x} \in L} \|\mathbf{x} - \mathbf{y}\|$ .  
 Let  $\mathbf{y}_0 = (\mathbf{I} - \mathbf{P})\vec{y}$  and  $\mathbf{x}_0 = 0$ .  
**for**  $j = 0, \dots$ , *until*  $2^{-j}\|\mathbf{P}\| \leq \frac{1}{2}$  **do**  
   If  $2^{-j}\|\mathbf{P}\| > 1$ , then let  $t_j = \frac{2^{-(j+2)}\|\mathbf{P}\|\text{OPT}}{K}$  and  $k_j = 3200\|\mathbf{P}\|^2 K$ .  
   If  $2^{-j}\|\mathbf{P}\| \leq 1$ , then let  $t_j = \frac{\varepsilon\text{OPT}}{2K}$  and  $k_j = \frac{800\|\mathbf{P}\|^2 K}{\varepsilon^2}$ .  
   Let  $g_j(\mathbf{x}) = f_{t_j}(\mathbf{P}\mathbf{x} - \mathbf{y}_j)$  and  $\mathbf{x}_0 = 0$ .  
   **for**  $i = 0, \dots, k_j - 1$  **do**  
      $\mathbf{x}_{i+1} = \mathbf{x}_i - \frac{t}{\|\mathbf{P}\|^2}(\nabla g_j(\mathbf{x}_i))^\#$ .  
   **end**  
   Let  $\mathbf{y}_{j+1} = \mathbf{y}_j - \mathbf{P}\mathbf{x}_{k_j}$ .  
**end**  
**Output:**  $\mathbf{y} - \mathbf{y}_{\text{last}}$ .

---

Note that this algorithm and its proof are quite similar to Theorem 12.3.4 but modified to scale parameters over an outer loop. By changing the parameter  $t$  we can decrease the dependence of the initial error.<sup>10</sup>

**Theorem 12.7.2.** *Assume the conditions in Assumption 12.7.1 are satisfied. Let  $\mathcal{T}$  be the time needed to compute  $\mathbf{P}\mathbf{x}$  and  $\mathbf{P}^T\mathbf{x}$  and  $\mathbf{x}^\#$ . Then, **NonlinearProjection** outputs a vector  $\mathbf{x}$  with  $\|\mathbf{x}\| \leq (1 + \varepsilon) \min_{\mathbf{x} \in L} \|\mathbf{x} - \mathbf{y}\|$  and the algorithm takes time*

$$O\left(\|\mathbf{P}\|^2 K (\mathcal{T} + m) \left(\frac{1}{\varepsilon^2} + \log \|\mathbf{P}\|\right)\right).$$

*Proof.* We prove by induction on  $j$  that when  $2^{-(j-1)}\|\mathbf{P}\| \geq 1$  we have  $\|\mathbf{y}_j\| \leq (1 + 2^{-j}\|\mathbf{P}\|) \text{OPT}$ .

For the base case ( $j = 0$ ), (12.7.2) shows that  $\|\mathbf{y}_0\| \leq (1 + \|\mathbf{P}\|) \text{OPT}$ .

For the inductive case we assume that the assertion holds for some  $j$ . We start by bounding the corresponding  $R$  in Theorem 12.3.1 for  $g_j$ , which we denote  $R_j$ . Note that

$$g_j(\mathbf{x}_0) = f_{t_j}(-\mathbf{y}_j) \leq \|\mathbf{y}_j\| + Kt_j \leq (1 + 2^{-j}\|\mathbf{P}\|) \text{OPT} + Kt_j.$$

Hence, the condition that  $g_j(\mathbf{x}) \leq g_j(\mathbf{x}_0)$  implies that

$$\|\mathbf{P}\mathbf{x} - \mathbf{y}_j\| \leq (1 + 2^{-j}\|\mathbf{P}\|) \text{OPT} + Kt_j.$$

Take any  $\mathbf{y} \in X^*$ , let  $\mathbf{c} = \mathbf{x} - \mathbf{P}\mathbf{x} + \mathbf{y}$ , and note that  $\mathbf{P}\mathbf{c} = \mathbf{P}\mathbf{y}$  and therefore  $\mathbf{c} \in X^*$ . Using these facts, we can bound  $R_j$  as follows

$$\begin{aligned} R_j &= \max_{\mathbf{x} \in \mathbb{R}^E : g_j(\mathbf{x}) \leq g_j(\mathbf{x}_0)} \|\mathbf{x} - \mathbf{c}\| \\ &\leq \max_{\mathbf{x} \in \mathbb{R}^E : g_j(\mathbf{x}) \leq g_j(\mathbf{x}_0)} \|\mathbf{P}\mathbf{x}\| + \|\mathbf{P}\mathbf{y}\| \\ &\leq 2\|\mathbf{y}_0\| + \|\mathbf{P}\mathbf{x} - \mathbf{y}_j\| + \|\mathbf{P}\mathbf{y} - \mathbf{y}_j\| \\ &\leq 4(1 + 2^{-j}\|\mathbf{P}\|) \text{OPT} + 2Kt_j. \end{aligned}$$

---

<sup>10</sup>This is an idea that has been applied previously to solve linear programming problems [207].

Similar to Lemma 12.3.3, the Lipschitz constant  $L_j$  of  $g_j$  is  $\|\mathbf{P}\|^2/t_j$ . Hence, Theorem 12.3.1 shows that

$$\begin{aligned} g_j(\mathbf{x}_{k_j}) &\leq \min_{\mathbf{x}} g_j(\mathbf{x}) + \frac{2 \cdot L_j \cdot R_j^2}{k_j + 4} \\ &\leq \min_{\mathbf{x}} \|\mathbf{P}\mathbf{x} - \mathbf{y}_j\| + \frac{2 \cdot L_j \cdot R_j^2}{k_j + 4} + Kt_j \end{aligned}$$

So, we have

$$\begin{aligned} \|\mathbf{P}\mathbf{x}_{k_j} - \mathbf{y}_j\| &\leq f_{t_j}(\mathbf{P}\mathbf{x}_{k_j} - \mathbf{y}_j) \\ &\leq \text{OPT} + Kt_j + \frac{2\|\mathbf{P}\|^2}{t_j(k_j + 4)} (4(1 + 2^{-j}\|\mathbf{P}\|)\text{OPT} + 2Kt_j)^2. \end{aligned}$$

When  $2^{-j}\|\mathbf{P}\| > 1$ , we have  $t_j = \frac{2^{-(j+2)}\|\mathbf{P}\|\text{OPT}}{K}$  and  $k_j = 3200\|\mathbf{P}\|^2K$  and hence

$$\|\mathbf{y}_{j+1}\| = \|\mathbf{P}\mathbf{x}_{k_j} - \mathbf{y}_j\| \leq (1 + 2^{-j-1}\|\mathbf{P}\|)\text{OPT}.$$

When  $2^{-j}\|\mathbf{P}\| \leq 1$ , we have  $t_j = \frac{\varepsilon\text{OPT}}{2K}$  and  $k_j = \frac{800\|\mathbf{P}\|^2K}{\varepsilon^2}$  and hence

$$\|\mathbf{y}_{j+1}\| = \|\mathbf{P}\mathbf{x}_{k_j} - \mathbf{y}_j\| \leq (1 + \varepsilon)\text{OPT}.$$

Since  $\mathbf{y}_{\text{last}}$  is  $\mathbf{y}$  plus some vectors in  $L$ ,  $\mathbf{y} - \mathbf{y}_{\text{last}} \in L$  and  $\|\mathbf{t} - \mathbf{y}_{\text{last}} - \mathbf{y}\| = \|\mathbf{y}_{\text{last}}\| \leq (1 + \varepsilon)\text{OPT}$ .  $\square$

## ■ 12.7.2 Maximum Concurrent Flow

For an arbitrary set of demands  $\chi_i \in \mathbb{R}^V$  with  $\sum_{v \in V} \chi_i(v) = 0$  for  $i = 1, \dots, k$ , we wish to solve the following *maximum concurrent flow* problem

$$\max_{\alpha \in \mathbb{R}, \mathbf{f} \in \mathbb{R}^E} \alpha \quad \text{subject to} \quad \mathbf{B}^T \mathbf{f}_i = \alpha \chi_i \quad \text{and} \quad \|\mathbf{U}^{-1} \sum_{i=1}^k \mathbf{f}_i\|_{\infty} \leq 1.$$

Similar to Section 12.3.2, it is equivalent to the problem

$$\min_{\alpha \in \mathbb{R}^{E \times [k]}} \left\| \sum_{i=1}^k |\alpha_i + (\mathbf{Q}\mathbf{x})_i| \right\|_{\infty}$$

where  $\mathbf{Q}$  is a projection matrix onto the subspace  $\{\mathbf{B}^T \mathbf{U}\mathbf{x}_i = 0\}$ , the output maximum concurrent flow is

$$\mathbf{f}_i(\mathbf{x}) = \mathbf{U}(\alpha_i + (\mathbf{Q}\mathbf{x})_i) / \left\| \sum_{i=1}^k |\alpha_i + (\mathbf{Q}\mathbf{x})_i| \right\|_{\infty},$$

and  $\mathbf{U}\alpha_i$  is any flow such that  $\mathbf{B}^T \mathbf{U}\alpha_i = \chi_i$ . In order to apply **NonlinearProjection**, we need to find a regularized norm and a good projection matrix. Let us define the norm

$$\|\mathbf{x}\|_{1,\infty} = \max_{e \in E} \sum_{i=1}^k |x_i(e)|.$$

The problem is simply  $\|\alpha + \mathbf{Q}\mathbf{x}\|_{1,\infty}$  where  $\mathbf{Q}$  is a projection matrix from  $\mathbb{R}^{E \times [k]}$  to  $\mathbb{R}^{E \times [k]}$  onto some subspace. Since each copy  $\mathbb{R}^E$  is same, there is no reason that there is coupling in  $\mathbf{Q}$  between different copies of  $\mathbb{R}^E$ . In the next lemma, we formalize this by the fact that any good projection

matrix  $\mathbf{P}$  onto the subspace  $\{\mathbf{B}^T \mathbf{U} \vec{x} = 0\} \subset \mathbb{R}^E$  extends to a good projection  $\mathbf{Q}$  onto the subspace  $\{\mathbf{B}^T \mathbf{U} \mathbf{x}_i = 0\} \subset \mathbb{R}^{E \times [k]}$ . Therefore, we can simply extend the good circulation projection  $\mathbf{P}$  by formula  $(\mathbf{Q}\mathbf{x})_i = \mathbf{P}\mathbf{x}_i$ . Thus, the only last piece needed is a regularized  $\|\cdot\|_{1;\infty}$ . However, it turns out that smoothing via conjugate does not work well in this case because the dual space of  $\|\cdot\|_{1;\infty}$  involves with  $\|\cdot\|_\infty$ , which is unfavorable for this kind of smoothing procedure. It can be shown that there is no such good regularized  $\|\cdot\|_{1;\infty}$ . Therefore, we could not do  $O(m^{1+o(1)}k/\varepsilon^2)$  using this approach, however,  $O(m^{1+o(1)}k^2/\varepsilon^2)$  is possible by using a bad regularized  $\|\cdot\|_{1;\infty}$ . We believe the dependence of  $k$  can be improved to  $3/2$  using this approach by suitable using Nesterov algorithm because the  $\|\cdot\|_1$  space caused by the multicommodity is a favorable geometry for accelerated methods.

**Lemma 12.7.3.** *Let  $\text{smax}L1_t(\mathbf{x}) = \text{smax}_t \left( \sum_{i=1}^k \sqrt{(x_i(e))^2 + t^2} \right)$ . It is a convex continuously differentiable function. The Lipschitz constant of  $\nabla \text{smax}L1_t$  is  $\frac{2}{t}$  and*

$$\|\mathbf{x}\|_{1;\infty} - t \ln(2m) \leq \text{smax}L1_t(\mathbf{x}) \leq \|\mathbf{x}\|_{1;\infty} + kt.$$

*Proof.* 1) It is clear that  $\text{smax}L1_t$  is smooth.

2)  $\text{smax}L1_t$  is convex.

Since  $\text{smax}_t$  is increasing for positive values and  $\sqrt{x^2 + t^2}$  is convex, for any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{E \times [k]}$  and  $0 \leq t \leq 1$ , we have

$$\begin{aligned} \text{smax}L1_t(t\mathbf{x} + (1-t)\mathbf{y}) &= \text{smax}_t \left( \sum_{i=1}^k \sqrt{((tx_i + (1-t)y_i)(e))^2 + t^2} \right) \\ &\leq \text{smax}_t \left( \sum_{i=1}^k \left( t\sqrt{(x_i(e))^2 + t^2} + (1-t)\sqrt{(y_i(e))^2 + t^2} \right) \right) \\ &\leq t\text{smax}L1_t(\mathbf{x}) + (1-t)\text{smax}L1_t(\mathbf{y}). \end{aligned}$$

3) The Lipschitz constant of  $\nabla \text{smax}L1_t$  is  $\frac{2}{t}$ .

Note that  $\text{smax}_t$  (not its gradient) has Lipschitz constant 1 because for any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^E$ ,

$$\begin{aligned} &|\text{smax}_t(\mathbf{x}) - \text{smax}_t(\mathbf{y})| \\ &= \left| t \ln \left( \frac{\sum_{e \in E} \left( \exp\left(-\frac{x(e)}{t}\right) + \exp\left(\frac{x(e)}{t}\right) \right)}{2m} \right) - t \ln \left( \frac{\sum_{e \in E} \left( \exp\left(-\frac{y(e)}{t}\right) + \exp\left(\frac{y(e)}{t}\right) \right)}{2m} \right) \right| \\ &= t \left| \ln \left( \frac{\sum_{e \in E} \left( \exp\left(-\frac{x(e)}{t}\right) + \exp\left(\frac{x(e)}{t}\right) \right)}{\sum_{e \in E} \left( \exp\left(-\frac{y(e)}{t}\right) + \exp\left(\frac{y(e)}{t}\right) \right)} \right) \right| \\ &\leq t \left| \ln \left( \max_{e \in E} \exp\left(\frac{|x - y|(e)}{t}\right) \right) \right| \\ &= \|\mathbf{x} - \mathbf{y}\|_\infty. \end{aligned}$$

Also, by the definition of derivative, for any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  and  $t \in \mathbb{R}$ , we have

$$\text{smax}_t(\mathbf{x} + t\mathbf{y}) - \text{smax}_t(\mathbf{x}) = t \langle \nabla \text{smax}_t(\mathbf{x}), \mathbf{y} \rangle + o(t).$$

and it implies  $|\langle \nabla \text{smax}_t(\mathbf{x}), \mathbf{y} \rangle| \leq \|\mathbf{y}\|_\infty$  for arbitrary  $\mathbf{y}$  and hence

$$\|\nabla \text{smax}_t(\mathbf{x})\|_1 \leq 1. \tag{12.5}$$



For notational simplicity, let  $s_1 = \text{smax}L1_t$ ,  $s_2 = \text{smax}_t$  and  $s_3(x) = \sqrt{x^2 + t^2}$ . Thus, we have

$$s_1(\vec{x}) = s_2 \left( \sum_{i=1}^k s_3(x_i(e)) \right).$$

Now, we want to prove

$$\|\nabla s_1(\mathbf{x}) - \nabla s_1(\mathbf{y})\|_{\infty;1} \leq \frac{2}{t} \|\mathbf{x} - \mathbf{y}\|_{1;\infty}.$$

Note that

$$\frac{\partial s_1(\mathbf{x})}{\partial x_i(e)} = \frac{\partial s_2}{\partial e} \left( \sum_j s_3(x_j(e)) \right) \frac{ds_3}{dx}(x_i(e)).$$

Hence, we have

$$\begin{aligned} & \|\nabla s_1(\mathbf{x}) - \nabla s_1(\mathbf{y})\|_{\infty;1} \\ &= \sum_e \max_i \left| \frac{\partial s_2}{\partial e} \left( \sum_j s_3(x_j(e)) \right) \frac{ds_3}{dx}(x_i(e)) - \frac{\partial s_2}{\partial e} \left( \sum_j s_3(y_j(e)) \right) \frac{ds_3}{dx}(y_i(e)) \right| \\ &\leq \sum_e \left| \frac{\partial s_2}{\partial e} \left( \sum_j s_3(x_j(e)) \right) \right| \max_i \left| \frac{ds_3}{dx}(x_i(e)) - \frac{ds_3}{dx}(y_i(e)) \right| \\ &\quad + \sum_e \max_i \frac{ds_3}{dx}(y_i(e)) \left| \frac{\partial s_2}{\partial e} \left( \sum_j s_3(x_j(e)) \right) - \frac{\partial s_2}{\partial e} \left( \sum_j s_3(y_j(e)) \right) \right|. \end{aligned}$$

Since  $s_3$  has  $\frac{1}{t}$ -Lipschitz gradient, we have

$$\left| \frac{ds_3}{dx}(\mathbf{x}) - \frac{ds_3}{dx}(\mathbf{y}) \right| \leq \frac{1}{t} |\mathbf{x} - \mathbf{y}|.$$

By (12.5), we have

$$\sum_e \left| \frac{\partial s_2}{\partial e}(\mathbf{x}(e)) \right| \leq 1.$$

Hence, we have

$$\begin{aligned} & \sum_e \max_i \left| \frac{ds_3}{dx}(x_i(e)) - \frac{ds_3}{dx}(y_i(e)) \right| \left| \frac{\partial s_2}{\partial e} \left( \sum_i s_3(x_i(e)) \right) \right| \\ &\leq \frac{1}{t} \max_{i,e} |x_i(e) - y_i(e)| \sum_e \left| \frac{\partial s_2}{\partial e} \left( \sum_i s_3(x_i(e)) \right) \right| \\ &= \frac{1}{t} \|\vec{x} - \vec{y}\|_{1;\infty}. \end{aligned}$$

Since  $s_3$  is 1-Lipschitz, we have

$$\left| \frac{ds_3}{dx} \right| \leq 1.$$

Since  $s_2$  has  $\frac{1}{t}$ -Lipschitz gradient in  $\|\cdot\|_\infty$ , we have

$$\sum_e \left| \frac{\partial s_2}{\partial e}(x) - \frac{\partial s_2}{\partial e}(y) \right| \leq \frac{1}{t} \|\mathbf{x} - \mathbf{y}\|_\infty.$$

Hence, we have

$$\begin{aligned} & \sum_e \max_i \frac{ds_3}{dx}(y_i(e)) \left| \frac{\partial s_2}{\partial e} \left( \sum_j s_3(x_j(e)) \right) - \frac{\partial s_2}{\partial e} \left( \sum_j s_3(y_j(e)) \right) \right| \\ & \leq \sum_e \left| \frac{\partial s_2}{\partial e} \left( \sum_j s_3(x_j(e)) \right) - \frac{\partial s_2}{\partial e} \left( \sum_j s_3(y_j(e)) \right) \right| \\ & \leq \frac{1}{t} \left\| \sum_i s_3(x_i(e)) - \sum_i s_3(y_i(e)) \right\|_\infty \\ & \leq \frac{1}{t} \left\| \sum_i |x_i(e) - y_i(e)| \right\| \\ & = \frac{1}{t} \|\mathbf{x} - \mathbf{y}\|_{1,\infty} \end{aligned}$$

Therefore, we have

$$\|\nabla s_1(\mathbf{x}) - \nabla s_1(\mathbf{y})\|_{\infty;1} \leq \frac{2}{t} \|\mathbf{x} - \mathbf{y}\|_{1,\infty}.$$

4) Using the fact that

$$\|x(e)\| \leq \sum_{i=1}^k \sqrt{(x_i(e))^2 + t^2} \leq \|x(e)\|_1 + kt$$

and  $\text{smax}$  is 1-Lipschitz, we have

$$\|\mathbf{x}\|_{1,\infty} - t \ln(2m) \leq \text{smax} L 1_t(\mathbf{x}) \leq \|\mathbf{x}\|_{1,\infty} + kt.$$

□

The last thing needed is to check is that the  $\#$  operator is easy to compute.

**Lemma 12.7.4.** *In  $\|\cdot\|_{1,\infty}$ , the  $\#$  operator is given by an explicit formula*

$$\left(\mathbf{x}^\#\right)_i(e) = \begin{cases} \|\mathbf{x}\|_{\infty;1} \text{sign}(x_i(e)) & \text{if } i \text{ is the smallest index such that } \max_j x_j(e) = x_i(e) \\ 0 & \text{otherwise} \end{cases}.$$

*Proof.* It can be proved by direct computation. □

Now, all the conditions in the Assumption 12.7.1 are satisfied. Therefore, Theorem 12.7.2 and Theorem 12.4.12 gives us the following theorem:

**Theorem 12.7.5.** *Given an undirected capacitated graph  $G = (V, E, \boldsymbol{\mu})$  with capacity ratio  $U$ . Assume  $U = \text{poly}(|V|)$ . There is an algorithm finds an  $(1 - \varepsilon)$  approximate Maximum Concurrent Flow in time*

$$O\left(\frac{k^2}{\varepsilon^2} |E| 2^{O(\sqrt{\log |V| \log \log |V|})}\right).$$

*Proof.* Let  $\mathbf{A}$  be the oblivious routing algorithm given by Theorem 12.4.12. And we have  $\rho(\mathbf{A}) \leq 2^{O(\sqrt{\log|V|\log\log|V|})}$ . Let us define the scaled circulation projection matrix  $\mathbf{P} = \mathbf{I} - \mathbf{U}\mathbf{A}\mathbf{B}^T\mathbf{U}^{-1}$ . Lemma 12.4.5 shows that  $\|\mathbf{P}\|_\infty \leq 1 + 2^{O(\sqrt{\log|V|\log\log|V|})}$ .

Let the multi-commodity circulation projection matrix  $\mathbf{Q} : \mathbb{R}^{E \times [k]} \rightarrow \mathbb{R}^{E \times [k]}$  defined by  $(\mathbf{Q}\mathbf{x})_i = \mathbf{P}\mathbf{x}_i$ . Note that the definition of  $\|\mathbf{Q}\|_{1;\infty}$  is similar to  $\rho(\mathbf{Q})$ . By similar proof as Lemma 12.4.3, we have  $\|\mathbf{Q}\|_{1;\infty} = \|\mathbf{P}\|_\infty$ . Hence, we have  $\|\mathbf{Q}\|_{1;\infty} \leq 1 + 2^{O(\sqrt{\log|V|\log\log|V|})}$ . Also, since  $\mathbf{P}$  is a projection matrix on the subspace  $\{\mathbf{x} \in \mathbb{R}^E : \mathbf{B}^T\mathbf{U}\mathbf{x} = 0\}$ ,  $\mathbf{Q}$  is a projection matrix on the subspace  $\{\mathbf{x} \in \mathbb{R}^{E \times [k]} : \mathbf{B}^T\mathbf{U}\mathbf{x}_i = 0\}$ .

By Lemma 12.7.3, the function  $\text{smax}L1_t(\mathbf{x})$  is a convex continuously differentiable function such that the Lipschitz constant of  $\nabla\text{smax}L1_t$  is  $\frac{2}{t}$  and

$$\|\mathbf{x}\|_{1;\infty} - t \ln(2m) \leq \text{smax}L1_t(\mathbf{x}) \leq \|\mathbf{x}\|_{1;\infty} + kt.$$

Given an arbitrary set of demands  $\chi_i \in \mathbb{R}^V$ , we find a vector  $\vec{y}$  such that

$$\mathbf{B}^T\mathbf{U}\mathbf{y} = -\chi_i.$$

Then, we use the **NonlinearProjection** to solve

$$\min_{\mathbf{B}^T\mathbf{U}\mathbf{x}=0} \|\mathbf{x} - \mathbf{y}\|_{1;\infty}$$

using a family of functions  $\text{smax}L1_t(\mathbf{x}) + t \ln(2n)$  and the projection matrix  $\mathbf{Q}$ . Since each iteration involves calculation of gradients and  $\#$  operator, it takes  $O(mk)$  each iteration. And it takes  $\tilde{O}(\|\mathbf{Q}\|_{1;\infty}^2 K/\varepsilon^2)$  iterations in total where  $K = k + \ln(2m)$ . In total, it **NonlinearProjection** outputs a  $(1 + \varepsilon)$  approximate minimizer  $\mathbf{x}$  in time

$$O\left(\frac{k^2}{\varepsilon^2} |E| 2^{O(\sqrt{\log|V|\log\log|V|})}\right).$$

And it gives a  $(1 - \varepsilon)$  approximate maximum concurrent flow  $\vec{f}_i$  by a direct formula.  $\square$

## ■ 12.8 Appendix

In this section, we present some basic fact used in this chapter about norm and dual norm. Also, we presented some lemmas about convex functions with Lipschitz gradient. See [35, 206] for comprehensive discussion.

### ■ 12.8.1 Norms

**Fact 12.8.1.**  $\mathbf{x} = 0 \Leftrightarrow \mathbf{x}^\# = 0$ .

*Proof.* If  $\mathbf{x} = 0$  then  $\forall \mathbf{s} \neq 0$  we have  $\langle \mathbf{x}, \mathbf{s} \rangle - \frac{1}{2}\|\mathbf{s}\|^2 < 0$  but  $\langle \mathbf{x}, \mathbf{x} \rangle - \frac{1}{2}\|\mathbf{x}\|^2 = 0$ . So we have  $\mathbf{x}^\# = 0$ . If  $\mathbf{x} \neq 0$  then let  $\mathbf{s} = \frac{\langle \mathbf{x}, \mathbf{x} \rangle}{\|\mathbf{x}\|^2} \mathbf{x}$  with this choice we have  $\langle \mathbf{x}, \mathbf{s} \rangle - \frac{1}{2}\|\mathbf{s}\|^2 = \frac{1}{2} \frac{\langle \mathbf{x}, \mathbf{x} \rangle^2}{\|\mathbf{x}\|^2} > 0$ . However, for  $\mathbf{s} = 0$  we have that  $\langle \mathbf{x}, \mathbf{s} \rangle - \frac{1}{2}\|\mathbf{s}\|^2 = 0$  therefore we have  $\mathbf{x}^\# \neq 0$ .  $\square$

**Fact 12.8.2.**  $\forall \mathbf{x} \in \mathbb{R}^n \quad : \quad \langle \mathbf{x}, \mathbf{x}^\# \rangle = \|\mathbf{x}^\#\|^2$ .

*Proof.* If  $\mathbf{x} = 0$  then  $\mathbf{x}^\# = 0$  by Fact 12.8.1 and we have the result. Otherwise, again by Fact 12.8.1

we know that  $\mathbf{x}^\# \neq 0$  and therefore by the definition of  $\mathbf{x}^\#$  we have

$$1 = \arg \max_{c \in \mathbb{R}} \langle x, c \cdot \mathbf{x}^\# \rangle - \frac{1}{2} \|c \cdot \mathbf{x}^\#\|^2 = \arg \max_{c \in \mathbb{R}} c \cdot \langle x, \mathbf{x}^\# \rangle - \frac{c^2}{2} \|\mathbf{x}^\#\|^2$$

Setting the derivative of with respect to  $c$  to 0 we get that  $1 = c = \frac{\langle x, \mathbf{x}^\# \rangle}{\|\mathbf{x}^\#\|^2}$ .  $\square$

**Fact 12.8.3.**  $\forall \mathbf{x} \in \mathbb{R}^n$  :  $\|\mathbf{x}\|^* = \|\mathbf{x}^\#\|$ .

*Proof.* Note that if  $\mathbf{x} = 0$  then the fact follows from Fact (12.8.1) otherwise we have

$$\|\mathbf{x}\|^* = \max_{\|y\| \leq 1} \langle x, y \rangle = \max_{\|y\|=1} \langle x, y \rangle \leq \max_{y \in \mathbb{R}^n} \frac{\langle x, y \rangle}{\|y\|}$$

From this it is clear that  $\|\mathbf{x}\|^* \geq \|\mathbf{x}^\#\|$ . To see the other direction consider a  $\mathbf{y}$  that maximizes the above and let  $\mathbf{z} = \frac{\langle x, \mathbf{y} \rangle}{\|\mathbf{y}\|^2} \mathbf{y}$

$$\langle \mathbf{x}, \mathbf{z} \rangle - \frac{1}{2} \|\mathbf{z}\|^2 \leq \langle \mathbf{x}, \mathbf{x}^\# \rangle - \frac{1}{2} \|\mathbf{x}^\#\|^2$$

and therefore

$$\|\mathbf{x}\|^{*2} - \frac{1}{2} \|\mathbf{x}\|^{*2} \leq \frac{1}{2} \|\mathbf{x}^\#\|^2$$

$\square$

**Fact 12.8.4** (Cauchy Schwarz).  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  :  $\langle \mathbf{y}, \mathbf{x} \rangle \leq \|\mathbf{y}\|^* \|\mathbf{x}\|$ .

*Proof.* By the definition of dual norm, for all  $\|\mathbf{x}\| = 1$ , we have  $\langle \mathbf{y}, \mathbf{x} \rangle \leq \|\mathbf{y}\|^*$ . Hence, it follows by linearity of both side.  $\square$

## ■ 12.8.2 Functions with Lipschitz Gradient

**Lemma 12.8.5.** *Let  $f$  be a continuously differentiable convex function. Then, the following are equivalence:*

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n \quad : \quad \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^* \leq L \cdot \|\mathbf{x} - \mathbf{y}\|$$

and

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n \quad : \quad f(\mathbf{x}) \leq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

For any such  $f$  and any  $\mathbf{x} \in \mathbb{R}^n$ , we have

$$f\left(\mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x})^\#\right) \leq f(\mathbf{x}) - \frac{1}{2L} \|\nabla f(\mathbf{x})\|^{*2}.$$

*Proof.* From the first condition, we have

$$\begin{aligned} f(\mathbf{y}) &= f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \int_0^1 \langle \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle dt \\ &\leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \int_0^1 \|\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})\|^* \|\mathbf{y} - \mathbf{x}\| dt \\ &\leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \int_0^1 Lt \|\mathbf{y} - \mathbf{x}\|^2 dt \\ &= f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2. \end{aligned}$$

Given the second condition. For any  $\mathbf{x} \in \mathbb{R}^n$ . let  $\phi_{\mathbf{x}}(\mathbf{y}) = f(\mathbf{y}) - \langle \nabla f(\mathbf{x}), \mathbf{y} \rangle$ . From the convexity of  $f$ , for any  $\mathbf{y} \in \mathbb{R}^n$

$$f(\mathbf{y}) - f(\mathbf{x}) \geq \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle.$$

Hence,  $\mathbf{x}$  is a minimizer of  $\phi_{\mathbf{x}}$ . Hence, we have

$$\begin{aligned} \phi_{\mathbf{x}}(\mathbf{x}) &\leq \phi_{\mathbf{x}}(\mathbf{y} - \frac{1}{L} \nabla \phi_{\mathbf{x}}(\mathbf{y})^{\#}) \\ &\leq \phi_{\mathbf{x}}(\mathbf{y}) - \langle \nabla \phi_{\mathbf{x}}(\mathbf{y}), \frac{1}{L} \nabla \phi_{\mathbf{x}}(\mathbf{y})^{\#} \rangle + \frac{L}{2} \|\frac{1}{L} \nabla \phi_{\mathbf{x}}(\mathbf{y})^{\#}\|^2 && \text{(First part of this lemma)} \\ &= \phi_{\mathbf{x}}(\mathbf{y}) - \frac{1}{2L} \|\nabla \phi_{\mathbf{x}}(\mathbf{y})^{\#}\|^2 \\ &= \phi_{\mathbf{x}}(\mathbf{y}) - \frac{1}{2L} (\|\nabla \phi_{\mathbf{x}}(\mathbf{y})\|^*)^2. \end{aligned}$$

Hence,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2L} (\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|^*)^2.$$

Adding up this inequality with  $\mathbf{x}$  and  $\mathbf{y}$  interchanged, we have

$$\begin{aligned} \frac{1}{L} (\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|^*)^2 &\leq \langle \nabla f(\mathbf{y}) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \\ &\leq \|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|^* \|\mathbf{y} - \mathbf{x}\|. \end{aligned}$$

The last inequality follows from similar proof in above for  $\phi_{\mathbf{x}}$ . □

The next lemma relate the Hessian of function with the Lipschitz parameter  $L$  and this lemma gives us a easy way to compute  $L$ .

**Lemma 12.8.6.** *Let  $f$  be a twice differentiable function such that for any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$*

$$0 \leq \mathbf{y}^T (\nabla^2 f(\mathbf{x})) \mathbf{y} \leq L \|\mathbf{y}\|^2.$$

*Then,  $f$  is convex and the gradient of  $f$  is Lipschitz continuous with Lipschitz parameter  $L$ .*

*Proof.* Similarly to Lemma 12.8.5, we have

$$\begin{aligned} f(\mathbf{y}) &= f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \int_0^1 \langle \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle dt \\ &= f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \int_0^1 t(\mathbf{y} - \mathbf{x})^T \nabla^2 f(\mathbf{x} + \theta_t(\mathbf{y} - \mathbf{x})) (\mathbf{y} - \mathbf{x}) dt \end{aligned}$$

where the  $0 \leq \theta_t \leq t$  comes from mean value theorem. By the assumption, we have

$$\begin{aligned} f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle &\leq f(\mathbf{y}) \\ &\leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \int_0^1 tL \|\mathbf{y} - \mathbf{x}\|^2 dt \\ &\leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2. \end{aligned}$$

And the conclusion follows from Lemma 12.8.5. □



---

---

## Connection Laplacian In Nearly-Linear Time

### ■ 13.1 Introduction

In this chapter, we introduce the *sparsified Cholesky* and *sparsified multigrid* algorithms for solving systems of linear equations. Two advantages of these algorithms over other nearly-linear time algorithms for solving systems of equations in Laplacian matrices [256, 243, 147, 146, 141, 57] are:

1. They give nearly-linear time algorithms for solving systems of equations in a much broader class of matrices—the connection Laplacians and Hermitian block diagonally dominant matrices. Connection Laplacians [238, 29] are a generalization of graph Laplacians that arise in many applications, including celebrated work on cryo-electron microscopy [237, 235, 275], phase retrieval [7, 186], and many image processing problems (e.g. [221, 16]). Previous algorithms for solving systems of equations in graph Laplacians cannot be extended to solve equations in connection Laplacians because the previous algorithms relied on some form of low stretch spanning trees, a concept that has no analog for the more general connection Laplacians.
2. They provide linear-sized approximate inverses of connection Laplacian matrices. That is, for every  $n$ -dimensional connection Laplacian  $\mathbf{M}$ , the sparsified Cholesky factorization algorithm produces a block-upper-triangular matrix  $\mathbf{U}$  and block diagonal  $\mathbf{D}$  with  $O(n)$  nonzero entries so that  $\mathbf{U}^T \mathbf{D} \mathbf{U}$  is a constant-factor approximation of  $\mathbf{M}$ . Such matrices  $\mathbf{U}$  and  $\mathbf{D}$  allow one to solve systems of equations in  $\mathbf{M}$  to accuracy  $\varepsilon$  in time  $O(m \log \varepsilon^{-1})$ , where  $m$  is the number of nonzero entries of  $\mathbf{M}$ . Even for ordinary Laplacian matrices, the existence of such approximate inverses is entirely new. The one caveat of this result is that we do not yet know how to compute those approximate inverses in nearly linear time.

The sparsified Cholesky and sparsified multigrid algorithms work by sparsifying the matrices produced during Gaussian elimination. Recall that Cholesky factorization is the version of Gaussian elimination for symmetric matrices, and that the high running time of Gaussian elimination comes from *fill*—new nonzero matrix entries that are created by row operations. If Gaussian elimination never produced rows with a super-constant number of entries, then it would run in linear time. The sparsified Cholesky algorithm accelerates Gaussian elimination by sparsifying the rows that are produced by elimination, thereby guaranteeing that the elimination will be fast. Sparsified Cholesky is inspired by one of the major advances in algorithms for solving linear equations in Laplacian matrices—the Incomplete Cholesky factorization (ICC) [191]. However, ICC merely drops entries produced by elimination, whereas sparsification also carefully increases ones that remain. The difference is crucial, and is why ICC does not provide a nearly-linear time solver.

To control the error introduced by sparsification, we have to be careful not to do it too often. This means that our algorithm actually chooses a large set of rows and columns to eliminate at once, and then sparsifies the result. This is the basis of our first algorithms, which establish the existence of linear time solvers and linear sized approximate inverses, after precomputation. This precomputation

is analogous to the procedure of computing a matrix inverse: the approximate inverse takes much less time to apply than to compute.

To produce entire algorithms that run in nearly linear time (both to compute and apply the approximate inverse) requires a little more work. To avoid the work of computing the matrix obtained by eliminating the large set of rows and columns, we design a fast algorithm for approximating it quickly. This resulting algorithm produces a solver routine that does a little more work at each level, and so resembles a multigrid V-cycle [252]. We call the resulting algorithm the *sparsified multigrid*. We note that Krishnan, Fattal, and Szeliski [150] present experimental results from the use of a sparsification heuristic in a multigrid algorithm for solving problems in computer vision.

Our new algorithms are most closely related to the Laplacian solver recently introduced by Peng and Spielman [223]: unlike the other Laplacian solvers, they rely only on sparsification and do not directly rely on “support theory” preconditioners or any form of low stretch spanning trees. This is why our algorithms can solve a much broader family of linear systems. To sparsify without using graph theoretic algorithms, we employ the algorithm in Chapter 3 that allows us to sparsify a matrix by solving systems of equations in a subsampled matrix.

### ■ 13.1.1 Connection Laplacians and Block DD Matrices

In this section, we define block diagonally dominant (bDD) matrices—the most general family of matrices such that the associated systems of linear equations can be solved by our algorithms. We begin by defining our motivating case: the connection Laplacians

Connection Laplacians may be thought of as a generalization of graph Laplacians where every vertex is associated with a vector, instead of a real number, and every edge is associated with a unitary matrix. Like graph Laplacians, they describe a natural quadratic form. Let  $\mathbf{M}_{[i,j]} \in \mathbb{C}^{r \times r}$  be the unitary matrix associated with edge  $(i, j)$ , and let  $w_{i,j}$  be the (nonnegative) weight of edge  $(i, j)$ . We require that  $\mathbf{M}_{[i,j]} = \mathbf{M}_{[j,i]}^*$ , where  $*$  denotes the conjugate transpose. The quadratic form associated with this connection Laplacian is a function of vectors  $\mathbf{v}^{(i)} \in \mathbb{C}^r$ , one for each vertex  $i$ , that equals

$$\sum_{(i,j) \in E} w_{i,j} \|\mathbf{v}^{(i)} - \mathbf{M}_{[i,j]} \mathbf{v}^{(j)}\|^2.$$

The matrix corresponding to this quadratic form is a block matrix with blocks  $\mathbf{M}_{[i,j]}$ . Most applications of the connection Laplacian require one to either solve systems of linear equations in this matrix, or to compute approximations of its smallest eigenvalues and eigenvectors. By applying the inverse power method (or inverse Lanczos), we can perform these eigenvector calculations by solving a logarithmic number of linear systems in the matrix (see [243, Section 7]).

The matrices obtained from the connection Laplacian are a special case of block diagonally dominant (bDD) matrices, which we now define. Throughout this chapter, we consider block matrices having entries in  $\mathbb{C}^{r \times r}$ , where  $r > 0$  is a fixed integer. We say that  $\mathbf{M} \in (\mathbb{C}^{r \times r})^{m \times n}$  has  $m$  block-rows, and  $n$  block-columns. For  $i \in [m], j \in [n]$ , we let  $\mathbf{M}_{[i,j]} \in \mathbb{C}^{r \times r}$  denote the  $i, j$ -block in  $\mathbf{M}$ , and  $\mathbf{M}_{[j]}$  denote the  $j^{\text{th}}$  block-column, i.e.,  $\mathbf{M}_{[j]} = [(\mathbf{M}_{[1,j]})^*, (\mathbf{M}_{[2,j]})^*, \dots, (\mathbf{M}_{[m,j]})^*]^*$ . For sets  $F \subseteq [m], C \subseteq [n]$ , we let  $\mathbf{M}_{[F,C]} \in (\mathbb{C}^{r \times r})^{|F| \times |C|}$  denote the block-submatrix with blocks  $\mathbf{M}_{[i,j]}$  for  $i \in F, j \in C$ .  $\mathbf{M}$  is block-diagonal, if  $\mathbf{M}_{[i,j]} = 0$  for  $i \neq j$ . *We emphasize that all computations are done over  $\mathbb{C}$ , not over a matrix group.*

**Definition 13.1.1.** A Hermitian block-matrix  $\mathbf{M} \in (\mathbb{C}^{r \times r})^{n \times n}$  is block diagonally dominant (or bDD) if

$$\text{for all } i \in [n], \quad \mathbf{M}_{[i,i]} \succeq \mathbf{I}_r \cdot \sum_{j:j \neq i} \|\mathbf{M}_{[i,j]}\|_2,$$

where  $\mathbf{I}_r \in \mathbb{C}^{r \times r}$  denotes the identity matrix.



For the rest of this chapter, we use  $\|\cdot\|$  to denote  $\|\cdot\|_2$ . Equivalently, a Hermitian  $\mathbf{M}$  is a bDD matrix iff it can be written as  $\mathbf{D} - \mathbf{A}$  where  $\mathbf{D}$  is block-diagonal, and  $\mathbf{D}_{[i,i]} \succeq \mathbf{I}_r \sum_j \|\mathbf{A}_{[i,j]}\|$  (see Lemma 13.6.1).

Throughout this chapter we treat  $r$  as a constant. The hidden dependence on  $r$ , of the running times of the algorithms we present, is polynomial. The major results of this chapter are the following.

**Theorem 13.1.2** (Sparsified Multigrid). *There is an algorithm that, when given a bDD matrix  $\mathbf{M}$  with  $n$  block rows and  $m$  nonzero blocks, produces a solver for  $\mathbf{M}$  in  $O(m \log n + n \log^{2+o(1)} n)$  work and  $O(n^{o(1)})$  depth so that the solver finds  $\varepsilon$ -approximate solutions to systems of equations in  $\mathbf{M}$  in  $O((m + n \log^{1+o(1)} n) \log(1/\varepsilon))$  work and  $O(\log^2 n \log \log n \log(1/\varepsilon))$  depth.*

Previously, the existence of nearly-linear time solvers was even unknown for the 1-dimensional case ( $r = 1$ ) when the off-diagonal entries were allowed to be complex numbers.

We can use the above algorithm to find approximations of the smallest eigenvalues and eigenvectors of such matrices at an additional logarithmic cost.

**Theorem 13.1.3** (Sparsified Cholesky). *For every bDD matrix  $\mathbf{M}$  with  $n$  block-rows there exists a diagonal matrix  $\mathbf{D}$  and an upper triangular matrix  $\mathbf{U}$  with  $O(n)$  nonzero blocks so that*

$$\mathbf{U}^T \mathbf{D} \mathbf{U} \approx_{3/4} \mathbf{M}.$$

Moreover, linear equations in  $\mathbf{U}$ ,  $\mathbf{U}^T$ , and  $\mathbf{D}$  can be solved with linear work in depth  $O(\log^2 n)$ , and these matrices can be computed in polynomial time.

These matrices allow one to solve systems of linear equations in  $\mathbf{M}$  to  $\varepsilon$ -accuracy in parallel time  $O(\log^2 n \log^{-1} \varepsilon)$  and work  $O(m \log^{-1} \varepsilon)$ . Results of this form were previously unknown even for graph Laplacians.

In the next two sections we explain the ideas used to prove these theorems. Proofs may be found in the appendices that follow.

## ■ 13.2 Preliminaries

We say that  $\mathbf{A}$  is an  $\varepsilon$ -approximation of  $\mathbf{B}$ , written  $\mathbf{A} \approx_\varepsilon \mathbf{B}$ , if  $e^\varepsilon \mathbf{B} \succcurlyeq \mathbf{A} \succcurlyeq e^{-\varepsilon} \mathbf{B}$ . This relation is symmetric. We say that  $\tilde{\mathbf{x}}$  is an  $\varepsilon$ -approximate solution to the system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  if  $\|\tilde{\mathbf{x}} - \mathbf{A}^{-1}\mathbf{b}\|_{\mathbf{A}} \leq \varepsilon \|\mathbf{x}\|_{\mathbf{A}}$ , where  $\|\mathbf{x}\|_{\mathbf{A}} = (\mathbf{x}^T \mathbf{A} \mathbf{x})^{1/2}$ . If  $\varepsilon < 1/2$ ,  $\mathbf{A} \approx_\varepsilon \mathbf{B}$  and  $\mathbf{B}\tilde{\mathbf{x}} = \mathbf{b}$ , then  $\tilde{\mathbf{x}}$  is a  $2\varepsilon$  approximate solution to  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .

**Fact 13.2.1.** *If  $\mathbf{A} \approx_\varepsilon \mathbf{B}$  and  $\mathbf{B} \approx_\delta \mathbf{C}$ , then  $\mathbf{A} \approx_{\varepsilon+\delta} \mathbf{C}$ .*

We now address one technicality of dealing with bDD matrices: it is not immediate whether or not a bDD matrix is singular. Moreover, if it is singular, the structure of its null space is not immediately clear either. Throughout the rest of this chapter, we will consider bDD matrices to which a small multiple of the identity have been added. These matrices will be nonsingular. To reduce the problem of solving equations in a general bDD matrix  $\mathbf{M}$  to that of solving equations in a nonsingular matrix, we require an estimate of the smallest nonzero eigenvalue of  $\mathbf{M}$ .

**Fact 13.2.2.** *Suppose that all nonzero eigenvalues of  $\mathbf{M}$  are at least  $\mu$  and  $\mathbf{Z} \approx_\varepsilon (\mathbf{M} + \varepsilon\mu\mathbf{I})^{-1}$  for some  $0 < \varepsilon < 1/2$ . Given any  $\mathbf{b}$  such that  $\mathbf{M}\mathbf{x} = \mathbf{b}$  for some  $\mathbf{x}$ , we have*

$$\|\mathbf{x} - \mathbf{Z}\mathbf{b}\|_{\mathbf{M}}^2 \leq 6\varepsilon \|\mathbf{x}\|_{\mathbf{M}}^2.$$

Hence, we can solve systems in  $\mathbf{M}$  by approximately solving systems in  $\mathbf{M} + \varepsilon\mu\mathbf{I}$ . Any lower bound  $\mu$  on the smallest nonzero eigenvalue of  $\mathbf{M}$  will suffice. It only impacts the running times of our algorithms in the numerical precision with which one must carry out the computations.

The above fact allows us to solve systems in  $\mathbf{M}$  that have a solution. If  $\mathbf{M}$  is singular and we want to apply its pseudoinverse (that is, to find the closest solution in the range of  $\mathbf{M}$ ), we can do so by pre and post multiplying by  $\mathbf{M}$  to project onto its range. The resulting algorithm, which is implicit in the following fact, requires applying a solver for  $\mathbf{M}$  three times. It also takes  $O(\log \kappa)$  times as long to run, where  $\kappa$  is the finite condition number<sup>1</sup> of  $\mathbf{M}$ .

**Fact 13.2.3.** *Let  $\kappa$  be an upper bound on the finite condition number of  $\mathbf{M}$ . Given an error parameter  $0 < \varepsilon < 1$ , let  $\delta = \varepsilon/56\kappa^3$ . If all nonzero eigenvalues of  $\mathbf{M}$  are at least  $\mu$ , and if  $\mathbf{Z} \approx_\delta (\mathbf{M} + \varepsilon\mu\mathbf{I})^{-1}$ , then*

$$\mathbf{M}\mathbf{Z}^3\mathbf{M} \approx_{4\varepsilon} \mathbf{M}^+.$$

### ■ 13.3 Overview of the Algorithms

We index the block rows and columns of a block matrix by a set of vertices (or indices)  $V$ . When we perform an elimination, we eliminate a set  $F \subset V$ , and let  $C = V - F$ . Here,  $F$  stands for “fine” and  $C$  stands for “coarse”. In contrast with standard multigrid methods, we will have  $|F| \leq |C|$ .

To describe block-Cholesky factorization, we write the matrix  $\mathbf{M}$  with the rows and columns in  $F$  first:

$$\mathbf{M} = \begin{bmatrix} \mathbf{M}_{[F,F]} & \mathbf{M}_{[F,C]} \\ \mathbf{M}_{[C,F]} & \mathbf{M}_{[C,C]} \end{bmatrix}.$$

Cholesky factorization writes the inverse of this matrix as

$$\mathbf{M}^{-1} = \begin{bmatrix} \mathbf{I} & -\mathbf{M}_{[F,F]}^{-1}\mathbf{M}_{[F,C]} \\ 0 & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{M}_{[F,F]}^{-1} & 0 \\ 0 & \text{Sc}(\mathbf{M}, F)^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{I} & 0 \\ -\mathbf{M}_{[C,F]}\mathbf{M}_{[F,F]}^{-1} & \mathbf{I} \end{bmatrix}, \quad (13.1)$$

where

$$\text{Sc}(\mathbf{M}, F) \stackrel{\text{def}}{=} \mathbf{M}_{[C,C]} - \mathbf{M}_{[C,F]}\mathbf{M}_{[F,F]}^{-1}\mathbf{M}_{[F,C]}$$

is the Schur complement of  $\mathbf{M}$  with respect to  $F$ .

Our algorithms rely on two fundamental facts about *bDD* matrices: that the Schur complement of a *bDD* matrix is a *bDD* matrix (Lemma 13.6.3) and that one can sparsify *bDD* matrices. The following theorem is implicit in [46].

**Theorem 13.3.1.** *For every  $\varepsilon \leq 1$ , every *bDD* matrix  $\mathbf{M}$  with  $n$   $r$ -dimensional block rows and columns can be  $\varepsilon$ -approximated by a *bDD* matrix having at most  $10nr/\varepsilon^2$  nonzero blocks.*

We use the identity (13.1) to reduce the problem of solving a system of equations in  $\mathbf{M}$  to that of solving equations in its Schur complement. The easiest part of this is multiplication by  $\mathbf{M}_{[C,F]}$ : the time is proportional to the number of nonzero entries in the submatrix, and we can sparsify  $\mathbf{M}$  to guarantee that this is small.

The costlier part of the reduction is the application of the inverse of  $\mathbf{M}_{[F,F]}$  three times. This would be fast if  $\mathbf{M}_{[F,F]}$  were block-diagonal, which corresponds to  $F$  being an independent set. We cannot find a sufficiently large independent set  $F$ , but we can find a large set that is almost independent. This results in a matrix  $\mathbf{M}_{[F,F]}$  that is well approximated by its diagonal, and thus linear equations

<sup>1</sup>The finite condition number is the ratio of the largest singular value to the smallest *nonzero* singular value.

in this matrix can be quickly solved to high accuracy by a few Jacobi iterations (see Theorem 13.3.1 and Section 13.9).

We can prove the existence of linear sized approximate inverses by explicitly computing  $Sc(\mathbf{M}, F)$ , sparsifying it, and then recursively applying the algorithm just described. To make this algorithm efficient, we must compute a sparse approximation to  $Sc(\mathbf{M}, F)$  without constructing  $Sc(\mathbf{M}, F)$ . This is the problem of *spectral vertex sparsification*, and we provide a fast algorithm for this task in Sections 13.3.4 and 13.11.

### ■ 13.3.1 Schur Complement Chains

We encode this recursive algorithm by a *Schur complement chain* (SCC). An SCC defines a linear operator that can be used to approximately solve equations in an initial matrix  $\mathbf{M}^{(0)}$ . If the *bDD* matrix  $\mathbf{M}$  is sparse, then it is the same as  $\mathbf{M}^{(0)}$ ; if not, then  $\mathbf{M}^{(0)}$  is a sparse approximation to  $\mathbf{M}$ . Let  $F_1$  be the first set of vertices eliminated,  $\mathbf{M}^{(1)}$  a sparse approximation to the Schur complement of  $\mathbf{M}^{(0)}$  with respect to  $F_1$ , and  $\mathbf{Z}^{(1)}$  an operators that approximates the inverse of  $\mathbf{M}_{[F_1, F_1]}^{(0)}$ .

**Definition 13.3.2** (Schur Complement Chain). An  $\varepsilon$ -Schur complement chain ( $\varepsilon$ -SCC) for a matrix  $\mathbf{M}^0$  indexed by vertex set  $V$  is a sequence of operators and subsets,

$$\left( (\mathbf{M}^{(1)}, \mathbf{Z}^{(1)}), \dots, (\mathbf{M}^{(d)}, \mathbf{Z}^{(d)}); F_1, \dots, F_d \right),$$

so that for  $C_0 = V$  and  $C_{i+1} = C_i \setminus F_{i+1}$ ,  $|C_d| \leq 1000$  and for  $1 \leq i \leq d$ ,

$$\mathbf{M}^{(i)} \approx_{\varepsilon_i} Sc\left(\mathbf{M}^{(i-1)}, F_i\right) \quad \text{and} \quad 0 \preceq (\mathbf{Z}^{(i)})^{-1} - \mathbf{M}_{[F_i, F_i]}^{(i-1)} \preceq \varepsilon_i \cdot Sc\left(\mathbf{M}^{(i-1)}, C_i\right).$$

The algorithm APPLYCHAIN, described in Section 13.7, applies an SCC to solve equations in  $\mathbf{M}^{(0)}$  in the natural way and satisfies the following guarantee.

*Lemma 13.3.1.* Consider an  $\varepsilon$ -SCC for  $\mathbf{M}^{(0)}$  where  $\mathbf{M}^{(i)}$  and  $\mathbf{Z}^{(i)}$  can be applied to a vector in work  $W_{\mathbf{M}^{(i)}}$ ,  $W_{\mathbf{Z}^{(i)}}$  and depth  $D_{\mathbf{M}^{(i)}}$ ,  $D_{\mathbf{Z}^{(i)}}$  respectively.

The algorithm APPLYCHAIN  $\left( (\mathbf{M}^{(1)}, \mathbf{Z}^{(1)}), \dots, (\mathbf{M}^{(d)}, \mathbf{Z}^{(d)}); F_1, \dots, F_d; \mathbf{b} \right)$  corresponds to a linear operator  $\mathbf{W}$  acting on  $\mathbf{b}$  such that

1.  $\mathbf{W}^{-1} \approx_{\sum_{i=1}^d 2\varepsilon_i} \mathbf{M}^{(0)}$ , and
2. for any vector  $\mathbf{b}$ , APPLYCHAIN  $\left( (\mathbf{M}^{(1)}, \mathbf{Z}^{(1)}), \dots, (\mathbf{M}^{(d)}, \mathbf{Z}^{(d)}); F_1, \dots, F_d; \mathbf{b} \right)$  runs in  $O\left(\sum_{i=1}^d (D_{\mathbf{M}^{(i)}} + D_{\mathbf{Z}^{(i)}})\right)$  depth and  $O\left(\sum_{i=1}^d (W_{\mathbf{M}^{(i)}} + W_{\mathbf{Z}^{(i)}})\right)$  work.

For an  $\varepsilon$ -SCC chain for  $\mathbf{M}^{(0)}$ , we define the depth and work of the chain to be the work and depth required by APPLYCHAIN to apply the SCC.

### ■ 13.3.2 Choosing $F_i$ : $\alpha$ -bDD Matrices

We must choose the set of vertices  $F_i$  so that we can approximate the inverse of  $\mathbf{M}_{[F_i, F_i]}^{(i)}$  by an operator  $\mathbf{Z}^{(i)}$  that is efficiently computable. We do this by requiring that the matrix  $\mathbf{M}_{[F_i, F_i]}^{(i)}$  be  $\alpha$ -block diagonally dominant ( $\alpha$ -bDD), a term that we now define.

**Definition 13.3.3.** A Hermitian block-matrix  $\mathbf{M}$  is  $\alpha$ -bDD if

$$\forall i, \quad \mathbf{M}_{[i, i]} \succcurlyeq (1 + \alpha) \mathbf{I}_r \cdot \sum_{j: j \neq i} \|\mathbf{M}_{[i, j]}\|. \quad (13.2)$$

We remark that a 0-bDD matrix is simply a bDD matrix. In particular, for  $r = 1$ , Laplacian matrices are 0-bDD.

By picking a subset of rows at random and discarding those that violate condition (13.2), the algorithm `BDDSUBSET` (described in Section 13.8) finds a linear sized subset  $F$  of the block-rows of a bDD matrix  $\mathbf{M}$  so that  $\mathbf{M}_{[F,F]}$  is  $\alpha$ -bDD.

*Lemma 13.3.2.* Given a bDD matrix  $\mathbf{M}$  with  $n$  block-rows, and an  $\alpha \geq 0$ , `BDDSUBSET` computes a subset  $F$  of size least  $n/(8(1+\alpha))$  such that  $\mathbf{M}_{[F,F]}$  is  $\alpha$ -bDD. It runs in  $O(m)$  expected work and  $O(\log n)$  expected depth, where  $m$  is the number of nonzero blocks in  $\mathbf{M}$ .

We can express an  $\alpha$ -bDD matrix as a sum of a block diagonal matrix and a bDD matrix so that it is well-approximated by the diagonal.

*Lemma 13.3.3.* Every  $\alpha$ -bDD matrix  $\mathbf{M}$  can be written in the form  $\mathbf{X} + \mathbf{L}$  where  $\mathbf{X}$  is block-diagonal,  $\mathbf{L}$  is bDD, and  $\mathbf{X} \succcurlyeq \frac{\alpha}{2}\mathbf{L}$ .

As  $\mathbf{L}$  is positive semidefinite,  $(1 + 2/\alpha)\mathbf{X} \succcurlyeq \mathbf{M} \succcurlyeq \mathbf{X}$ , which means that  $\mathbf{X}$  is a good approximation of  $\mathbf{M}$  when  $\alpha$  is reasonably big. As block-diagonal matrices like  $\mathbf{X}$  are easy to invert, systems in these well-conditioned matrices can be solved rapidly using preconditioned iterative methods. In Section 13.9, we show that a variant of Jacobi iteration provides an operator that satisfies the requirements of Definition 13.3.2.

**Theorem 13.3.1.** *Let  $\mathbf{M}$  be a bDD matrix with index set  $V$ , and let  $F \subseteq V$  such that  $\mathbf{M}_{[F,F]}$  is  $\alpha$ -bDD for some  $\alpha \geq 4$ , and has  $m_{FF}$  nonzero blocks. The algorithm `JACOBI`( $\varepsilon, \mathbf{M}_{[F,F]}, \mathbf{b}$ ) acts as a linear operator  $\mathbf{Z}$  on  $\mathbf{b}$  that satisfies*

$$0 \preceq (\mathbf{Z})^{-1} - \mathbf{M}_{[F,F]} \preceq \varepsilon \cdot \text{Sc}(\mathbf{M}, V \setminus F).$$

*The algorithm takes  $O(m_{FF} \log(\frac{1}{\varepsilon}))$  work and  $O(\log n \log(\frac{1}{\varepsilon}))$  depth.*

We now explain how Theorems 13.3.1 and 13.3.1 allow us to construct Schur complement chains that can be applied in nearly linear time. We optimize the construction in the next section.

Theorem 13.3.1 tells us that there is a bDD matrix  $\mathbf{M}^{(0)}$  with  $O(n/\varepsilon^2)$  nonzero blocks that  $\varepsilon$ -approximates  $\mathbf{M}$ , and that for every  $i$  there is a bDD matrix  $\mathbf{M}^{(i+1)}$  with  $O(|C_i|/\varepsilon^2)$  nonzero blocks that is an  $\varepsilon$ -approximation of  $\text{Sc}(\mathbf{M}^{(i)}, F_i)$ . We will pick  $\varepsilon$  later. Lemma 13.3.2 provides a set  $F_i$  containing a constant fraction of the block-rows of  $\mathbf{M}^{(i)}$  so that  $\mathbf{M}_{[F_i, F_i]}^{(i)}$  is 4-bDD. Theorem 13.3.1 then provides an operator that solves systems in  $\mathbf{M}_{[F_i, F_i]}^{(i)}$  to  $\varepsilon$  accuracy in time  $O(\log \varepsilon^{-1})$  times the number of nonzero entries in  $\mathbf{M}_{[F_i, F_i]}^{(i)}$ . This is at most the number of nonzero entries in  $\mathbf{M}^{(i)}$ , and thus at most  $O(|C_i|/\varepsilon^2)$ . As each  $F_i$  contains at least a constant fraction of the rows of  $\mathbf{M}^{(i)}$ , the depth of the recursion,  $d$ , is  $O(\log n)$ . Thus, we can obtain constant accuracy by setting  $\varepsilon = \Theta(1/\log n)$ . The time required to apply the resulting SCC would thus be  $O(n \log^2 n \log \log n)$ .

We can reduce the running time by setting  $\varepsilon_1$  to a constant and allowing it to shrink as  $i$  increases. For example, setting  $\varepsilon_i = 1/2(i+1)^2$  results in a linear time algorithm that produces a constant-factor approximation of the inverse of  $\mathbf{M}$ . We refine this idea in the next section.

### ■ 13.3.3 Linear Sized Approximate Inverses

In this section we sketch the proof of Theorem 13.1.3, which tells us that every bDD matrix has a linear-sized approximate inverse. The rest of the details appear in Section 13.10.

The linear-sized approximate inverse of a bDD matrix  $\mathbf{M}$  with  $n$  block rows and columns is provided by a 3/4-approximate of the form  $\mathbf{U}^T \mathbf{D} \mathbf{U}$  where  $\mathbf{D}$  is block diagonal and  $\mathbf{U}$  is block upper-triangular

and has  $O(n)$  nonzero blocks. As systems of linear equations in block-triangular matrices like  $\mathbf{U}$  and  $\mathbf{U}^T$  can be solved in time proportional to their number of nonzero blocks, this provides a linear time algorithm for computing a 3/4 approximation of the inverse of  $\mathbf{M}$ . By iteratively refining the solutions provided by the approximate inverse, this allows one to find  $\varepsilon$ -accurate solutions to systems of equations in  $\mathbf{M}$  in time  $O(m \log \varepsilon^{-1})$ .

The matrix  $\mathbf{U}$  that we construct has meta-block structure that allows solves in  $\mathbf{U}$  and  $\mathbf{U}^T$  to be performed with linear work and depth  $O(\log^2 n)$ . This results in a parallel algorithm for solving equations in  $\mathbf{M}$  in work  $O(m \log \varepsilon^{-1})$  and depth  $O(\log^2 n \log \varepsilon^{-1})$ . We remark that with some additional work (analogous to Section 7 of [160]), one could reduce this depth to  $O(\log n \log \log n \log \varepsilon^{-1})$ .

The key to constructing  $\mathbf{U}$  is realizing that the algorithm JACOBI corresponds to multiplying by the matrix  $\mathbf{Z}^{(k)}$  defined in equation (13.7). Moreover, this matrix is a polynomial in  $\mathbf{X}$  and  $\mathbf{L}$  of degree  $\log(3/\varepsilon)$ , where  $\mathbf{M}_{[F_i, F_i]}^{(i)} = \mathbf{X} + \mathbf{L}$  where  $\mathbf{L}$  is bDD, and  $\mathbf{X}$  block diagonal such that  $\mathbf{X} \succcurlyeq 2\mathbf{L}$ . To force  $\mathbf{Z}^k$  to be a sparse matrix, we require  $\mathbf{L}$  be sparse.

If we use the algorithm BSDDSUBSET to choose  $F_i$ , then  $\mathbf{L}$  need not be sparse. However, this problem is easily remedied by forbidding algorithm BSDDSUBSET from choosing any vertex of more than twice average degree. Thus, we can ensure that all  $\mathbf{Z}^{(i)}$  are sparse.

**Lemma 13.3.4.** *For every bDD matrix  $\mathbf{M}$  and every  $\alpha \geq 0$ , there is a subset  $F$  of size at least  $\frac{n}{16(1+\alpha)}$  such that  $\mathbf{M}_{[F, F]}$  is  $\alpha$ -bDD and the number of nonzeros blocks in each block-row of  $F$  at most twice the average number of nonzero blocks in a block-row of  $\mathbf{M}$ .*

*Proof.* Discard every block-row of  $\mathbf{M}$  that has more than twice the average number of nonzeros blocks per row-block. Then remove the corresponding row blocks. The remaining matrix has dimension at least  $n/2$ . We can now use Lemma 13.3.2 to find an  $\alpha$ -bDD submatrix.  $\square$

We obtain the  $\mathbf{U}^T \mathbf{D} \mathbf{U}$  factorization by applying the inverse of the factorization (13.1):

$$\mathbf{M} = \begin{bmatrix} \mathbf{I} & 0 \\ \mathbf{M}_{[F, F]}^{-1} \mathbf{M}_{[F, C]} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{M}_{[F, F]} & 0 \\ 0 & \text{Sc}(\mathbf{M}, F) \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{M}_{[C, F]} \mathbf{M}_{[F, F]}^{-1} \\ 0 & \mathbf{I} \end{bmatrix}.$$

In the left and right triangular matrices we replace  $\mathbf{M}_{[F, F]}^{-1}$  with the polynomial we obtain from JACOBI. In the middle matrix, it suffices to approximate  $\mathbf{M}_{[F, F]}$  by a block diagonal matrix, and  $\text{Sc}(\mathbf{M}, F)$  by a factorization of its sparse approximation given by Theorem 13.3.1. The details, along with a careful setting of  $\varepsilon_i$ , are carried out in Section 13.10.

### ■ 13.3.4 Spectral Vertex Sparsification Algorithm

In this section, we outline a procedure APPROXSCHUR that efficiently approximates the Schur complement of a bDD matrix  $\mathbf{M}$  w.r.t. a set of indices  $F$  s.t.  $\mathbf{M}_{[F, F]}$  is  $\alpha$ -bDD. The following lemma summarizes the guarantees of APPROXSCHUR.

*Lemma 13.3.4.* Let  $\mathbf{M}$  be a bDD matrix with index set  $V$ , and  $m$  nonzero blocks. Let  $F \subseteq V$  be such that  $\mathbf{M}_{[F, F]}$  is  $\alpha$ -bDD for some  $\alpha \geq 4$ . The algorithm APPROXSCHUR( $\mathbf{M}, F, \varepsilon$ ), returns a matrix  $\widetilde{\mathbf{M}}_{SC}$  s.t.

1.  $\widetilde{\mathbf{M}}_{SC}$  has  $O(m(\varepsilon^{-1} \log \log \varepsilon^{-1})^{O(\log \log \varepsilon^{-1})})$  nonzero blocks, and
2.  $\widetilde{\mathbf{M}}_{SC} \approx_{\varepsilon} \text{Sc}(\mathbf{M}, F)$ ,

in  $O(m(\varepsilon^{-1} \log \log \varepsilon^{-1})^{O(\log \log \varepsilon^{-1})})$  work and  $O(\log n(\log \log \varepsilon^{-1}))$  depth.

We sketch a proof of the above lemma in this section. A complete proof and pseudocode for APPROXSCHUR are given in Section 13.11.

First consider a very simple special case: where  $F$  is a singleton,  $F = \{i\}$ . Let  $C = V \setminus F$  be the remaining indices.

$$Sc(\mathbf{M}, i) = \mathbf{M}_{[C,C]} - \mathbf{M}_{[C,i]} \mathbf{M}_{[i,i]}^{-1} \mathbf{M}_{[i,C]}$$

Thus, if  $\mathbf{M}_{[C,i]}$  has  $k$  nonzero blocks,  $Sc(\mathbf{M}, i)$  could have  $k^2$  additional nonzero blocks compared to  $\mathbf{M}_{[C,C]}$ , potentially making it dense. If  $\mathbf{M}$  were a graph Laplacian, then  $\mathbf{M}_{[C,i]} \mathbf{M}_{[i,i]}^{-1} \mathbf{M}_{[i,C]}$  would represent the adjacency matrix of a weighted clique. In Section 13.14, we construct weighted expanders that allow us to  $\varepsilon$ -approximate  $Sc(\mathbf{M}, i)$  in this case using  $m + O(k\varepsilon^{-4})$  edges. In Section 13.11.4, we show how to use such weighted expanders to sparsify  $Sc(\mathbf{M}, i)$  when  $\mathbf{M}$  is a bDD matrix.

This reduction can also be performed in parallel. If  $F$  is such that  $\mathbf{M}_{[F,F]}$  is block diagonal, we can approximate  $Sc(\mathbf{M}, F)$  by expressing  $\mathbf{M}_{[C,i]} \mathbf{M}_{[i,i]}^{-1} \mathbf{M}_{[i,C]}$  as  $\sum_{i \in F} \mathbf{M}_{[C,i]} \mathbf{M}_{[i,i]}^{-1} \mathbf{M}_{[i,C]}$ , and using  $|F|$  weighted expanders. However,  $\mathbf{M}_{[F,F]}$  may not be diagonal. Instead, we give a procedure SCHURSQUARE that generates  $\mathbf{M}'$  that is better approximated by its diagonal.

Invoking SCHURSQUARE a few times leads a sequence of matrices  $\mathbf{M}_{[F,F]}^{(0)}, \mathbf{M}_{[F,F]}^{(1)}, \dots, \mathbf{M}_{[F,F]}^{(i)}$ . We will show that  $\mathbf{M}_{[F,F]}^{(i)}$  is  $\varepsilon$ -approximated by its diagonal and we call the procedure LASTSTEP to approximate  $Sc(\mathbf{M}^{(i)}, F)$ . An additional caveat is that replacing  $\mathbf{M}_{[F,F]}^{(i)}$  by its diagonal at this step gives errors that are difficult to bound. We discuss the correct approximation below

SCHURSQUARE is based on a *squaring* identity for matrix inverse developed in [223]. Given a splitting of  $\mathbf{M}_{[F,F]}$  into  $\mathbf{D} - \mathbf{A}$ , where  $\mathbf{D}$  is block-diagonal, and  $\mathbf{A}$  has its diagonal blocks as zero, it relies on the fact that the matrix

$$\mathbf{M}_2 = \frac{1}{2} \begin{bmatrix} \mathbf{D} - \mathbf{A}\mathbf{D}^{-1}\mathbf{A} & \mathbf{M}_{[F,C]} + \mathbf{A}\mathbf{D}^{-1}\mathbf{M}_{[F,C]} \\ \mathbf{M}_{[C,F]} + \mathbf{M}_{[C,F]}\mathbf{D}^{-1}\mathbf{A} & 2\mathbf{M}_{[C,C]} - \mathbf{M}_{[C,F]}\mathbf{D}^{-1}\mathbf{M}_{[F,C]} \end{bmatrix} \quad (13.3)$$

satisfies  $Sc(\mathbf{M}, F) = Sc(\mathbf{M}_2, F)$ . Furthermore, we can show that if  $\mathbf{M}$  is  $\alpha$ -bDD,  $\mathbf{D} - \mathbf{A}\mathbf{D}^{-1}\mathbf{A}$  is  $\alpha^2$ -bDD, which indicates that the block on  $[F, F]$  rapidly approaches being diagonal.

As  $\mathbf{M}_2$  may be dense, we construct a sparse approximation to it. Since  $\mathbf{D}$  is diagonal, we can construct sparse approximations to the blocks  $(\mathbf{M}_2)_{[F,F]}$  and  $(\mathbf{M}_2)_{[C,C]}$  in a manner analogous to the case of diagonal  $\mathbf{M}_{[F,F]}$ . Similarly, we use bipartite expanders to construct sparse approximations to  $(\mathbf{M}_2)_{[C,F]}$  and  $(\mathbf{M}_2)_{[F,C]}$ .

Our sequence of calls to SCHURSQUARE terminates with  $\mathbf{M}_{[F,F]}^{(i)}$  being roughly  $\varepsilon^{-1}$ -bDD. We then return LASTSTEP( $\mathbf{M}^{(i)}, F, \varepsilon^{-1}, \varepsilon$ ). As mentioned above, we cannot just replace  $\mathbf{M}_{[F,F]}^{(i)}$  by its diagonal. Instead, LASTSTEP performs one step of *squaring* similar to SCHURSQUARE with a key difference: Rather than expressing  $\mathbf{M}_{[F,F]}^{(i)}$  as  $\mathbf{D} - \mathbf{A}$ , it expresses it as  $\mathbf{X} + \mathbf{L}$ , where  $\mathbf{X}$  is block-diagonal, and  $\mathbf{L}$  is just barely bDD. With this splitting, it constructs  $\mathbf{M}^{(last)}$  after performing one iteration similar to Eq. (13.3). After this step, it replaces the  $\mathbf{M}_{[F,F]}^{(last)}$  block with the block-diagonal matrix  $\mathbf{X}$ . Again, we directly produce sparse approximations to  $\mathbf{M}^{(last)}$  and its Schur complement via weighted (bipartite) expanders. A precise description and proofs are given in Section 13.11.2.

### ■ 13.3.5 Sparsifying bDD Matrices

The main technical hurdle left to address is how we sparsify bDD matrices. We to do this both to approximate  $\mathbf{M}$  by  $\mathbf{M}^{(0)}$ , if  $\mathbf{M}$  is not already sparse, and to ensure that all the matrices  $\mathbf{M}^{(i)}$  remain sparse. While the spectral vertex sparsification algorithm described in the previous section allows us

to compute an approximation to a Schur complement  $Sc(\mathbf{M}^{(i)}, F_i)$ , it is sparse only when  $\mathbf{M}^{(i)}$  is already sparse. As we iteratively apply this procedure, the density of the matrices produced will grow unacceptably. We overcome this problem by occasionally applying another sparsification routine that substantially decreases the number of nonzero blocks. The cost of this sparsification routine is that it requires solving systems of equations in sparse  $bDD$  matrices. We, of course, do this recursively.

Our sparsification procedure begins by generalizing the observation that graph Laplacians can be sparsified by sampling edges with probabilities determined by their effective resistances [242]. There is a block analog of leverage scores (13.36) that provides probabilities of sampling blocks so that the resulting sampled matrix approximates the original and has  $O(n \log n)$  nonzero blocks with high probability. To compute this block analog of leverage scores we employ the procedure developed in Chapter 3.

Once we generalize their results to block matrices, they show that we can obtain sufficiently good estimates of the block leverage scores by computing leverage scores in a  $bDD$  matrix obtained by randomly subsampling blocks of the original. The block leverage scores in this matrix are obtained by solving a logarithmic number of linear equations in this subsampled matrix. Thus, our sparsification procedure requires constructing a solver for a subsampled matrix and then applying that solver a logarithmic number of times. We compute this solver recursively.

There is a tradeoff between the number of nonzero blocks in the subsampled system and in the resulting approximation of the original matrix. If the original matrix has  $m$  nonzero blocks and we subsample to a system of  $m/K$  nonzero blocks, then we obtain an  $\varepsilon$ -approximation of the original matrix with  $O(K\varepsilon^{-2}n \log n)$  nonzero blocks.

The details of the analysis of the undersampling procedure appear in Section 13.12.

### ■ 13.3.6 Our Algorithm

We now explain how we prove Theorem 13.1.2. The details supporting this exposition appear in Section 13.13. Our main goal is to control the density of the Schur complement chain as we repeatedly invoke Lemma 13.3.4.

Starting from some  $\mathbf{M}^{(0)}$ , we compute sets  $F_i$  (via calls to `BDDSUBSET`), approximate solvers  $\mathbf{Z}^{(i)}$  (via `JACOBI`), and approximations of Schur complements  $\mathbf{M}^{(i)}$  (via `APPROXSCHUR`), until we obtain a matrix  $\mathbf{M}^{(i)}$  such that its dimension is a smaller than that of  $\mathbf{M}^{(0)}$  by a large constant factor (like 4). While the dimension of  $\mathbf{M}^{(i)}$  is much smaller, its number of nonzero blocks is potentially larger by an even larger factor. We use the procedure described in the previous section to sparsify it. This sparsification procedure produces a sparse approximation of the matrix at the cost of solving systems of equations in a subsampled version of that matrix. Some care is required to balance the cost of the resulting recursion.

We now sketch an analysis of a nearly linear time algorithm that results from a simple choice of parameters. We optimize the parameter choice and analysis in Section 13.13. Let  $n$  be the dimension of  $\mathbf{M}^{(0)}$  and let  $m$  be its number of nonzero blocks. To begin, assume that  $m \leq n\Delta \log^3 n$ , for a  $\Delta$  to be specified later. We call this the *sparse* case, and address the case of dense  $\mathbf{M}$  later.

We consider fixing  $\varepsilon_i = c/\log n$  for all  $i$ , for some constant  $c$ . As the depth of the Schur complement chain is  $O(\log n)$ , this results in a solver with constant accuracy. A constant number of iterations of the procedure described above are required to produce an  $\mathbf{M}^{(i)}$  whose dimension is a factor of 4 smaller than  $\mathbf{M}^{(0)}$ . Lemma 13.3.4 tells us that the edge density of this  $\mathbf{M}^{(i)}$  is potentially higher than that of  $\mathbf{M}^{(0)}$  by a factor of

$$O\left(\left(\varepsilon^{-1} \log \log \varepsilon^{-1}\right)^{O(\log \log \varepsilon^{-1})}\right) = \exp\left(O(\log \log^2 n)\right) = n^{o(1)}.$$

Set  $\Delta$  to be this factor. We use the sparsification procedure from the previous section to guarantee that no matrix in the chain has density higher than that of this matrix, which is upper bounded by  $(m/n)\Delta = \Delta^2 \log^3 n$ .

Setting  $K = 2\Delta$ , the subsampling produces a matrix of density half that of  $\mathbf{M}^{(0)}$ , and it produces a sparse approximation of  $\mathbf{M}^{(i)}$  of density  $O(K\varepsilon_i^{-2} \log n) \leq O(\Delta \log^3 n)$ , which, by setting constants appropriately, we can also force to be half that of  $\mathbf{M}^{(0)}$ . In order to perform the sparsification procedure, we need to construct a Schur complement chain for the subsampled matrix, and then use it to solve  $O(\log n)$  systems of linear equations. The cost of using this chain to solve equations in the subsampled system is at most  $O(n\Delta^2 \log^4 n)$ , and the cost of using the solutions to these equations to sparsify  $\mathbf{M}^{(i)}$  is  $O(m \log n)$ . The cost of the calls to `BDDSUBSET` and `APPROXSCHUR` are proportional to the number of edges in the matrices, which is  $O(n\Delta \log^4 n)$ .

We repeat this procedure all the way down the chain, only using sparsification when the dimension of  $\mathbf{M}^{(i)}$  shrinks by a factor of 4. Since none of the matrices that we generate have density higher than  $\Delta \log^3 n$ , we remain in the sparse case. Let  $T_{sparse}(n)$  be the time required to construct a solver chain on systems of size  $n$  with  $m \leq n\Delta \log^3 n$ . We obtain the following recurrence

$$T_{sparse}(n) \leq 2T_{sparse}(n/4) + n\Delta^2 \log^4 n + m \log n + m\Delta \leq 2T_{sparse}(n/4) + n\Delta^2 \log^4 n + n\Delta \log n + n\Delta^2,$$

which gives

$$T(n)_{sparse} \leq O(n\Delta^2 + n\Delta^2 \log^4 n) \leq n^{1+o(1)}.$$

To handle the case of dense  $\mathbf{M}$ , we repeatedly sparsify while keeping  $n$  fixed until we obtain a matrix with fewer than  $n\Delta \log^3 n$  edges, at which point we switch to the algorithm described above. The running time of this algorithm on a graph with  $m$  edges,  $T_{dense}(m)$ , satisfies the recurrence

$$T_{dense}(m) \leq \begin{cases} T_{sparse}(n) & \text{if } m \leq n\Delta \log^3 n, \text{ and} \\ 2T_{dense}(m/2) + n\Delta^2 \log^4 n + m \log n + m\Delta & \text{otherwise.} \end{cases}$$

Thus  $T_{dense}(m)$  is upper bounded by  $O(mn^{o(1)} + n^{1+o(1)})$ .

We tighten this bound in Section 13.13 to prove Theorem 13.1.2 by carefully choosing the parameters to accompany a sequence  $\varepsilon$  that starts constant and decreases slowly.

## ■ 13.4 Summary

We introduce a new approach to solving systems of linear equations that gives the first nearly linear time algorithms for solving systems proof that connection Laplacians have linear-sized approximate inverses. This was unknown even for graph Laplacians.

Our algorithms build on ideas introduced in [223] and are a break from those used in the previous work on solving systems of equations in graph Laplacians [256, 243, 147, 146, 141, 57]. Those algorithms all rest on *support theory* [34], originally introduced by Vaidya [256], and rely on the fact that the Laplacian of one edge is approximated by the Laplacian of a path between its endpoints. No analogous fact is true for connection Laplacians, even those with complex entries for  $r = 1$ .

Instead, our algorithms rely on many new ideas, the first being that of sparsifying the matrices that appear during elimination. Other critical ideas are finding  $\alpha$ -bDD subsets of vertices to eliminate in bulk, approximating Schur complements without computing them explicitly, and the use of sub-sampling to sparsify in a recursive fashion. To efficiently compute approximations of the Schur complements, we introduce a new operation that transforms a matrix into one with the same Schur complement but a much better conditioned upper block (13.3). To obtain the sharp bounds in Theorem 13.1.2, we exploit a new linear-time algorithm for constructing linear-sized sparse approximations to



implicitly represented weighted cliques whose edge weights are products of weights at vertices (Section 13.14), and extend this to the analog for bDD matrices (Section 13.11.4).

## ■ 13.5 Background

*Proof.* (of Fact 13.2.2) Note that

$$\|\mathbf{x} - \mathbf{Z}\mathbf{b}\|_M^2 = \mathbf{x}^* \mathbf{M} \mathbf{x} - 2\mathbf{x}^* \mathbf{M} \mathbf{Z} \mathbf{M} \mathbf{x} + \mathbf{x}^* \mathbf{M} \mathbf{Z} \mathbf{M} \mathbf{Z} \mathbf{M} \mathbf{x}$$

Since all nonzero eigenvalues of  $\mathbf{M}$  are at least  $\mu$ , the eigenvalues of  $\mathbf{M}^{1/2}(\mathbf{M} + \varepsilon\mu\mathbf{I})^{-1}\mathbf{M}^{1/2}$  lie between  $1/(1 + \varepsilon)$  and 1. Using  $\mathbf{Z} \approx_\varepsilon (\mathbf{M} + \varepsilon\mu\mathbf{I})^{-1}$ , we see that the eigenvalues of  $\mathbf{M}^{1/2}\mathbf{Z}\mathbf{M}^{1/2}$  lie between  $e^{-2\varepsilon}$  and  $e^\varepsilon$ . Using  $0 < \varepsilon < 1/2$ , we have

$$\|\mathbf{x} - \mathbf{Z}\mathbf{b}\|_M^2 \leq (1 - 2e^{-2\varepsilon} + e^{2\varepsilon})\mathbf{x}^* \mathbf{M} \mathbf{x} \leq 6\varepsilon\|\mathbf{x}\|_M^2.$$

□

**Fact 13.5.1.** Let  $\mathbf{A}$  be a matrix of condition number  $\kappa$  and let  $\mathbf{A} \approx_\varepsilon \mathbf{B}$  for  $\varepsilon \leq (56\kappa^3)^{-1}$ . Then,  $\mathbf{A}^3 \approx_{28\kappa^3\varepsilon} \mathbf{B}^3$ .

*Proof.* First, observe that  $\mathbf{A} \approx_\varepsilon \mathbf{B}$  implies that  $\|\mathbf{B}\| \leq (1 + e^\varepsilon)\|\mathbf{A}\| \leq 2\|\mathbf{A}\|$ . It also implies that  $\|\mathbf{A} - \mathbf{B}\| \leq 2\varepsilon\|\mathbf{A}\|$ . As

$$\mathbf{A}^2 - \mathbf{B}^2 = \frac{1}{2}(\mathbf{A} - \mathbf{B})(\mathbf{A} + \mathbf{B}) + \frac{1}{2}(\mathbf{A} + \mathbf{B})(\mathbf{A} - \mathbf{B}),$$

$\|\mathbf{A}^2 - \mathbf{B}^2\| \leq 6\varepsilon\|\mathbf{A}\|^2$ . Similarly, as

$$\mathbf{A}^3 - \mathbf{B}^3 = \frac{1}{2}(\mathbf{A} - \mathbf{B})(\mathbf{A}^2 + \mathbf{B}^2) + \frac{1}{2}(\mathbf{A} + \mathbf{B})(\mathbf{A}^2 - \mathbf{B}^2),$$

$\|\mathbf{A}^3 - \mathbf{B}^3\| \leq 28\varepsilon\|\mathbf{A}\|^3$ .

Let  $\kappa = \|\mathbf{A}\|/\lambda_{\min}(\mathbf{A})$ . The above relation implies that

$$\|\mathbf{A}^3 - \mathbf{B}^3\| \leq 28\varepsilon\kappa^3\lambda_{\min}(\mathbf{A})^3 = 28\varepsilon\kappa^3\lambda_{\min}(\mathbf{A}^3).$$

Thus, as  $28\varepsilon\kappa^3 \leq 1/2$ ,

$$\mathbf{A}^3 \approx_{56\varepsilon\kappa^3} \mathbf{B}^3.$$

□

*Proof.* (of Fact 13.2.3)

By Fact 13.5.1, using  $\delta = \frac{\varepsilon}{56\kappa^3}$ ,

$$\mathbf{M} \mathbf{Z}^3 \mathbf{M} \approx_\varepsilon \mathbf{M}(\mathbf{M} + \varepsilon\mu\mathbf{I})^{-3} \mathbf{M}.$$

$\mathbf{M}$  has an eigendecomposition in the same basis as  $(\mathbf{M} + \varepsilon\mu\mathbf{I})^{-3}$ , and so it follows that  $\mathbf{M}(\mathbf{M} + \varepsilon\mu\mathbf{I})^{-3}\mathbf{M}$  has the same eigenbasis, and the same null space as  $\mathbf{M}$ .

When  $\lambda^{-1}$  is the eigenvalue of  $\mathbf{M}^+$  of an eigenvector  $\mathbf{v}$ , the corresponding eigenvalue of  $\mathbf{M}(\mathbf{M} + \varepsilon\mu\mathbf{I})^{-3}\mathbf{M}$  is  $\beta \stackrel{\text{def}}{=} \frac{\lambda^2}{(\lambda + \varepsilon\mu)^3}$  and

$$\lambda^{-1} > \beta \geq \frac{\lambda^2}{(1 + \varepsilon)^3 \lambda^3} = e^{-3\varepsilon} \lambda^{-1}.$$

So  $\mathbf{M}(\mathbf{M} + \varepsilon\mu\mathbf{I})^{-3}\mathbf{M} \approx_{3\varepsilon} \mathbf{M}^+$ , and  $\mathbf{M} \mathbf{Z}^3 \mathbf{M} \approx_{4\varepsilon} \mathbf{M}^+$ . □

**Fact 13.5.2.** For every  $\mathbf{d} \in \mathbb{C}^{r \times r}$ , we can find  $\mathbf{Q}^{(1)}, \mathbf{Q}^{(2)} \in \mathbb{C}^{r \times r}$  where

$$(\mathbf{Q}^{(1)})(\mathbf{Q}^{(1)})^* = (\mathbf{Q}^{(1)})^*(\mathbf{Q}^{(1)}) = (\mathbf{Q}^{(2)})(\mathbf{Q}^{(2)})^* = (\mathbf{Q}^{(2)})^*(\mathbf{Q}^{(2)}) = \mathbf{I}_r,$$

such that

$$\mathbf{d} = \frac{1}{2} \|\mathbf{d}\| (\mathbf{Q}^{(1)} + \mathbf{Q}^{(2)}).$$

*Proof.* Let  $\hat{\mathbf{d}} = \frac{1}{\|\mathbf{d}\|} \mathbf{d}$ . Thus,  $\|\hat{\mathbf{d}}\| = 1$ . Write  $\hat{\mathbf{d}}$  using its singular value decomposition as  $\mathbf{U} \mathbf{D} \mathbf{V}^*$ , where  $\mathbf{D}, \mathbf{U}, \mathbf{V} \in \mathbb{C}^{r \times r}$ ,  $\mathbf{D}$  is a real diagonal matrix with the singular values of  $\mathbf{d}$  on the diagonal, and  $\mathbf{U} \mathbf{U}^* = \mathbf{U}^* \mathbf{U} = \mathbf{V} \mathbf{V}^* = \mathbf{V}^* \mathbf{V} = \mathbf{I}_r$ . Since  $\|\hat{\mathbf{d}}\| = 1$ , we have  $\mathbf{D}_{j,j} \in [0, 1]$  for all  $j \in [r]$ . Thus, there exists a real  $\theta_j$  such that  $\cos \theta_j = \mathbf{D}_{j,j}$ . If we let  $\mathbf{D}^{(1)}, \mathbf{D}^{(2)}$  be diagonal matrices such that  $\mathbf{D}_{j,j}^{(1)} = \exp(i\theta_j), \mathbf{D}_{j,j}^{(2)} = \exp(-i\theta_j)$ , we have  $\mathbf{D} = \frac{1}{2}(\mathbf{D}^{(1)} + \mathbf{D}^{(2)})$ . Moreover,

$$(\mathbf{D}^{(1)})(\mathbf{D}^{(1)})^* = (\mathbf{D}^{(1)})^*(\mathbf{D}^{(1)}) = (\mathbf{D}^{(2)})(\mathbf{D}^{(2)})^* = (\mathbf{D}^{(2)})^*(\mathbf{D}^{(2)}) = \mathbf{I}_r.$$

Letting  $\mathbf{Q}^{(k)} = \mathbf{U} \mathbf{D}^{(k)} \mathbf{V}^*$  for  $k = 1, 2$ , we get  $\hat{\mathbf{d}} = \frac{1}{2}(\mathbf{Q}^{(1)} + \mathbf{Q}^{(2)})$  and hence  $\mathbf{d} = \frac{1}{2} \|\mathbf{d}\| (\mathbf{Q}^{(1)} + \mathbf{Q}^{(2)})$ . Moreover,

$$(\mathbf{Q}^{(1)})(\mathbf{Q}^{(1)})^* = (\mathbf{Q}^{(1)})^*(\mathbf{Q}^{(1)}) = (\mathbf{Q}^{(2)})(\mathbf{Q}^{(2)})^* = (\mathbf{Q}^{(2)})^*(\mathbf{Q}^{(2)}) = \mathbf{I}_r. \quad \square$$

**Fact 13.5.3** (Lemma B.1. from [194]). If  $\mathbf{M}$  and  $\widetilde{\mathbf{M}}$  are positive semidefinite matrices satisfying  $\mathbf{M} \preceq \widetilde{\mathbf{M}}$ , then

$$Sc(\mathbf{M}, F) \preceq Sc(\widetilde{\mathbf{M}}, F).$$

This fact can be proven via an energy minimization definition of Schur complement. More details on this formulation can be found in [194].

## ■ 13.6 Block Diagonally Dominant Matrices

In this section, we prove a few basic facts about bDD matrices. The following lemma gives an equivalent definition of bDD matrices.

**Lemma 13.6.1.** A Hermitian block-matrix  $\mathbf{M} \in (\mathbb{C}^{r \times r})^{n \times n}$  is bDD iff it can be written as  $\mathbf{D} - \mathbf{A}$  where  $\mathbf{D}$  is block-diagonal,  $\mathbf{A}$  is Hermitian, and  $\mathbf{D}_{[i,i]} \succeq \mathbf{I}_r \sum_j \|\mathbf{A}_{[i,j]}\|$ , for all  $i \in V$ .

*Proof.* The only if direction is easy. For bDD matrix  $\mathbf{M}$ , if we let  $\mathbf{D}$  be the block diagonal matrix such that  $\mathbf{D}_{[i,i]} = \mathbf{M}_{[i,i]}$ , and  $\mathbf{A}$  be the matrix such that

$$\mathbf{A}_{[i,j]} = \begin{cases} 0 & i = j \\ \mathbf{M}_{[i,j]} & i \neq j \end{cases}$$

It is immediate that  $\mathbf{A}$  is Hermitian. Moreover,  $\mathbf{M}$  is bDD implies that  $\mathbf{D}_{[i,i]} \succeq \mathbf{I}_r \sum_j \|\mathbf{A}_{[i,j]}\|$ .

For the if direction. Suppose we have  $\mathbf{D}, \mathbf{A}$  such that  $\mathbf{D}$  is block-diagonal, and for all  $i$   $\mathbf{D}_{[i,i]} \succeq \mathbf{I}_r \sum_j \|\mathbf{A}_{[i,j]}\|$ . Thus, letting  $\mathbf{M} = \mathbf{D} - \mathbf{A}$ , we have, for all  $i \in V$ ,

$$\mathbf{I}_r \cdot \sum_{j \neq i} \|\mathbf{M}_{[i,j]}\| = \mathbf{I}_r \cdot \sum_{j \neq i} \|\mathbf{A}_{[i,j]}\| \preceq \mathbf{D}_{[i,i]} - \mathbf{I}_r \cdot \|\mathbf{A}_{[i,i]}\| \preceq \mathbf{D}_{[i,i]} - \mathbf{A}_{[i,i]} = \mathbf{M}_{[i,i]},$$

where we used that since  $\mathbf{A}, \mathbf{M}$  are Hermitian,  $\mathbf{D}$  is also Hermitian, and thus  $\mathbf{D}_{[i,i]}, \mathbf{A}_{[i,i]}$  are Hermitian.  $\square$

This immediately implies the corollary that flipping the sign of off-diagonal blocks preserves bDD-ness.

**Corollary 13.6.2.** *Given a bDD matrix  $\mathbf{M}$ , write it as  $\mathbf{D} - \mathbf{A}$ , where  $\mathbf{D}$  is a block-diagonal, and  $\mathbf{A}$  has its diagonal blocks as zero. Then,  $\mathbf{D} + \mathbf{A}$  is also PSD.*

*Proof.* First observe that for all  $i$ ,  $(\mathbf{D} + \mathbf{A})_{[i,i]} = (\mathbf{D} - \mathbf{A})_{[i,i]} = \mathbf{M}_{[i,i]}$ , i.e., their diagonal blocks are identical. Moreover, for all  $i \neq j$ , we have  $\|(\mathbf{D} + \mathbf{A})_{[i,j]}\| = \|(\mathbf{D} - \mathbf{A})_{[i,j]}\| = \|\mathbf{M}_{[i,j]}\|$ . Thus  $\mathbf{D} + \mathbf{A}$  is also bDD, and hence PSD.  $\square$

Next, we show that the class of bDD matrices is closed under Schur complement.

**Lemma 13.6.3.** *The class of bDD matrices is closed under Schur complement.*

*Proof.* Since Schur complementation does not depend on the order of indices eliminated, it suffices to prove that for any bDD matrix  $\mathbf{M} \in (\mathbb{C}^{r \times r})^{n \times n}$ ,  $Sc(\mathbf{M}, 1)$  is a bDD matrix. Let  $C = V \setminus \{1\}$ .

We have  $Sc(\mathbf{M}, 1) = \mathbf{M}_{[C,C]} - \mathbf{M}_{[C,1]} \mathbf{M}_{[1,1]}^{-1} \mathbf{M}_{[1,C]}$ . Let  $\mathbf{D} \in (\mathbb{C}^{r \times r})^{C \times C}$  be the block diagonal matrix such that  $\mathbf{D}_{[i,i]} = \mathbf{M}_{[i,i]}$  for  $i \in C$ . Expressing  $Sc(\mathbf{M}, 1)$  as  $\mathbf{D} + (\mathbf{M}_{[C,C]} - \mathbf{D} + \mathbf{M}_{[C,1]} \mathbf{M}_{[1,1]}^{-1} \mathbf{M}_{[1,C]})$ , we have for any  $i \in C$ ,

$$\begin{aligned} \sum_{j \in C} \|(\mathbf{M}_{[C,C]} - \mathbf{D} + \mathbf{M}_{[C,1]} \mathbf{M}_{[1,1]}^{-1} \mathbf{M}_{[1,C]})_{[i,j]}\| &\leq \left( \sum_{j \in C: j \neq i} \|\mathbf{M}_{[i,j]}\| \right) + \sum_{j \in C} \|\mathbf{M}_{[i,1]} \mathbf{M}_{[1,1]}^{-1} \mathbf{M}_{[1,j]}\| \\ &\leq \left( \sum_{j \in C: j \neq i} \|\mathbf{M}_{[i,j]}\| \right) + \|\mathbf{M}_{[i,1]}\| \|\mathbf{M}_{[1,1]}^{-1}\| \sum_{j \in C} \|\mathbf{M}_{[1,j]}\| \\ &\leq \left( \sum_{j \in C: j \neq i} \|\mathbf{M}_{[i,j]}\| \right) + \|\mathbf{M}_{[i,1]}\| = \sum_{j \in V: j \neq i} \|\mathbf{M}_{[i,j]}\|, \end{aligned}$$

where the last inequality uses  $\|\mathbf{M}_{[1,1]}^{-1}\| \left( \sum_{j \neq 1} \|\mathbf{M}_{[1,j]}\| \right) \leq 1$ , since  $\mathbf{M}_{[1,1]} \succcurlyeq \mathbf{I}_r \cdot \sum_{j \neq 1} \|\mathbf{M}_{[1,j]}\|$ . Thus using Lemma 13.6.1, we have  $Sc(\mathbf{M}, 1) = \mathbf{D} - (-\mathbf{M}_{[C,C]} + \mathbf{D} - \mathbf{M}_{[C,1]} \mathbf{M}_{[1,1]}^{-1} \mathbf{M}_{[1,C]})$  is bDD.  $\square$

The next definition describes a special form that we can express any bDD matrix in, which will occasionally be useful.

**Definition 13.6.4.** A matrix  $\mathbf{B} \in (\mathbb{C}^{r \times r})^{n \times m}$  is called a *unitary edge-vertex transfer matrix*, when each block column of  $\mathbf{B}_{[e]}$  has exactly two nonzero blocks  $\mathbf{U}_e, \mathbf{V}_e \in \mathbb{C}^{r \times r}$  s.t.  $\mathbf{U}_e \mathbf{U}_e^* = \mathbf{V}_e \mathbf{V}_e^* = w_e \mathbf{I}_r$  where  $w_e \geq 0$ .

**Lemma 13.6.5.** *Every bDD matrix  $\mathbf{M} \in (\mathbb{C}^{r \times r})^{n \times n}$  with  $m$  nonzero off-diagonal blocks can be written as  $\mathbf{X} + \mathbf{B}\mathbf{B}^*$  where  $\mathbf{B} \in (\mathbb{C}^{r \times r})^{n \times 2m}$  is a unitary edge-vertex transfer matrix and  $\mathbf{X}$  is a block diagonal PSD matrix. This implies that every bDD matrix is PSD. Furthermore, for every block diagonal matrix  $\mathbf{Y}$  s.t.  $\mathbf{M} - \mathbf{Y}$  is bDD, we have  $\mathbf{X} \succcurlyeq \mathbf{Y}$ . This decomposition can be found in  $O(m)$  time and  $O(\log n)$  depth.*

*Proof.* Consider a pair  $\{i, j\} \in V \times V$  such that  $i \neq j$ , and  $\mathbf{M}_{[i,j]} \neq 0$ . Using Fact 13.5.2, we can write such a  $\mathbf{M}_{[i,j]}$  as  $\frac{1}{2} \|\mathbf{M}_{[i,j]}\| (\mathbf{Q}_{\{i,j\}}^{(1)} + \mathbf{Q}_{\{i,j\}}^{(2)})$ , where  $\mathbf{Q}_{\{i,j\}}^{(1)} (\mathbf{Q}_{\{i,j\}}^{(1)})^* = \mathbf{Q}_{\{i,j\}}^{(2)} (\mathbf{Q}_{\{i,j\}}^{(2)})^* = \mathbf{I}_r$ . we construct two vectors vector  $\mathbf{B}_{\{i,j\}}^{(1)}, \mathbf{B}_{\{i,j\}}^{(2)} \in (\mathbb{C}^{r \times r})^n$  such that for  $k = 1, 2$ ,  $(\mathbf{B}_{\{i,j\}}^{(k)})_{[i]} = \frac{1}{\sqrt{2}} \|\mathbf{M}_{[i,j]}\|^{1/2} \mathbf{I}_r$ ,

$(\mathbf{B}_{\{i,j\}}^{(k)})_{[j]} = \frac{1}{\sqrt{2}} \|\mathbf{M}_{[i,j]}\|^{1/2} (\mathbf{Q}_{\{i,j\}}^{(k)})^*$ , and all other blocks are zero. We can verify that for all  $k, \ell \in V$ ,

$$\left( \mathbf{B}_{\{i,j\}}^{(1)} (\mathbf{B}_{\{i,j\}}^{(1)})^* + \mathbf{B}_{\{i,j\}}^{(2)} (\mathbf{B}_{\{i,j\}}^{(2)})^* \right)_{[k,\ell]} = \begin{cases} \|\mathbf{M}_{[i,j]}\| \mathbf{I}_r & k = \ell = i, \\ \|\mathbf{M}_{[i,j]}\| \mathbf{I}_r & k = \ell = j, \\ \mathbf{M}_{[i,j]} & k = i, \ell = j, \\ \mathbf{M}_{[i,j]}^* = \mathbf{M}_{[j,i]} & k = j, \ell = i, \\ 0 & \text{otherwise.} \end{cases}$$

We let  $\mathbf{X} = \mathbf{M} - \sum_{\{i,j\}:i \neq j} \left( \mathbf{B}_{\{i,j\}}^{(1)} (\mathbf{B}_{\{i,j\}}^{(1)})^* + \mathbf{B}_{\{i,j\}}^{(2)} (\mathbf{B}_{\{i,j\}}^{(2)})^* \right)$ , which must be block-diagonal. We now show that for all  $i \in V$ , the block  $\mathbf{X}_{[i,i]}$  is PSD. We have for all  $i \in V$ ,

$$\mathbf{X}_{[i,i]} = \mathbf{M}_{[i,i]} - \mathbf{I}_r \cdot \sum_{j:j \neq i} \|\mathbf{M}_{[i,j]}\| \succcurlyeq 0,$$

where the last inequality holds since  $\mathbf{M}$  is bDD.

Thus, if we define  $\mathbf{B} \in (\mathbb{C}^{r \times r})^{n \times 2m}$  such that its columns are all the vectors  $\mathbf{B}_{\{i,j\}}^{(1)}, \mathbf{B}_{\{i,j\}}^{(2)}$  defined above, we have  $\mathbf{M} = \mathbf{X} + \mathbf{B}\mathbf{B}^*$ , and every column of  $\mathbf{B}$  has exactly 2 nonzero blocks.

To show that for every block diagonal  $\mathbf{Y}$  s.t.  $\mathbf{M} - \mathbf{Y}$  is bDD,  $\mathbf{X} \succcurlyeq \mathbf{Y}$ , first consider applying the decomposition described above to  $\mathbf{M} - \mathbf{Y}$  instead of  $\mathbf{M}$ . Since the construction of  $\mathbf{B}$  only depends on the off-diagonal blocks, we get  $\mathbf{M} - \mathbf{Y} = \mathbf{\Lambda} + \mathbf{B}\mathbf{B}^*$ , where  $\mathbf{\Lambda}$  is block diagonal and PSD. So,  $\mathbf{X} - \mathbf{Y} = \mathbf{M} - \mathbf{B}\mathbf{B}^* - \mathbf{Y} = \mathbf{\Lambda} \succcurlyeq 0$ .

It is immediate that the decomposition can be found in  $O(m)$  time and  $O(\log n)$  depth.  $\square$

## ■ 13.7 Schur Complement Chains

In this section, we give a proof of Lemma 13.3.1. We restate the lemma here for convenience.

*Lemma 13.3.1.* Consider an  $\varepsilon$ -SCC for  $\mathbf{M}^{(0)}$  where  $\mathbf{M}^{(i)}$  and  $\mathbf{Z}^{(i)}$  can be applied to a vector in work  $W_{\mathbf{M}^{(i)}}, W_{\mathbf{D}^{(i)}}$  and depth  $D_{\mathbf{M}^{(i)}}, D_{\mathbf{Z}^{(i)}}$  respectively.

The algorithm  $\text{APPLYCHAIN} \left( (\mathbf{M}^{(1)}, \mathbf{Z}^{(1)}), \dots, (\mathbf{M}^{(d)}, \mathbf{Z}^{(d)}); F_1, \dots, F_d; \mathbf{b} \right)$  corresponds to a linear operator  $\mathbf{W}$  acting on  $\mathbf{b}$  such that

1.  $\mathbf{W}^{-1} \approx_{\sum_{i=1}^d 2\varepsilon_i} \mathbf{M}^{(0)}$ , and
2. for any vector  $\mathbf{b}$ ,  $\text{APPLYCHAIN} \left( (\mathbf{M}^{(1)}, \mathbf{Z}^{(1)}), \dots, (\mathbf{M}^{(d)}, \mathbf{Z}^{(d)}); F_1, \dots, F_d; \mathbf{b} \right)$  runs in  $O \left( \sum_{i=1}^d (D_{\mathbf{M}^{(i)}} + D_{\mathbf{Z}^{(i)}}) \right)$  depth and  $O \left( \sum_{i=1}^d (W_{\mathbf{M}^{(i)}} + W_{\mathbf{Z}^{(i)}}) \right)$  work.

The pseudocode for procedure  $\text{APPLYCHAIN}$  that uses an  $\varepsilon$ -vertex sparsifier chain to approximately solve a system of equations in  $\mathbf{M}^{(0)}$  is given in Figure 13-1.

*Proof.* We begin by observing that the output vector  $\mathbf{x}^{(0)}$  is a linear transformation of the input vector  $\mathbf{b}^{(0)}$ . Let  $\mathbf{W}^{(0)}$  be the matrix that realizes this transformation. Similarly, for  $1 \leq i \leq d$ , define  $\mathbf{W}^{(i)}$  to be the matrix so that

$$\mathbf{x}^{(i)} = \mathbf{W}^{(i)} \mathbf{b}^{(i)}.$$

An examination of the algorithm reveals that

$$\mathbf{W}^{(d)} = \left( \mathbf{M}^{(d)} \right)^{-1}, \tag{13.4}$$

---

**Algorithm 35:**  $\mathbf{x}^{(0)} = \text{APPLYCHAIN} \left( (\mathbf{M}^{(1)}, \mathbf{Z}^{(1)}), \dots, (\mathbf{M}^{(d)}, \mathbf{Z}^{(d)}); F_1, \dots, F_d \right)$

---

```

Initialize  $\mathbf{b}^{(0)} \leftarrow \mathbf{b}$ .
for  $i = 1, \dots, d$  do
     $\mathbf{x}_{[F_i]}^{(i-1)} \leftarrow \mathbf{Z}^{(i)} \mathbf{b}_{[F_i]}^{(i-1)}$ ,
     $\mathbf{b}_{[C_i]}^{(i)} \leftarrow \mathbf{b}_{[C_i]}^{(i-1)} - \mathbf{M}_{[C_i, F_i]}^{(i-1)} \mathbf{x}_{[F_i]}^{(i-1)}$ .
end
 $\mathbf{x}^{(d)} \leftarrow (\mathbf{M}^{(d)})^{-1} \mathbf{b}^{(d)}$ .
for  $i = d, \dots, 1$  do
     $\mathbf{x}_{[C_i]}^{(i-1)} \leftarrow \mathbf{x}^{(i)}$ .
     $\mathbf{x}_{[F_i]}^{(i-1)} \leftarrow \mathbf{x}_{[F_i]}^{(i-1)} - \mathbf{Z}^{(i)} \mathbf{M}_{[F_i, C_i]}^{(i-1)} \mathbf{x}^{(i)}$ .
end

```

---

**Figure 13-1:** Solver Algorithm using Vertex Sparsifier Chain

and for  $1 \leq i \leq d$ ,

$$\mathbf{W}^{(i-1)} = \begin{bmatrix} \mathbf{I} & -\mathbf{Z}^{(i)} \mathbf{M}_{[F_i, C_i]}^{(i-1)} \\ 0 & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{Z}^{(i)} & 0 \\ 0 & \mathbf{W}^{(i)} \end{bmatrix} \begin{bmatrix} \mathbf{I} & 0 \\ -\mathbf{M}_{[C_i, F_i]}^{(i-1)} \mathbf{Z}^{(i)} & \mathbf{I} \end{bmatrix}. \quad (13.5)$$

We will now prove by backwards induction on  $i$  that

$$\left( \mathbf{W}^{(i)} \right)^{-1} \approx_{\sum_{j=i+1}^d 2\varepsilon_j} \mathbf{M}^{(i)}.$$

The base case of  $i = d$  follows from (13.4). Using the definition of an  $\varepsilon$ -SCC, we know that  $0 \preceq (\mathbf{Z}^{(i)})^{-1} - \mathbf{M}_{[F_i, F_i]}^{(i-1)} \preceq \varepsilon_i \cdot \text{Sc} \left( \mathbf{M}^{(i-1)}, C_i \right)$ . We show in Lemma 13.7.1 that this implies

$$\begin{bmatrix} \mathbf{I} & -\mathbf{Z}^{(i)} \mathbf{M}_{[F_i, C_i]}^{(i-1)} \\ 0 & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{Z}^{(i)} & 0 \\ 0 & \text{Sc} \left( \mathbf{M}^{(i-1)}, F_i \right)^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{I} & 0 \\ -\mathbf{M}_{[C_i, F_i]}^{(i-1)} \mathbf{Z}^{(i)} & \mathbf{I} \end{bmatrix} \approx_{\varepsilon_i} \left( \mathbf{M}^{(i-1)} \right)^{-1}.$$

As  $\mathbf{M}^{(i)} \approx_{\varepsilon_i} \text{Sc} \left( \mathbf{M}^{(i-1)}, F_i \right)$ ,

$$\begin{bmatrix} \mathbf{I} & -\mathbf{Z}^{(i)} \mathbf{M}_{[F_i, C_i]}^{(i-1)} \\ 0 & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{Z}^{(i)} & 0 \\ 0 & \left( \mathbf{M}^{(i)} \right)^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{I} & 0 \\ -\mathbf{M}_{[C_i, F_i]}^{(i-1)} \mathbf{Z}^{(i)} & \mathbf{I} \end{bmatrix} \approx_{2\varepsilon_i} \left( \mathbf{M}^{(i-1)} \right)^{-1}.$$

By combining this identity with (13.5) and our inductive hypothesis, we obtain

$$\mathbf{W}^{(i-1)} \approx_{\sum_{j=i}^d 2\varepsilon_j} \left( \mathbf{M}^{(i-1)} \right)^{-1}.$$

Thus, by induction, we obtain

$$\mathbf{W}^{(0)} \approx_{\sum_{j=1}^d 2\varepsilon_j} \left( \mathbf{M}^{(0)} \right)^{-1}.$$

The whole algorithm involves a constant number of applications of  $\mathbf{Z}^{(i)}$ ,  $\mathbf{M}_{[F_i, C_i]}^{(i-1)}$ , and  $\mathbf{M}_{[C_i, F_i]}^{(i-1)}$ . We observe that in order to compute  $\mathbf{M}_{[F_i, C_i]}^{(i-1)} \mathbf{x}_{C_i}$ , using a multiplication procedure for  $\mathbf{M}^{(i-1)}$ , we can

pad  $x_{C_i}$  with zeros, multiply by  $\mathbf{M}^{(i-1)}$ , and read off the answer on the indices in  $C_i$ . Similarly, we can multiply vectors with  $\mathbf{M}_{[C_i, F_i]}^{(i-1)}$ . This immediately gives the claimed bounds on work and depth.  $\square$

We now prove the deferred claims from the above proof.

**Lemma 13.7.1.** *Let  $\mathbf{M} \in (\mathbb{C}^{r \times r})^{n \times n}$  be a bDD matrix,  $F \subseteq V$  be a subset of the indices, and  $\mathbf{Z} \in (\mathbb{C}^{r \times r})^{|F| \times |F|}$  be an hermitian operator satisfying  $0 \preceq \mathbf{Z}^{-1} - \mathbf{M}_{[F, F]} \preceq \varepsilon \cdot \text{Sc}(\mathbf{M}, C)$ . Then,*

$$\begin{bmatrix} \mathbf{I} & -\mathbf{Z}\mathbf{M}_{[F, C]} \\ 0 & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{Z} & 0 \\ 0 & \text{Sc}(\mathbf{M}, F)^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{I} & 0 \\ -\mathbf{M}_{[C, F]}\mathbf{Z} & \mathbf{I} \end{bmatrix} \approx_{\varepsilon} \mathbf{M}^{-1}.$$

*Proof.* Define

$$\widehat{\mathbf{M}} = \begin{bmatrix} (\mathbf{Z})^{-1} & \mathbf{M}_{[F, C]} \\ \mathbf{M}_{[C, F]} & \mathbf{M}_{[C, C]} \end{bmatrix}.$$

Using Lemma 13.7.2, we know that the assumption on  $\mathbf{Z}$  is equivalent to

$$\mathbf{M} \preceq \widehat{\mathbf{M}} \preceq (1 + \varepsilon) \mathbf{M}.$$

By Eq. (13.1), this implies

$$\mathbf{M}^{-1} \succcurlyeq \begin{bmatrix} \mathbf{I} & -\mathbf{Z}\mathbf{M}_{[F, C]} \\ 0 & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{Z} & 0 \\ 0 & \text{Sc}(\widehat{\mathbf{M}}, F)^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{I} & 0 \\ -\mathbf{M}_{[C, F]}\mathbf{Z} & \mathbf{I} \end{bmatrix} \succcurlyeq (1 + \varepsilon)^{-1} \mathbf{M}^{-1}.$$

From Facts 13.5.3, we know that

$$\text{Sc}(\mathbf{M}, F)^{-1} \succcurlyeq \text{Sc}(\widehat{\mathbf{M}}, F)^{-1} \succcurlyeq (1 + \varepsilon)^{-1} \text{Sc}(\mathbf{M}, F)^{-1}.$$

Now, we substitute this inequality into the one above and obtain

$$(1 + \varepsilon) \mathbf{M}^{-1} \succcurlyeq \begin{bmatrix} \mathbf{I} & -\mathbf{Z}\mathbf{M}_{[F, C]} \\ 0 & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{Z} & 0 \\ 0 & \text{Sc}(\mathbf{M}, F)^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{I} & 0 \\ -\mathbf{M}_{[C, F]}\mathbf{Z} & \mathbf{I} \end{bmatrix} \succcurlyeq (1 + \varepsilon)^{-1} \mathbf{M}^{-1},$$

which implies the lemma.  $\square$

**Lemma 13.7.2.** *Given a bDD matrix  $\mathbf{M}$ , a partition of its indices  $(F, C)$  such that  $\mathbf{M}_{[F, F]}$ ,  $\mathbf{M}_{[C, C]}$  are invertible, and an invertible hermitian operator  $\mathbf{Z} \in (\mathbb{C}^{r \times r})^{|F| \times |F|}$ , the following two conditions are equivalent:*

1.  $0 \preceq (\mathbf{Z})^{-1} - \mathbf{M}_{[F, F]} \preceq \varepsilon \cdot \text{Sc}(\mathbf{M}, C)$ .

- 2.

$$\mathbf{M} \preceq \begin{bmatrix} (\mathbf{Z})^{-1} & \mathbf{M}_{[F, C]} \\ \mathbf{M}_{[C, F]} & \mathbf{M}_{[C, C]} \end{bmatrix} \preceq (1 + \varepsilon) \mathbf{M}.$$

*Proof.* Writing

$$\mathbf{M} = \begin{bmatrix} \mathbf{M}_{[F, F]} & \mathbf{M}_{[F, C]} \\ \mathbf{M}_{[C, F]} & \mathbf{M}_{[C, C]} \end{bmatrix},$$

condition 2 is equivalent to

$$0 \preceq \begin{bmatrix} (\mathbf{Z})^{-1} - \mathbf{M}_{[F, F]} & 0 \\ 0 & 0 \end{bmatrix} \preceq \varepsilon \mathbf{M}.$$

The left inequality in this statement is equivalent to the left inequality in condition 1. Thus, it suffices to prove the right sides are equivalent.

To this end, using  $\mathbf{M}_{[F,C]} = \mathbf{M}_{[C,F]}^*$ , it suffices to prove that  $\forall x \in (\mathbb{C}^r)^{|F|}, y \in (\mathbb{C}^r)^{|C|}$ ,

$$x^*((\mathbf{Z})^{-1} - \mathbf{M}_{[F,F]})x \leq \varepsilon(x^*\mathbf{M}_{[F,F]}x + 2x^*\mathbf{M}_{[F,C]}y + y^*\mathbf{M}_{[C,C]}y).$$

This is equivalent to proving  $\forall x \in (\mathbb{C}^r)^{|F|}$ ,

$$x^*((\mathbf{Z})^{-1} - \mathbf{M}_{[F,F]})x \leq \varepsilon \inf_{y \in (\mathbb{C}^r)^{|C|}} (x^*\mathbf{M}_{[F,F]}x + 2x^*\mathbf{M}_{[F,C]}y + y^*\mathbf{M}_{[C,C]}y).$$

Since  $\mathbf{M}_{[C,C]} \succcurlyeq 0$ , the rhs is a convex function of  $y$ . The minimum is achieved at  $y = -\mathbf{M}_{[C,C]}^{-1}\mathbf{M}_{[F,C]}^*x$ , and we obtain the equivalent condition

$$x^*((\mathbf{Z})^{-1} - \mathbf{M}_{[F,F]})x \leq \varepsilon x^*(\mathbf{M}_{[F,F]} - \mathbf{M}_{[F,C]}\mathbf{M}_{[C,C]}^{-1}\mathbf{M}_{[C,F]}^*)x.$$

Since  $Sc(\mathbf{M}, C) = \mathbf{M}_{[F,F]} - \mathbf{M}_{[F,C]}\mathbf{M}_{[C,C]}^{-1}\mathbf{M}_{[C,F]}^* = \mathbf{M}_{[F,F]} - \mathbf{M}_{[F,C]}\mathbf{M}_{[C,C]}^{-1}\mathbf{M}_{[F,C]}^*$ , we obtain our claim.  $\square$

## ■ 13.8 Finding $\alpha$ -bDD Subsets

In this section we check that a simple randomized sampling procedure leads to  $\alpha$ -bDD subsets. Specifically we will prove Lemma 13.3.2:

*Lemma 13.3.2.* Given a bDD matrix  $\mathbf{M}$  with  $n$  block-rows, and an  $\alpha \geq 0$ , BDDSUBSET computes a subset  $F$  of size least  $n/(8(1+\alpha))$  such that  $\mathbf{M}_{[F,F]}$  is  $\alpha$ -bDD. It runs in runs in  $O(m)$  expected work and  $O(\log n)$  expected depth, where  $m$  is the number of nonzero blocks in  $\mathbf{M}$ .

Pseudocode for this routine is given in Figure 13-2.

---

**Algorithm 36:**  $F = \text{BDDSUBSET}(\mathbf{M}, \alpha)$ , where  $\mathbf{M}$  is a bDD matrix with  $n$  rows.

---

Let  $F'$  be a uniform random subset of  $\{1, \dots, n\}$  of size  $\frac{n}{4(1+\alpha)}$ .

Set

$$F = \left\{ i \in F' \text{ such that } \sum_{j:j \neq i} \|\mathbf{M}_{[i,j]}\| \geq (1+\alpha) \sum_{j \in F': j \neq i} \|\mathbf{M}_{[i,j]}\| \right\}.$$

If  $|F| < \frac{n}{8(1+\alpha)}$ , goto Step 36.

**Output**  $F$ .

---

**Figure 13-2:** Routine for finding an  $\alpha$ -strongly block diagonally dominant submatrix

We first show that the set returned is guaranteed to be  $\alpha$ -strongly block diagonally dominant.

**Lemma 13.8.1.** *If BDDSUBSET terminates, it returns  $F$  such that  $\mathbf{M}_{[F,F]}$  is  $\alpha$ -bDD.*

*Proof.* Consider some  $i \in F$ , the criteria for including  $i$  in  $F$  in Step 36 gives:

$$\sum_{j:j \neq i} \|\mathbf{M}_{[i,j]}\| \geq (1+\alpha) \sum_{j \in F', j \neq i} \|\mathbf{M}_{[i,j]}\| \geq (1+\alpha) \sum_{j \in F, j \neq i} \|\mathbf{M}_{[i,j]}\|,$$

where the last inequality follows since  $F$  is a subset of  $F'$ .

Incorporating this into the definition of  $\mathbf{M}$  being bDD gives

$$\mathbf{M}_{[i,i]} \succcurlyeq \mathbf{I}_r \cdot \sum_{j:j \neq i} \|\mathbf{M}_{[i,j]}\| \geq (1+\alpha)\mathbf{I}_r \cdot \sum_{j \in F, j \neq i} \|\mathbf{M}_{[i,j]}\|.$$

which means  $\mathbf{M}_{[F,F]}$  is  $\alpha$ -bDD.  $\square$

It remains to show that the algorithm finds a big  $F$  quickly. This can be done by upper bounding the expected size of  $F$ , or the probability of a single index being in  $F$ .

**Lemma 13.8.2.** *For any  $i$ , we have*

$$\Pr [i \in F' \text{ and } i \notin F] \leq \frac{1}{16(1+\alpha)}.$$

*Proof.* This event only happens if  $i \in F'$  and

$$\sum_{j \in F', j \neq i} \|\mathbf{M}_{[i,j]}\| > \frac{1}{1+\alpha} \sum_{j:j \neq i} \|\mathbf{M}_{[i,j]}\|. \quad (13.6)$$

Conditioning on  $i$  being selected initially, or  $i \in F'$ , the probability that each other  $j \neq i$  is in  $F'$  is

$$\frac{1}{n-1} \left( \frac{n}{4(1+\alpha)} - 1 \right) < \frac{1}{4(\alpha+1)},$$

which gives:

$$\mathbb{E} \left[ \sum_{j \in F', j \neq i} \|\mathbf{M}_{[i,j]}\| \middle| i \in F' \right] < \frac{1}{4(1+\alpha)} \sum_{j:j \neq i} \|\mathbf{M}_{[i,j]}\|.$$

Markov's inequality then gives:

$$\Pr \left[ \sum_{j \in F', j \neq i} \|\mathbf{M}_{[i,j]}\| > \frac{1}{1+\alpha} \sum_{j:j \neq i} \|\mathbf{M}_{[i,j]}\| \middle| i \in F' \right] < 1/4,$$

and thus

$$\Pr [i \in F' \text{ and } i \notin F] = \Pr [i \notin F | i \in F'] \Pr [i \in F'] < \frac{1}{4} \frac{1}{4(1+\alpha)} = \frac{1}{16(1+\alpha)}.$$

$\square$

Combining these two bounds gives Lemma 13.3.2.

*Proof.* (of Lemma 13.3.2) Applying Linearity of Expectation to Lemma 13.8.2 gives

$$\mathbb{E} [ |F' \setminus F| ] \leq \frac{n}{16(1+\alpha)}.$$

Markov's inequality then gives

$$\Pr \left[ |F' \setminus F| \geq \frac{n}{8(1+\alpha)} \right] < 1/2.$$

So, with probability at least  $1/2$ ,  $|F| \geq n/(8(1+\alpha))$ , and the algorithm will pass the test in line 3. Thus, the expected number of iterations made by the algorithm is at most 2. The claimed bounds on the expected work and depth of the algorithm follow.  $\square$



## ■ 13.9 Jacobi Iteration on $\alpha$ -bDD Matrices

From an  $\alpha$ -bDD set  $F$ , we will construct an operator  $\mathbf{Z}$  that approximates  $\mathbf{M}_{[F,F]}^{-1}$  and that can be applied quickly. Specifically, we will show:

**Theorem 13.3.1.** *Let  $\mathbf{M}$  be a bDD matrix with index set  $V$ , and let  $F \subseteq V$  such that  $\mathbf{M}_{[F,F]}$  is  $\alpha$ -bDD for some  $\alpha \geq 4$ , and has  $m_{FF}$  nonzero blocks. The algorithm  $\text{JACOBI}(\varepsilon, \mathbf{M}_{[F,F]}, \mathbf{b})$  acts as a linear operator  $\mathbf{Z}$  on  $\mathbf{b}$  that satisfies*

$$0 \preceq (\mathbf{Z})^{-1} - \mathbf{M}_{[F,F]} \preceq \varepsilon \cdot \text{Sc}(\mathbf{M}, V \setminus F).$$

The algorithm takes  $O(m_{FF} \log(\frac{1}{\varepsilon}))$  work and  $O(\log n \log(\frac{1}{\varepsilon}))$  depth.

Pseudocode of this routine is given in Figure 13-3.

---

**Algorithm 37:**  $\mathbf{x} = \text{JACOBI}(\varepsilon, \mathbf{M}, \mathbf{b})$

---

Create the matrix  $\mathbf{L}$  where

$$\mathbf{L}_{[i,j]} = \begin{cases} \mathbf{I}_r \cdot \sum_{j:j \neq i} \|\mathbf{M}_{[i,j]}\| & \text{if } i = j, \\ \mathbf{M}_{[i,j]} & \text{otherwise.} \end{cases}$$

Set  $\mathbf{X} = \mathbf{M} - \mathbf{L}$ .

Set  $k$  to be an odd integer that is greater than  $\log(3/\varepsilon)$ .

Set  $\mathbf{x}^{(0)} = \mathbf{X}^{-1}\mathbf{b}$ .

**for**  $i = 1 \dots k$  **do**

  | Set  $\mathbf{x}^{(i)} = -\mathbf{X}^{-1}\mathbf{L}\mathbf{x}^{(i-1)} + \mathbf{X}^{-1}\mathbf{b}$ .

**end**

**Output**  $\mathbf{x}^{(k)}$ .

---

**Figure 13-3:** Jacobi Iteration for Solving Linear Systems in an  $\alpha$ -bDD Matrix

We first verify that any  $\alpha$ -bDD matrix has a good block-diagonal preconditioner.

*Lemma 13.3.3.* Every  $\alpha$ -bDD matrix  $\mathbf{M}$  can be written in the form  $\mathbf{X} + \mathbf{L}$  where  $\mathbf{X}$  is block-diagonal,  $\mathbf{L}$  is bDD, and  $\mathbf{X} \succcurlyeq \frac{\alpha}{2}\mathbf{L}$ .

*Proof.* Write  $\mathbf{L} = \mathbf{Y} - \mathbf{A}$  where  $\mathbf{Y}$  is block-diagonal and  $\mathbf{A}$  has its diagonal blocks as zeros. Note that  $\mathbf{L}$  is bDD by definition. Thus, by Corollary 13.6.5,  $\mathbf{L} = \mathbf{Y} - \mathbf{A} \succcurlyeq 0$ .

Using Lemma 13.6.2, we know that  $\mathbf{Y} + \mathbf{A} \succcurlyeq 0$ , or  $\mathbf{Y} \succcurlyeq -\mathbf{A}$ . This implies  $2\mathbf{Y} \succcurlyeq \mathbf{L}$ .

As  $\mathbf{M}$  is  $\alpha$ -strongly diagonally dominant and  $\mathbf{M}_{[i,i]} = \mathbf{X}_{[i,i]} + \mathbf{Y}_{[i,i]}$ , we have

$$(\mathbf{X} + \mathbf{Y})_{[i,i]} \succcurlyeq (1 + \alpha) \mathbf{I}_r \cdot \sum_{j \neq i} \|\mathbf{M}_{[i,j]}\| = (1 + \alpha) \mathbf{Y}_{[i,i]}.$$

Manipulating this then gives:

$$\mathbf{X} \succcurlyeq \alpha \mathbf{Y} \succcurlyeq \frac{\alpha}{2} \mathbf{L}.$$

□

We now move on to measuring the quality of the operator generated by JACOBI. It can be checked

that running it  $k$  steps gives the operator

$$\mathbf{Z}^{(k)} \stackrel{\text{def}}{=} \sum_{i=0}^k \mathbf{X}_{[F,F]}^{-1} \left( -\mathbf{L}_{[F,F]} \mathbf{X}_{[F,F]}^{-1} \right)^i, \quad (13.7)$$

which is equivalent to evaluating a truncation of the Neumann series for  $\mathbf{M}^{-1}$

**Lemma 13.9.1.** *Let  $\mathbf{M}$  be a matrix with splitting  $\mathbf{M} = \mathbf{X} + \mathbf{L}$  where  $0 \preceq \mathbf{L} \preceq \beta \mathbf{X}$  for some parameter  $1 > \beta > 0$ . Then, for odd  $k$  and for  $\mathbf{Z}^{(k)}$  as defined in (13.7) we have:*

$$\mathbf{X} + \mathbf{L} \preceq (\mathbf{Z}^{(k)})^{-1} \preceq \mathbf{X} + (1 + \delta) \mathbf{L} \quad (13.8)$$

where

$$\delta = \beta^k \frac{1 + \beta}{1 - \beta^{k+1}}.$$

*Proof.* The left-hand inequality is equivalent to the statement that all the eigenvalues of  $\mathbf{Z}^{(k)}(\mathbf{X} + \mathbf{L})$  are at most 1 (see [34, Lemma 2.2] or [243, Proposition 3.3]). To see that this is the case, expand

$$\begin{aligned} \mathbf{Z}^{(k)}(\mathbf{X} + \mathbf{L}) &= \left( \sum_{i=0}^k \mathbf{X}^{-1} (-\mathbf{L} \mathbf{X}^{-1})^i \right) (\mathbf{X} + \mathbf{L}) \\ &= \sum_{i=0}^k (-\mathbf{X}^{-1} \mathbf{L})^i - \sum_{i=1}^{k+1} (-\mathbf{X}^{-1} \mathbf{L})^i \\ &= \mathbf{I} - (\mathbf{X}^{-1} \mathbf{L})^{k+1}. \end{aligned}$$

As all the eigenvalues of an even power of a matrix are nonnegative, all of the eigenvalues of this last matrix are at most 1.

Similarly, the other inequality is equivalent to the assertion that all of the eigenvalues of  $\mathbf{Z}^{(k)}(\mathbf{X} + (1 + \delta)\mathbf{L})$  are at least one. Expanding this product yields

$$\left( \sum_{i=0}^k \mathbf{X}^{-1} (-\mathbf{L} \mathbf{X}^{-1})^i \right) (\mathbf{X} + (1 + \delta)\mathbf{L}) = \mathbf{I} - (\mathbf{X}^{-1} \mathbf{L})^{k+1} + \delta \sum_{i=0}^k (-1)^i (\mathbf{X}^{-1} \mathbf{L})^{i+1}$$

The eigenvalues of this matrix are precisely the numbers

$$1 - \lambda^{k+1} + \delta \sum_{i=0}^k (-1)^i \lambda^{i+1}, \quad (13.9)$$

where  $\lambda$  ranges over the eigenvalues of  $\mathbf{X}^{-1} \mathbf{L}$ . The assumption  $\mathbf{L} \preceq \beta \mathbf{X}$  implies that the eigenvalues of  $\mathbf{X}^{-1} \mathbf{L}$  are at most  $\beta$ , so  $0 \leq \lambda \leq \beta$ . We have chosen the value of  $\delta$  precisely to guarantee that, under this condition on  $\lambda$ , the value of (13.9) is at least 1.  $\square$

This error crucially depends only on  $\mathbf{L}$ , which for any choice of  $F$  can be upper bounded by  $\mathbf{M}$ . Propagating the error this way allows us to prove the guarantees for JACOBI

*Proof.* (of Theorem 13.3.1) Consider the matrix  $\mathbf{L}_{[F,F]}$  generated when calling JACOBI with  $\mathbf{M}_{[F,F]}$ .  $\mathbf{M}$  being bDD means for each  $i$  we have

$$\mathbf{M}_{[i,i]} \succcurlyeq \mathbf{I}_r \cdot \sum_{j \neq i} \|\mathbf{M}_{[i,j]}\|,$$

which means for all  $i \in F$

$$\mathbf{M}_{[i,i]} - \mathbf{I}_r \cdot \sum_{j \neq i, j \in F} \|\mathbf{M}_{[i,j]}\| \succcurlyeq \mathbf{I}_r \cdot \sum_{j \notin F} \|\mathbf{M}_{[i,j]}\|.$$

Therefore if we extend  $\mathbf{L}_{[F,F]}$  onto the full matrix by putting zeros everywhere else, we have  $\mathbf{L} \preccurlyeq \mathbf{M}$ . Fact 13.5.3 then gives  $\mathbf{L}_{[F,F]} \preccurlyeq \text{Sc}(\mathbf{M}, V \setminus F)$ .

Lemma 13.3.3 gives that  $\mathbf{X}_{[F,F]} \succcurlyeq \frac{\alpha}{2} \mathbf{L}_{[F,F]}$ . As  $\alpha \geq 4$ , we can invoke Lemma 13.9.1 with  $\beta = \frac{1}{2}$ . Since  $\frac{1+\beta}{1-\beta^{k+1}} \leq (\frac{3}{2})/(\frac{1}{2}) \leq 3$ , our choice of  $k = \log(3/\varepsilon)$  gives the desired error. Each of these steps involve a matrix vector multiplication in  $\mathbf{L}_{[F,F]}$  and two linear system solves in  $\mathbf{X}_{[F,F]}$ . The former takes  $O(m)$  work and  $O(\log n)$  depth since the blocks of  $\mathbf{L}_{[F,F]}$  are a subset of the blocks of  $\mathbf{M}$ , while the latter takes  $O(n)$  work and  $O(\log n)$  depth due to  $\mathbf{X}_{[F,F]}$  being block-diagonal.  $\square$

## ■ 13.10 Existence of Linear-Sized Approximate Inverses

In this section we prove Theorem 13.1.3, which tells us that every *bDD* matrix has a linear-sized approximate inverse. In particular, this implies that for every *bDD* matrix there is a linear-time algorithm that approximately solves systems of equations in that matrix. To save space, we will not dwell on this algorithm, but rather will develop the linear-sized approximate inverses directly. There is some cost in doing so: there are very large constants in the linear-sized approximate inverses that are not present in a simpler linear-time solver.

To obtain a  $\mathbf{U}^T \mathbf{D} \mathbf{U}$  factorization from an  $\varepsilon$ -SCC in which each  $\mathbf{M}^{(i)}_{[F_i F_i]}$  is 4-bDD, we employ the procedure in Figure 13-4.

---

**Algorithm 38:**  $(\mathbf{D}, \mathbf{U}) = \text{DECOMPOSE}(\mathbf{M}^{(1)}, \dots, \mathbf{M}^{(d)}, F_1, \dots, F_{d-1})$ , where each  $\mathbf{M}^{(i)}$  is a bDD matrix and each  $\mathbf{M}^{(i)}_{[F_i F_i]}$  is 4-bDD.

---

Use (13.7) to compute the matrix  $\mathbf{Z}^{(i)}$  such that  $\text{JACOBI}(\varepsilon_i, \mathbf{M}, \mathbf{b}) = \mathbf{Z}^{(i)} \mathbf{b}$ .

For each  $i < d$ , write  $\mathbf{M}^{(i)} = \mathbf{X}^{(i)} + \mathbf{L}^{(i)}$ , where  $\mathbf{X}^{(i)}$  is block diagonal and  $\mathbf{L}^{(i)}$  is *bDD*, as in JACOBI.

Let  $\mathbf{X}^{(d)} = \mathbf{I}_{|C_{d-1}|}$  and let  $\widehat{\mathbf{U}}$  be the upper-triangular Cholesky factor of  $\mathbf{M}^{(d)}$ .

Let  $\mathbf{D}$  be the block diagonal matrix with  $\mathbf{D}_{[F_i, F_i]} = \mathbf{X}^{(i)}$ , for  $1 \leq i < d$ , and

$$\mathbf{D}_{[C_{d-1}, C_{d-1}]} = \mathbf{I}_{|C_{d-1}|}.$$

Let  $\mathbf{U}$  be the upper-triangular matrix with 1s on the diagonal,  $\mathbf{U}_{[C_{d-1}, C_{d-1}]} = \widehat{\mathbf{U}}$ , and

$$\mathbf{U}_{[F_i, C_i]} = \mathbf{Z}^{(i)} \mathbf{M}^{(i)}_{[F_i, C_i]}, \text{ for } 1 \leq i < d.$$


---

**Figure 13-4:** Converting a vertex sparsifier chain into  $\mathbf{U}$  and  $\mathbf{D}$ .

**Lemma 13.10.1.** *On input an  $\varepsilon$ -SCC of  $\mathbf{M}^{(0)}$  in which each  $\mathbf{M}^{(i)}_{[F_i F_i]}$  is 4-bDD, the algorithm DECOMPOSE produces matrices  $\mathbf{D}$  and  $\mathbf{U}$  such that*

$$\mathbf{U}^T \mathbf{D} \mathbf{U} \approx_{\gamma} \mathbf{M},$$

where

$$\gamma \leq 2 \sum_{i=0}^{d-1} \varepsilon_i + \max_i \varepsilon_i + 1/2.$$

*Proof.* Consider the inverse of the operator  $\mathbf{W} = \mathbf{W}^{(1)}$  realized by the algorithm APPLYCHAIN, and the operators  $\mathbf{W}^{(i)}$  that appear in the proof of Lemma 13.3.1.

We have

$$\left(\mathbf{W}^{(i)}\right)^{-1} = \begin{bmatrix} \mathbf{I} & 0 \\ \mathbf{M}_{[C_i, F_i]} \mathbf{Z}^{(i)} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \left(\mathbf{Z}^{(i)}\right)^{-1} & 0 \\ 0 & \left(\mathbf{W}^{(i+1)}\right)^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{Z}^{(i)} \mathbf{M}_{[F_i, C_i]} \\ 0 & \mathbf{I} \end{bmatrix},$$

and

$$\left(\mathbf{W}^{(d)}\right)^{-1} = \mathbf{M}^{(d)} = \widehat{\mathbf{U}}^T \widehat{\mathbf{U}}.$$

After expanding and multiplying the matrices in this recursive factorization, we obtain

$$\left(\mathbf{W}^{(1)}\right)^{-1} = \mathbf{U}^T \begin{bmatrix} \left(\mathbf{Z}^{(1)}\right)^{-1} & \dots & 0 & 0 \\ 0 & \ddots & 0 & 0 \\ 0 & \dots & \left(\mathbf{Z}^{(d-1)}\right)^{-1} & 0 \\ 0 & \dots & 0 & \mathbf{I}_{|C_{d-1}|} \end{bmatrix} \mathbf{U}.$$

Moreover, we know that this latter matrix is a  $2 \sum_{i=0}^{d-1} \varepsilon_i$  approximation of  $\mathbf{M}$ . It remains to determine the impact of replacing the matrix in the middle of this expression with  $\mathbf{D}$ .

Lemma 13.9.1 implies that each  $\mathbf{M}_{[F_i, F_i]} \approx_{\varepsilon_i} \left(\mathbf{Z}^{(i)}\right)^{-1}$  and Lemma 13.3.3 implies that  $\mathbf{X}^{(i)} \approx_{1/2} \mathbf{M}_{[F_i, F_i]}$ . So, the loss in approximation quality when we substitute the diagonal matrices is  $\max_i \varepsilon_i + 1/2$ .  $\square$

Invoking this decomposition procedure in conjunction with the the near-optimal sparsification routine from Theorem 13.3.1 gives a nearly-linear work routine. Repeatedly picking subsets using Lemma 13.3.4 gives then gives the linear sized decomposition.

**Theorem 13.1.3** (Sparsified Cholesky). *For every bDD matrix  $\mathbf{M}$  with  $n$  block-rows there exists a diagonal matrix  $\mathbf{D}$  and an upper triangular matrix  $\mathbf{U}$  with  $O(n)$  nonzero blocks so that*

$$\mathbf{U}^T \mathbf{D} \mathbf{U} \approx_{3/4} \mathbf{M}.$$

*Moreover, linear equations in  $\mathbf{U}$ ,  $\mathbf{U}^T$ , and  $\mathbf{D}$  can be solved with linear work in depth  $O(\log^2 n)$ , and these matrices can be computed in polynomial time.*

*Proof.* We set  $\alpha = 4$  throughout and  $\varepsilon_i = 1/8(i+2)^2$ . Theorem 13.3.1 then guarantees that the average number of nonzero blocks in each column of  $\mathbf{M}^{(i)}$  is at most  $10r/\varepsilon_i^2 = 640r(i+2)^4$ . If we now apply Lemma 13.3.4 to find 4-diagonally dominant subsets  $F_i$  of each  $\mathbf{M}^{(i)}$ , we find that each such subset contains at least a  $1/80$  fraction of the block columns of its matrix and that each column and row of  $\mathbf{M}^{(i)}$  indexed by  $F$  has at most  $1280r(i+2)^4$  nonzero entries. This implies that each row of  $\mathbf{Z}^{(i)} \mathbf{M}_{[F_i, C_i]}$  has at most  $(1280r(i+2)^4)^{k_i+1}$  nonzero entries.

Let  $n_i$  denote the number of block columns of  $\mathbf{M}^{(i)}$ . By induction, we know that

$$n_i \leq n \left(1 - \frac{1}{80}\right)^{i-1}.$$

So, the total number of nonzero blocks in  $\mathbf{U}$  is at most

$$\sum_{i=1}^d n_i (1280r(i+2)^4)^{k_i+1} \leq n \sum_{i=1}^d \left(1 - \frac{1}{80}\right)^{i-1} (1280r(i+2)^4)^{k_i+1}.$$

We will show that the term multiplying  $n$  in this later expression is upper bounded by a constant. To see this, note that  $k_i \leq 1 + \log(2\varepsilon_i^{-1}) \leq \nu \log(i+1)$  for some constant  $\nu$ . So, there is some other constant  $\mu$  for which

$$(1280r(i+2)^4)^{k_i+1} \leq \exp(\mu \log^2(i+1)).$$

This implies that the sum is at most

$$\sum_{i \geq 1} \exp(\mu \log^2(i+1) - i/80),$$

which is bounded by a constant.

To bound the quality of the approximation, we compute

$$2 \sum_i \varepsilon_i + \max_i \varepsilon_i + 1/2 = 2 \sum_i 1/8(i+2)^2 + 1/72 + 1/2 < 3/4.$$

The claimed bound on the work to perform backwards and forwards substitution with  $\mathbf{U}$  is standard: these operations require work linear in the number of nonzero entries of  $\mathbf{U}$ . The bound on the depth follows from the fact that the substitutions can be performed level-by-level, take depth  $O(\log n)$  for each level, and the number of levels,  $d$ , is logarithmic in  $n$ .  $\square$

## ■ 13.11 Spectral Vertex Sparsification Algorithm

In this section, we give a proof of the following lemma that immediately implies Lemma 13.3.4.

**Lemma 13.11.1.** *Let  $\mathbf{M}$  be a bDD matrix with index set  $V$ , and  $m$  nonzero blocks. Let  $F \subseteq V$  be such that  $\mathbf{M}_{[F,F]}$  is  $\alpha$ -bDD for some  $\alpha \geq 4$ . The algorithm APPROXSCHUR( $\mathbf{M}, F, \varepsilon$ ), returns a matrix  $\widetilde{\mathbf{M}}_{SC}$  s.t.*

1.  $\widetilde{\mathbf{M}}_{SC}$  has  $O(m(\varepsilon^{-1} \log \log \varepsilon^{-1})^{O(\log \log \varepsilon^{-1})})$  nonzero blocks, and
2.  $\widetilde{\mathbf{M}}_{SC} \approx_\varepsilon Sc(\mathbf{M}, F)$ ,

in  $O(m(\varepsilon^{-1} \log \log \varepsilon^{-1})^{O(\log \log \varepsilon^{-1})})$  work and  $O(\log n(\log \log \varepsilon^{-1}))$  depth.

Moreover, if  $\mathbf{T}$  is a matrix such that only the submatrix  $\mathbf{T}_{[C,C]}$  is nonzero, and  $\mathbf{M} + \mathbf{T}$  is bDD, then APPROXSCHUR( $\mathbf{M} + \mathbf{T}, F, \varepsilon$ ) = APPROXSCHUR( $\mathbf{M}, F, \varepsilon$ ) +  $\mathbf{T}_{[C,C]}$ .

We show how to sparsify the Schur complement of  $\mathbf{M}$  after eliminating a set of indices  $F$  such that  $\mathbf{M}_{[F,F]}$  is an  $\alpha$ -bDD matrix. The procedure APPROXSCHUR is described in Figure 13-5 and uses two key subroutines SCHUR SQUARE and LASTSTEP. SCHUR SQUARE allows us to approximate  $Sc(\mathbf{M}, F)$  as the Schur complement of another matrix  $\mathbf{M}_1$  such that  $\mathbf{M}_{1[F,F]}$  is roughly  $\alpha^2$ -bDD. LASTSTEP allows us to approximate  $Sc(\mathbf{M}, F)$  with  $O(\varepsilon)$  error, when  $\mathbf{M}_{[F,F]}$  is roughly  $1/\varepsilon$ -bDD.

**Lemma 13.11.2.** *Let  $\mathbf{M}$  be an bDD matrix with index set  $V$ , and  $m$  nonzero blocks and let  $F \subseteq V$  be such that  $\mathbf{M}_{[F,F]}$  is  $\alpha$ -bDD for some  $\alpha \geq 4$ . Given  $\varepsilon < 1/2$ , the algorithm SCHUR SQUARE( $\mathbf{M}, F, \varepsilon$ ), returns a bDD matrix  $\mathbf{M}_1$  in  $O(m\varepsilon^{-4})$  work and  $O(\log n)$  depth, such that*

1.  $Sc(\mathbf{M}, F) \approx_\varepsilon Sc(\mathbf{M}_1, F)$ , and
2.  $(\mathbf{M}_1)_{[F,F]}$  is  $\alpha^2/2$ -bDD.
3.  $\mathbf{M}_1$  has  $O(m\varepsilon^{-4})$  nonzero blocks,

---

**Algorithm 39:**  $\widetilde{\mathbf{M}}_{SC} = \text{APPROXSCHUR}(\mathbf{M}, F, \alpha, \varepsilon)$ 


---

 Initialize  $\mathbf{M}^{(0)} \leftarrow \mathbf{M}$ ,  $d = 2 \log_2 \log_{\alpha/2}(4\varepsilon^{-1})$ 

 for  $i = 1, \dots, d$  do

 $\quad \left| \mathbf{M}^{(i)} \leftarrow \text{SCHUR SQUARE} \left( \mathbf{M}^{(i-1)}, F, \frac{\varepsilon}{2^i} \right).$ 

end

 $\widetilde{\mathbf{M}}_{SC} \leftarrow \text{LASTSTEP} \left( \mathbf{M}^{(d)}, F, 4\varepsilon^{-1}, \frac{\varepsilon}{4} \right).$ 
**Output**  $\widetilde{\mathbf{M}}_{SC}$ .
 

---

**Figure 13-5:** Pseudocode for Computing Spectral Vertex Sparsifiers

If  $\mathbf{T}$  is a matrix such that only the submatrix  $\mathbf{T}_{[C,C]}$  is nonzero, and  $\mathbf{M} + \mathbf{T}$  is bDD, then

$$\text{SCHUR SQUARE}(\mathbf{M} + \mathbf{T}, F, \varepsilon) = \text{SCHUR SQUARE}(\mathbf{M}, F, \varepsilon) + \mathbf{T}_{[C,C]}.$$

We can repeatedly applying the above lemma, to approximate  $Sc(\mathbf{M}, F)$  as  $Sc(\mathbf{M}_1, F)$ , where  $\mathbf{M}_1$  is  $O(\varepsilon^{-1})$ -bDD. LASTSTEP allows us to approximate the Schur complement for such a strongly block diagonally dominant matrix  $\mathbf{M}_1$ . The guarantees of LASTSTEP are given by the following lemma.

**Lemma 13.11.3.** *Let  $\mathbf{M}$  be an bDD matrix with index set  $V$ , and  $m$  nonzero blocks and let  $F \subseteq V$  be such that  $\mathbf{M}_{[F,F]}$  is  $\alpha$ -bDD for some  $\alpha \geq 4$ . There exist a procedure LASTSTEP such that LASTSTEP( $\mathbf{M}, F, \alpha, \varepsilon$ ) returns in  $O(m\varepsilon^{-8})$  work and  $O(\log n)$  depth a matrix  $\widetilde{\mathbf{M}}_{SC}$  s.t.  $\widetilde{\mathbf{M}}_{SC}$  has  $O(m\varepsilon^{-8})$  nonzero blocks and  $\widetilde{\mathbf{M}}_{SC} \approx_{\varepsilon+2/\alpha} Sc(\mathbf{M}, F)$ . If  $\mathbf{T}$  is a matrix such that only the submatrix  $\mathbf{T}_{[C,C]}$  is nonzero, and  $\mathbf{M} + \mathbf{T}$  is bDD, then LASTSTEP( $\mathbf{M} + \mathbf{T}, F, \alpha, \varepsilon$ ) = LASTSTEP( $\mathbf{M}, F, \alpha, \varepsilon$ ) +  $\mathbf{T}_{[C,C]}$ .*

Combining the above two lemmas, we obtain a proof of Lemma 13.3.4.

*Proof.* (of Lemma 13.11.1). By induction, after  $i$  steps of the main loop in APPROXSCHUR,

$$Sc(\mathbf{M}, F) \approx_{\frac{\varepsilon^i}{2^i}} Sc(\mathbf{M}^{(i)}, F).$$

Lemma 13.11.2 also implies that  $\mathbf{M}^{(i)}$  is  $2(\frac{\alpha}{2})^{2^i}$ -bDD. Thus, we have that  $\mathbf{M}_{FF}^{(d)}$  is  $8\varepsilon^{-1}$ -strongly diagonally dominant at the last step. Hence, Lemma 13.11.3 then gives

$$Sc(\mathbf{M}^{(d)}, F) \approx_{\frac{1}{2}\varepsilon} \text{LASTSTEP} \left( \mathbf{M}^{(d)}, F, 4\varepsilon^{-1}, \frac{\varepsilon}{4} \right).$$

Composing this bound with the guarantees of the iterations then gives the bound on overall error.

The property that if  $\mathbf{T}$  is a matrix such that only the submatrix  $\mathbf{T}_{[C,C]}$  is nonzero, and  $\mathbf{M} + \mathbf{T}$  is bDD, then APPROXSCHUR( $\mathbf{M} + \mathbf{T}, F, \varepsilon$ ) = APPROXSCHUR( $\mathbf{M}, F, \varepsilon$ ) +  $\mathbf{T}_{[C,C]}$  follows from Lemma 13.11.2 and 13.11.3, which ensure that this property holds for all our calls to SCHUR SQUARE and LASTSTEP.

The work of these steps, and the size of the output graph follow from Lemma 13.11.2 and 13.11.3.  $\square$

### ■ 13.11.1 Iterative Squaring and Sparsification

In this section, we give a proof of Lemma 13.11.2. At the core of procedure is a *squaring* identity that preserves Schur complements, and efficient sparsification of special classes of bDD matrices that we

call product demand block-Laplacians.

**Definition 13.11.4.** The product demand block-Laplacian of a vector  $\mathbf{d} \in (\mathbb{C}^{r \times r})^n$ , is a bDD matrix  $\mathbf{L}_{G(\mathbf{d})} \in (\mathbb{C}^{r \times r})^{n \times n}$ , defined as

$$(\mathbf{L}_{G(\mathbf{d})})_{[i,j]} = \begin{cases} -\mathbf{d}_{[i]} \mathbf{d}_{[j]}^* & i \neq j, \\ \mathbf{I}_r \cdot \|\mathbf{d}_{[i]}\| \sum_{k:k \neq i} \|\mathbf{d}_{[k]}\| & \text{otherwise.} \end{cases}$$

**Definition 13.11.5.** Given a vector  $\mathbf{d} \in (\mathbb{C}^{r \times r})^n$  and an index set  $F \subseteq V$ , let  $C = V \setminus F$ . The bipartite product demand block-Laplacian of  $(\mathbf{d}, F)$  is a bDD matrix  $\mathbf{L}_{G(\mathbf{d}, F)} \in (\mathbb{C}^{r \times r})^{n \times n}$ , defined as

$$(\mathbf{L}_{G(\mathbf{d}, F)})_{[i,j]} = \begin{cases} \mathbf{I}_r \cdot \|\mathbf{d}_{[i]}\| \sum_{k \in C} \|\mathbf{d}_{[k]}\| & i = j, \text{ and } i \in F \\ \mathbf{I}_r \cdot \|\mathbf{d}_{[i]}\| \sum_{k \in F} \|\mathbf{d}_{[k]}\| & i = j, \text{ and } i \in C \\ 0 & i \neq j, \text{ and } (i, j \in F \text{ or } i, j \in C) \\ -\mathbf{d}_{[i]} \mathbf{d}_{[j]}^* & i \neq j, \text{ otherwise.} \end{cases}$$

In Section 13.11.4, we prove the following lemmas that allows us to efficiently construct sparse approximations to these matrices.

**Lemma 13.11.6.** *There is a routine  $\text{CLIQUESPARSIFICATION}(\mathbf{d}, \varepsilon)$  such that for any demand vector  $\mathbf{d} \in (\mathbb{C}^{r \times r})^n$ , and  $\varepsilon > 0$ ,  $\text{CLIQUESPARSIFICATION}(\mathbf{d}, \varepsilon)$  returns in  $O(n\varepsilon^{-4})$  work and  $O(\log n)$  depth a bDD matrix  $\mathbf{L}_H$  with  $O(n\varepsilon^{-4})$  nonzero blocks such that*

$$\mathbf{L}_H \approx_\varepsilon \mathbf{L}_{G(\mathbf{d})}.$$

**Lemma 13.11.7.** *There is a routine  $\text{BIPARTITECLIQUESPARSIFICATION}(\mathbf{d}, F, \varepsilon)$  such that for any demand vector  $\mathbf{d} \in (\mathbb{C}^{r \times r})^n$ ,  $\varepsilon > 0$ , and  $F \subseteq V$ ,  $\text{BIPARTITECLIQUESPARSIFICATION}(\mathbf{d}, \varepsilon)$  returns in  $O(n\varepsilon^{-4})$  work and  $O(\log n)$  depth a bDD matrix  $\mathbf{L}_H$  with  $O(n\varepsilon^{-4})$  nonzero blocks such that*

$$\mathbf{L}_H \approx_\varepsilon \mathbf{L}_{G(\mathbf{d}, F)}.$$

Moreover,  $(\mathbf{L}_H)_{[F,F]}$ , and  $(\mathbf{L}_H)_{[C,C]}$  are block-diagonal, where  $C = V \setminus F$ .

We now use these efficient sparse approximations to give a proof of the guarantees of the procedure  $\text{SCHURSQUARE}$ .

---

**Algorithm 40:**  $M_1 = \text{SCHURSQUARE}(\mathbf{M}, F, \varepsilon)$

---

Construct sparse representations of  $\mathbf{d}_F^{(j)}$ ,  $\mathbf{d}_C^{(j)}$ ,  $\mathbf{f}^{(j)}$  defined by Eqs. (13.11), (13.12), and (13.13).

Construct  $M_4$  as given by Eq. (13.14).

Compute  $M_3$  as given by Eq. (13.15).

For all  $j \in F$ , let  $\tilde{\mathbf{L}}_{G(\mathbf{d}_F^{(j)})} \leftarrow \text{CLIQUESPARSIFICATION}(\mathbf{d}_F^{(j)}, \varepsilon)$ .

For all  $j \in F$ , let  $\tilde{\mathbf{L}}_{G(\mathbf{d}_C^{(j)})} \leftarrow \text{CLIQUESPARSIFICATION}(\mathbf{d}_C^{(j)}, \varepsilon)$ .

For all  $j \in F$ , let  $\tilde{\mathbf{L}}_{G(\mathbf{f}^{(j)}, F)} \leftarrow \text{BIPARTITECLIQUESPARSIFICATION}(\mathbf{f}^{(j)}, \varepsilon)$ .

Construct and return  $M_1$  defined by Eq. (13.16).

---

**Figure 13-6:** Pseudocode for procedure  $\text{SCHURSQUARE}$  Iterative Squaring and Sparsification

*Proof.* (of Lemma 13.11.2) Let  $C$  denote  $V \setminus F$ . Write  $\mathbf{M}_{[F,F]}$  as  $\mathbf{D} - \mathbf{A}$ , where  $\mathbf{D}$  is a block-diagonal, and  $\mathbf{A}$  has its diagonal blocks as zero. The proof is based on the identity  $Sc(\mathbf{M}, F) = Sc(\mathbf{M}_2, F)$ , where  $\mathbf{M}_2$  is the following matrix.

$$\mathbf{M}_2 = \frac{1}{2} \begin{bmatrix} \mathbf{D} - \mathbf{A}\mathbf{D}^{-1}\mathbf{A} & \mathbf{M}_{[F,C]} + \mathbf{A}\mathbf{D}^{-1}\mathbf{M}_{[F,C]} \\ \mathbf{M}_{[C,F]} + \mathbf{M}_{[C,F]}\mathbf{D}^{-1}\mathbf{A} & 2\mathbf{M}_{[C,C]} - \mathbf{M}_{[C,F]}\mathbf{D}^{-1}\mathbf{M}_{[F,C]} \end{bmatrix}. \quad (13.10)$$

It is straightforward to prove that  $\mathbf{M}_2$  satisfies  $Sc(\mathbf{M}, F) = Sc(\mathbf{M}_2, F)$  (see Lemma 13.11.8). It is also straightforward to show that  $\mathbf{D} - \mathbf{A}\mathbf{D}^{-1}\mathbf{A}$  is  $((\alpha + 1)^2 - 1)$ -bDD. Thus,  $\mathbf{M}_2$  satisfies the first two requirements.

However,  $\mathbf{M}_2$  is likely to be a dense matrix, and we will not construct it in full. The key observation is that  $\mathbf{M}_2$  can be written as a sum of an explicit sparse bDD matrix, and several product demand block-Laplacians. Formally, for every  $j \in F$ , we define  $\mathbf{d}_F^{(j)}, \mathbf{d}_C^{(j)} \in (\mathbb{C}^{r \times r})^n$  as follows

$$(\mathbf{d}_F^{(j)})_{[i]} = \begin{cases} \mathbf{A}_{[i,j]}\mathbf{D}_{[j,j]}^{-1/2} & i \in F, \\ 0 & i \in C. \end{cases} \quad (13.11)$$

$$(\mathbf{d}_C^{(j)})_{[i]} = \begin{cases} 0 & i \in F, \\ \mathbf{M}_{[i,j]}\mathbf{D}_{[j,j]}^{-1/2} & i \in C. \end{cases} \quad (13.12)$$

For every  $j \in F$ , we need to define a bipartite product demand block-Laplacian given by  $\mathbf{f}^{(j)} \in (\mathbb{C}^{r \times r})^n$  defined as

$$\mathbf{f}_{[i]}^{(j)} = \begin{cases} \mathbf{A}_{[i,j]}\mathbf{D}_{[j,j]}^{-1/2} & i \in F, \\ -\mathbf{M}_{[i,j]}\mathbf{D}_{[j,j]}^{-1/2} & i \in C. \end{cases} \quad (13.13)$$

Letting  $\deg_j$  denote the number of nonzero blocks in  $\mathbf{M}_{[j]}$ , the number of nonzero blocks in each of  $\mathbf{d}_F^{(j)}, \mathbf{d}_C^{(j)}, \mathbf{f}^{(j)}$  is at most  $\deg_j$ . We construct a sparse representation of each of them using  $O(\deg_j)$  work. Thus, the total number of nonzero blocks in all these block-vectors is at most  $O(m)$  and we can explicitly construct sparse representations for them using  $O(m)$  work and  $O(\log n)$  depth.

We can now express  $\mathbf{M}_2$  as

$$\mathbf{M}_2 = \mathbf{M}_3 + \frac{1}{2} \sum_{j \in F} \mathbf{L}_{G(\mathbf{d}_F^{(j)})} + \frac{1}{2} \sum_{j \in F} \mathbf{L}_{G(\mathbf{d}_C^{(j)})} + \frac{1}{2} \sum_{j \in F} \mathbf{L}_{G(\mathbf{f}^{(j)}, F)}.$$

For all  $i, j \in V$ , we compute  $\beta_{i,j} = \|\mathbf{f}_i^{(j)}\|$ . Define  $\mathbf{M}_4$  to be the block diagonal matrix such that

$$\mathbf{M}_{4[i,i]} = \frac{1}{2} \mathbf{I}_r \cdot \sum_{j \in F} \beta_{i,j} \left( \left( \sum_{k \in V} \beta_{k,j} \right) - \beta_{i,j} \right). \quad (13.14)$$

Since at most  $m$  of  $\beta_{i,j}$  are nonzero, and we can compute  $\beta_{i,j}$  and construct  $\mathbf{M}_4$  using  $O(m)$  work and  $O(\log n)$  depth. It is easy to verify that we can express  $\mathbf{M}_3$  explicitly as

$$\mathbf{M}_3 = \frac{1}{2} \begin{bmatrix} \mathbf{D} & \mathbf{M}_{[F,C]} \\ \mathbf{M}_{[C,F]} & 2\mathbf{M}_{[C,C]} \end{bmatrix} - \mathbf{M}_4. \quad (13.15)$$

For each  $j$ , we use Lemmas 13.11.6 and 13.11.7 to construct, in  $O(\deg_j \varepsilon^{-4})$  work and  $O(\log n)$  depth, bDD matrices  $\tilde{\mathbf{L}}_{G(\mathbf{d}_F^{(j)})}, \tilde{\mathbf{L}}_{G(\mathbf{d}_C^{(j)})}, \tilde{\mathbf{L}}_{G(\mathbf{f}^{(j)}, F)}$ , each with  $O(\deg_j \varepsilon^{-4})$  nonzero blocks such that

$$\tilde{\mathbf{L}}_{G(\mathbf{d}_F^{(j)})} \approx_\varepsilon \mathbf{L}_{G(\mathbf{d}_F^{(j)})}, \quad \tilde{\mathbf{L}}_{G(\mathbf{d}_C^{(j)})} \approx_\varepsilon \mathbf{L}_{G(\mathbf{d}_C^{(j)})}, \quad \text{and} \quad \tilde{\mathbf{L}}_{G(\mathbf{f}^{(j)}, F)} \approx_\varepsilon \mathbf{L}_{G(\mathbf{f}^{(j)}, F)}.$$

We can now construct the matrix

$$\mathbf{M}_1 = \mathbf{M}_3 + \frac{1}{2} \sum_{j \in F} \tilde{\mathbf{L}}_{G(\mathbf{d}_F^{(j)})} + \frac{1}{2} \sum_{j \in F} \tilde{\mathbf{L}}_{G(\mathbf{d}_C^{(j)})} + \frac{1}{2} \sum_{j \in F} \tilde{\mathbf{L}}_{G(\mathbf{f}^{(j)}, F)}, \quad (13.16)$$



in  $O(m\varepsilon^{-4})$  time and  $O(\log n)$  depth.  $\mathbf{M}_1$  has  $O(m\varepsilon^{-4})$  nonzero blocks and is our required matrix.

We first show that  $\mathbf{M}_3$  is bDD. Using  $\beta_{k,j} \leq \|D_{[j,j]}^{-1/2} \mathbf{M}_{[i,k]}\|$  for all  $j \in F, k \in V$ , we have for all  $i \in F$ ,

$$2\|\mathbf{M}_{4[i,i]}\| \leq \sum_{j \in F} \beta_{i,j} \sum_{k \in V} \beta_{k,j} \leq \sum_{j \in F} \|\mathbf{A}_{[i,j]}\| \|D_{[j,j]}^{-1/2}\|^2 \sum_{k \in V} \|\mathbf{M}_{[j,k]}\| \leq \sum_{j \in F} \|\mathbf{A}_{[i,j]}\|, \quad (13.17)$$

where the last inequality uses that  $\mathbf{M}$  is bDD. Again, using  $\mathbf{M}$  is bDD, for all  $i \in F$ , we have

$$D_{[i,i]} \succcurlyeq I_r \sum_{j \in F} \|\mathbf{A}_{[i,j]}\| + I_r \sum_{k \in C} \|\mathbf{M}_{[i,k]}\|.$$

Now combining with Eqs. (13.15) and (13.17), we obtain  $\mathbf{M}_{3[i,i]} \succcurlyeq I_r \sum_{k \in C} \|\mathbf{M}_{[i,k]}\|$ . Thus,  $\mathbf{M}_3$  is bDD, and hence  $\mathbf{M}_1$  is bDD.

We also have  $\mathbf{M}_1 \approx_\varepsilon \mathbf{M}_2$ . Thus, using Fact 13.5.3, we obtain  $Sc(\mathbf{M}_1, F) \approx_\varepsilon Sc(\mathbf{M}_2, F)$ . It remains to show that  $\mathbf{M}_{1[F,F]}$  is  $\alpha^2/2$ -bDD.

The key observation is that only the matrices  $\tilde{\mathbf{L}}_{G(d^{(j)})}$  contribute to off-diagonal blocks in  $\mathbf{M}_{1[F,F]}$ . By construction,

$$\forall i, j \in F, \quad \mathbf{L}_{G(d^{(j)})[i,i]} = I_r \sum_{k \in F: k \neq i} \beta_{i,j} \beta_{k,j}.$$

Thus, for all  $i \in F$

$$\begin{aligned} \sum_{j \in F} \|\mathbf{L}_{G(d^{(j)})[i,i]}\| &\leq \sum_{j \in F} \sum_{k \in F} \beta_{i,j} \beta_{k,j} \leq \sum_{j \in F} \|\mathbf{A}_{[i,j]}\| \|D_{[j,j]}^{-1/2}\|^2 \sum_{k \in F} \|\mathbf{A}_{[j,k]}\| \\ &\leq (1 + \alpha)^{-1} I_r \cdot \sum_{j \in F} \|\mathbf{A}_{[i,j]}\|, \end{aligned} \quad (13.18)$$

where the last inequality follows since  $D - A = \mathbf{M}_{[F,F]}$  is  $\alpha$ -bDD. Moreover  $D - A$  being  $\alpha$ -bDD implies that for each  $i \in F$ ,  $D_{[i,i]} \succcurlyeq (\alpha + 1)I_r \sum_{j \in F} \|\mathbf{A}_{[i,j]}\|$ . Thus, using Eqs. (13.15) and (13.17), we obtain,  $2\mathbf{M}_{3[i,i]} \succcurlyeq \alpha I_r \sum_{j \in F} \|\mathbf{A}_{[i,j]}\|$ .

Since for each  $j$ ,  $\tilde{\mathbf{L}}_{G(d^{(j)})} \approx_\varepsilon \tilde{\mathbf{L}}_{G(d^{(j)})}$ , we have  $(\tilde{\mathbf{L}}_{G(d^{(j)})})_{[i,i]} \preccurlyeq e^\varepsilon (\mathbf{L}_{G(d^{(j)})})_{[i,i]}$ . Thus,  $\sum_j \tilde{\mathbf{L}}_{G(d^{(j)})}$  is bDD with sums of norms of off-diagonal blocks at most  $e^\varepsilon (1 + \alpha)^{-1} \sum_{j \in F} \|\mathbf{A}_{[i,j]}\|$ . Since these are the only matrices that contribute to off-diagonal blocks in  $\mathbf{M}_{1[F,F]}$ , we get that  $\mathbf{M}_{1[F,F]}$  is at least  $\frac{\alpha}{e^\varepsilon(1+\alpha)^{-1}}$ -bDD. Using  $\varepsilon < 1/2$  gives us our claim.

It is easy to verify that our transformations maintain that if  $\mathbf{T}$  is a matrix that is only nonzero inside the submatrix  $\mathbf{T}_{[C,C]}$ , and  $\mathbf{M} + \mathbf{T}$  is bDD, then

$$\text{SCHUR SQUARE}(\mathbf{M} + \mathbf{T}, F, \varepsilon) = \text{SCHUR SQUARE}(\mathbf{M}, F, \varepsilon) + \mathbf{T}_{[C,C]}.$$

□

We now prove the claims deferred from the above proof.

**Lemma 13.11.8.** *The matrix  $\mathbf{M}_2$  defined by Eq. (13.10) satisfies  $Sc(\mathbf{M}_2, F) = Sc(\mathbf{M}, F)$ .*

*Proof.* We need the following identity from [223]:

$$(D - A)^{-1} = 1/2 \cdot (D^{-1} + (I + D^{-1}A)(D - AD^{-1}A)^{-1}(I + AD^{-1})). \quad (13.19)$$

We have,

$$\begin{aligned}
Sc(\mathbf{M}_2, F) &= \mathbf{M}_{[C,C]} - 1/2 \cdot \mathbf{M}_{[C,F]} \mathbf{D}^{-1} \mathbf{M}_{[F,C]} \\
&\quad - 1/2 \cdot \mathbf{M}_{[C,F]} (\mathbf{I} + \mathbf{D}^{-1} \mathbf{A}) (\mathbf{D} - \mathbf{A} \mathbf{D}^{-1} \mathbf{A})^{-1} (\mathbf{I} + \mathbf{A} \mathbf{D}^{-1}) \mathbf{M}_{[F,C]} \\
&= \mathbf{M}_{[C,C]} \\
&\quad - 1/2 \cdot \mathbf{M}_{[C,F]} (\mathbf{D}^{-1} + (\mathbf{I} + \mathbf{D}^{-1} \mathbf{A}) (\mathbf{D} - \mathbf{A} \mathbf{D}^{-1} \mathbf{A})^{-1} (\mathbf{I} + \mathbf{A} \mathbf{D}^{-1})) \mathbf{M}_{[F,C]} \\
&= \mathbf{M}_{[C,C]} - \mathbf{M}_{[C,F]} (\mathbf{D} - \mathbf{A})^{-1} \mathbf{M}_{[F,C]} \\
&= Sc(\mathbf{M}, F).
\end{aligned}$$

□

### ■ 13.11.2 Schur Complement w.r.t. Highly $\alpha$ -bDD Submatrices

---

**Algorithm 41:**  $\widetilde{\mathbf{M}}_{SC} = \text{LASTSTEP}(\mathbf{M}, F, \alpha, \varepsilon)$

---

Compute  $\mathbf{X}$ ,  $\mathbf{D}$ , and  $\mathbf{A}$  as given by equations (13.20), (13.21), and (13.22).

Construct  $\mathbf{Y}$  as defined by equations (13.33), (13.34), and (13.35).

Construct sparse vectors  $\mathbf{d}^{(j)}$ ,  $\mathbf{g}^{(j)}$  for all  $j \in F$  as defined by equations (13.26) and (13.27).

For all  $j \in F$ , let  $\tilde{\mathbf{L}}_{G(\mathbf{d}^{(j)}, F)} \leftarrow \text{BIPARTITECLIQUESPARSIFICATION}(\mathbf{d}^{(j)}, F, \varepsilon/2)$

For all  $j \in F$ , let  $\tilde{\mathbf{L}}_{G(\mathbf{g}^{(j)})} \leftarrow \text{CLIQUESPARSIFICATION}(\mathbf{g}^{(j)}, \varepsilon/2)$ .

Compute  $\mathbf{R}$  as given by equation (13.29).

Construct sparse vectors  $\mathbf{r}^{(j)}$  for all  $j \in F$  as defined by equations (13.30).

For all  $j \in F$ , let  $\tilde{\mathbf{L}}_{G(\mathbf{r}^{(j)})} \leftarrow \text{CLIQUESPARSIFICATION}(\mathbf{r}^{(j)}, \varepsilon/2)$ .

Compute  $\mathbf{S}$  as given by equation (13.31).

**Output**

$$\widetilde{\mathbf{M}}_{SC} = \mathbf{S} + \sum_i \tilde{\mathbf{L}}_{G(\mathbf{r}^{(i)})}.$$


---

**Figure 13-7:** Pseudocode for procedure LASTSTEP: Computing an approximate Schur complement w.r.t. a highly  $\alpha$ -bDD submatrices.

In this section, we describe the LASTSTEP procedure for computing an approximate Schur complement of a bDD matrix  $\mathbf{M}$  with a highly  $\alpha$ -bDD submatrix  $\mathbf{M}_{[F,F]}$ . LASTSTEP is the final step of the APPROXSCHUR algorithm. The key element of the procedure is a formula for approximating the inverse of  $\mathbf{M}_{[F,F]}$  that is leveraged to approximate the Schur complement of  $\mathbf{M}$  as the Schur complement of matrix with the  $FF$  submatrix being block diagonal.

We prove guarantees for the LASTSTEP algorithm as stated in Lemma 13.11.3.

One could attempt to deal with the highly  $\alpha$ -bDD matrix at the last step by directly replacing it with its diagonal, but this is problematic. Consider the case where  $F$  contains  $u$  and  $v$  with a weight  $\varepsilon$  edge between them, and  $u$  and  $v$  are connected to  $u'$  and  $v'$  in  $C$  by weight 1 edges respectively. Keeping only the diagonal results in a Schur complement that disconnects  $u'$  and  $v'$ . This however can be fixed by taking a step of random walk within  $F$ .

Given a bDD matrix  $\mathbf{M}$ , s.t.  $\mathbf{M}_{[F,F]}$  is  $\alpha$ -bDD we define a block diagonal matrix  $\mathbf{X} \in (\mathbb{C}^{r \times r})^{|F| \times |F|}$  s.t. for each  $i \in F$

$$\mathbf{X}_{[i,i]} = \frac{\alpha}{\alpha + 1} \mathbf{M}_{[i,i]} \quad (13.20)$$

and another block diagonal matrix  $\mathbf{D} \in (\mathbb{C}^{r \times r})^{|F| \times |F|}$  s.t. for each  $i \in F$

$$\mathbf{D}_{[i,i]} = \frac{1}{\alpha + 1} \mathbf{M}_{[i,i]}, \quad (13.21)$$

and we define a matrix  $\mathbf{A} \in (\mathbb{C}^{r \times r})^{|F| \times |F|}$

$$\mathbf{A}_{[i,j]} = \begin{cases} 0 & i = j \\ -\mathbf{M}_{[i,j]} & \text{otherwise.} \end{cases} \quad (13.22)$$

Thus  $\mathbf{M}_{[F,F]} = \mathbf{X} + \mathbf{D} - \mathbf{A}$ . One can check that because  $\mathbf{M}$  is bDD and  $\mathbf{M}_{[F,F]}$  is  $\alpha$ -bDD, it follows that  $\mathbf{D} - \mathbf{A}$  is bDD and the matrix

$$\begin{bmatrix} \mathbf{X} & \mathbf{M}_{[F,C]} \\ \mathbf{M}_{[C,F]} & \mathbf{M}_{[C,C]} \end{bmatrix}$$

is also bDD.

We will consider the linear operator

$$\mathbf{Z}^{(last)} \stackrel{\text{def}}{=} \frac{1}{2} \mathbf{X}^{-1} + \frac{1}{2} \mathbf{X}^{-1} (\mathbf{X} - \mathbf{D} + \mathbf{A}) \mathbf{X}^{-1} (\mathbf{X} - \mathbf{D} + \mathbf{A}) \mathbf{X}^{-1}. \quad (13.23)$$

We define

$$\mathbf{M}^{(last)} = \begin{pmatrix} (\mathbf{Z}^{(last)})^{-1} & \mathbf{M}_{[F,C]} \\ \mathbf{M}_{[F,C]} & \mathbf{M}_{[C,C]} \end{pmatrix}$$

**Lemma 13.11.9.**

$$\mathbf{M} \preceq \mathbf{M}^{(last)} \preceq \left(1 + \frac{2}{\alpha}\right) \mathbf{M}.$$

We defer the proof of Lemma 13.11.9 to Section 13.11.3.

To utilize  $\mathbf{Z}^{(last)}$ , define

$$\mathbf{M}_2^{(last)} \stackrel{\text{def}}{=} \begin{bmatrix} \frac{1}{2} \mathbf{X} & \frac{1}{2} (\mathbf{X} - \mathbf{D} + \mathbf{A}) \mathbf{X}^{-1} \mathbf{M}_{[F,C]} \\ \frac{1}{2} \mathbf{M}_{[C,F]} \mathbf{X}^{-1} (\mathbf{X} - \mathbf{D} + \mathbf{A}) & \mathbf{M}_{[C,C]} - \frac{1}{2} \mathbf{M}_{[C,F]} \mathbf{X}^{-1} \mathbf{M}_{[F,C]} \end{bmatrix} \quad (13.24)$$

and note

$$S_C(\mathbf{M}^{(last)}, F) = S_C(\mathbf{M}_2^{(last)}, F) \quad (13.25)$$

Lemma 13.11.9 tells us that for large enough  $\alpha$ , we can approximate the Schur complement of  $\mathbf{M}$  by approximating the the Schur complement of  $\mathbf{M}_2^{(last)}$ .

The next lemma tells us that  $\mathbf{M}_2$  is bDD and that we can write the matrix as a sum of an explicit bDD matrix and sparse implicitly represented product demand block-Laplacians and bipartite product demand block-Laplacians.

**Lemma 13.11.10.** *Consider a bDD matrix  $\mathbf{M}$ , where  $\mathbf{M}_{[F,F]}$  is  $\alpha$ -bDD for some  $\alpha \geq 4$ , and let  $\mathbf{M}_2^{(last)}$  be the associated matrix defined by equation (13.24). Let  $m$  be the number of nonzero blocks of  $\mathbf{M}$ .*

For  $j \in F$ , we define  $\mathbf{d}^{(j)} \in (\mathbb{C}^{r \times r})^n$

$$\mathbf{d}_{[i]}^{(j)} = \begin{cases} \mathbf{A}_{[i,j]} \mathbf{X}_{[j,j]}^{-1/2} & \text{for } i \in F \\ \mathbf{M}_{[i,j]} \mathbf{X}_{[j,j]}^{-1/2} & \text{for } i \in C \end{cases} \quad (13.26)$$

For  $j \in F$ , we define  $\mathbf{g}^{(j)} \in (\mathbb{C}^{r \times r})^n$

$$\mathbf{g}_{[i]}^{(j)} = \begin{cases} 0 & \text{for } i \in F \\ \mathbf{M}_{[i,j]} \mathbf{X}_{[j,j]}^{-1/2} & \text{for } i \in C \end{cases} \quad (13.27)$$

Then

$$\mathbf{M}_2^{(last)} = \mathbf{Y} + \frac{1}{2} \sum_{j \in F} \mathbf{L}_{G(\mathbf{g}^{(j)})} + \frac{1}{2} \sum_{j \in F} \mathbf{L}_{G(\mathbf{d}^{(j)}, F)}$$

where  $\mathbf{Y}$  is bDD and has  $O(m)$  nonzero blocks, and the total number of nonzero blocks in  $\mathbf{d}^{(i)}$  and  $\mathbf{g}^{(i)}$  for all  $i$  combined is also  $O(m)$ .

$\mathbf{Y}$  as well as  $\mathbf{d}^{(i)}$  and  $\mathbf{g}^{(i)}$  for all  $i$  can be computed in  $O(m)$  time and  $O(\log n)$  depth.

If  $\mathbf{T}$  is a matrix that is only nonzero inside the submatrix  $\mathbf{T}_{[C,C]}$ , then if we apply the transformation of Eq. 13.24, to  $\mathbf{M} + \mathbf{T}$  instead of  $\mathbf{M}$ , we find  $(\mathbf{M} + \mathbf{T})_2^{(last)} = \mathbf{M}_2^{(last)} + \mathbf{T}_{[C,C]}$ , and in particular  $\mathbf{Y}(\mathbf{M} + \mathbf{T}) = \mathbf{Y}(\mathbf{M}) + \mathbf{T}_{[C,C]}$ .

We defer the proof of Lemma 13.11.10 to Section 13.11.3.

*Proof.* (of Lemma 13.11.3) The procedure  $\text{LASTSTEP}(M, F, \alpha, \varepsilon)$  first computes  $\mathbf{Y}$  and  $\mathbf{d}^{(i)}$  and  $\mathbf{g}^{(i)}$  for all  $i$  s.t.

$$\mathbf{M}_2^{(last)} = \mathbf{Y} + \frac{1}{2} \sum_{j \in F} \mathbf{L}_{G(\mathbf{g}^{(j)})} + \frac{1}{2} \sum_{j \in F} \mathbf{L}_{G(\mathbf{d}^{(j)}, F)}, \quad (13.28)$$

where  $\mathbf{M}_2^{(last)} \approx_{2/\alpha} \mathbf{M}$ . Let  $n_{\mathbf{d}^{(i)}}$  and  $n_{\mathbf{g}^{(i)}}$  denote the number of nonzero blocks in each  $\mathbf{d}^{(i)}$  and  $\mathbf{g}^{(i)}$ . By Lemma 13.11.6, we may call  $\text{CLIQUESPARSIFICATION}(\mathbf{g}^{(i)}, \varepsilon/2)$  in  $O(n_{\mathbf{g}^{(i)}} \varepsilon^{-4})$  time to sparsify  $\mathbf{L}_{G(\mathbf{g}^{(i)})}$ , producing a bDD matrix  $\tilde{\mathbf{L}}_{G(\mathbf{g}^{(i)})} \approx_{\varepsilon/2} \mathbf{L}_{G(\mathbf{g}^{(i)})}$  with  $O(n_{\mathbf{g}^{(i)}} \varepsilon^{-4})$  nonzero blocks. By Lemma 13.11.7, we may call  $\text{BIPARTITECLIQUESPARSIFICATION}(\mathbf{d}^{(i)}, F, \varepsilon/2)$  in  $O(n_{\mathbf{d}^{(i)}} \varepsilon^{-4})$  time to sparsify  $\mathbf{L}_{G(\mathbf{d}^{(i)}, F)}$ , producing a bDD matrix  $\tilde{\mathbf{L}}_{G(\mathbf{d}^{(i)}, F)} \approx_{\varepsilon/2} \mathbf{L}_{G(\mathbf{d}^{(i)}, F)}$  with  $O(n_{\mathbf{d}^{(i)}} \varepsilon^{-4})$  nonzero blocks. The total running time of this is  $O(\varepsilon^{-4} \sum_i (n_{\mathbf{d}^{(i)}} + n_{\mathbf{g}^{(i)}})) = O(\varepsilon^{-4} m)$ , and the total number of nonzero blocks in

$$\frac{1}{2} \sum_{j \in F} \tilde{\mathbf{L}}_{G(\mathbf{g}^{(j)})} + \frac{1}{2} \sum_{j \in F} \tilde{\mathbf{L}}_{G(\mathbf{d}^{(j)}, F)},$$

is  $O(\varepsilon^{-4} (\sum_i n_{\mathbf{d}^{(i)}} + n_{\mathbf{g}^{(i)}})) = O(\varepsilon^{-4} m)$ . We define

$$\mathbf{R} = \mathbf{Y} + \frac{1}{2} \sum_{j \in F} \tilde{\mathbf{L}}_{G(\mathbf{g}^{(j)})} + \frac{1}{2} \sum_{j \in F} \tilde{\mathbf{L}}_{G(\mathbf{d}^{(j)}, F)}. \quad (13.29)$$

which we can compute in  $O(\varepsilon^{-4} m)$  time and  $O(\log n)$  depth. We have  $\mathbf{M} \approx_{2/\alpha} \mathbf{M}_2^{(last)}$  and  $\mathbf{M}_2^{(last)} \approx_{\varepsilon/2} \mathbf{R}$ , so that  $\mathbf{M} \approx_{2/\alpha + \varepsilon/2} \mathbf{R}$ . It follows from Fact 13.5.3 that  $\text{Sc}(\mathbf{M}, F) \approx_{2/\alpha + \varepsilon/2} \text{Sc}(\mathbf{R}, F)$ .

Because the sparsifiers computed by  $\text{BIPARTITECLIQUESPARSIFICATION}$  preserve the graph bipartition,  $\mathbf{R}_{[F,F]}$  is block diagonal.

We can use the block diagonal structure of  $\mathbf{R}_{[F,F]}$  quickly compute a sparse approximation to  $\text{Sc}(\mathbf{R}, F)$ .

For  $j \in F$ , we define  $\mathbf{r}^{(j)} \in (\mathbb{C}^{r \times r})^{|C|}$

$$\mathbf{r}_i^{(j)} = \mathbf{R}_{[i,j]} \mathbf{R}_{[j,j]}^{-1/2} \text{ for } i \in C \quad (13.30)$$

Then

$$\text{Sc}(\mathbf{R}, F) = \mathbf{R}_{[C,C]} - \mathbf{R}_{[C,F]} \mathbf{R}_{[F,F]}^{-1} \mathbf{R}_{[F,C]} = \mathbf{S} + \sum_{j \in F} \mathbf{L}_{G(\mathbf{r}^{(j)})}.$$

where

$$\begin{aligned} \mathbf{S} = & \mathbf{R}_{[C,C]} - \text{Diag} \left( \mathbf{I}_r \sum_{i \in C} \sum_{j \in F} \sum_{k \in C \setminus \{i\}} \|\mathbf{R}_{[i,j]} \mathbf{R}_{[j,j]}^{-1/2}\| \|\mathbf{R}_{[k,j]} \mathbf{R}_{[j,j]}^{-1/2}\| \right) \\ & - \text{Diag} \sum_{i \in C} \left( \mathbf{R}_{[i,j]} \mathbf{R}_{[j,j]}^{-1} \mathbf{R}_{[j,i]} \right). \end{aligned} \quad (13.31)$$

Let us define, for each block row  $i \in F$ ,  $\rho_i = \sum_{j \in C} \|\mathbf{R}_{[i,j]}\|$ , and for each block row  $i \in C$ ,  $\rho_i = \sum_{j \in F} \|\mathbf{R}_{[i,j]}\|$ . From  $\mathbf{R}$  being bDD, we then conclude for each  $i \in F$

$$\mathbf{R}_{[i,i]} \succcurlyeq \mathbf{I}_r \rho_i$$

and for each  $i \in C$

$$\mathbf{R}_{[i,i]} \succcurlyeq \mathbf{I}_r \left( \rho_i + \sum_{j \in C} \|\mathbf{R}_{[i,j]}\| \right).$$

With this in mind, we check that each for  $i \in C$  of  $\mathbf{S}$  is bDD.

The diagonal block satisfies

$$\begin{aligned} \mathbf{S}_{[i,i]} & \succcurlyeq \mathbf{I}_r \left( \sum_{j \in C} \|\mathbf{R}_{[i,j]}\| + \rho_i - \sum_{j \in F} \sum_{k \in C} \|\mathbf{R}_{[i,j]} \mathbf{R}_{[j,j]}^{-1/2}\| \|\mathbf{R}_{[k,j]} \mathbf{R}_{[j,j]}^{-1/2}\| \right) \\ & \succcurlyeq \mathbf{I}_r \left( \sum_{j \in C} \|\mathbf{R}_{[i,j]}\| + \rho_i - \sum_{j \in F} \sum_{k \in C} \frac{1}{\rho_j} \|\mathbf{R}_{[i,j]}\| \|\mathbf{R}_{[k,j]}\| \right) \\ & \succcurlyeq \mathbf{I}_r \left( \sum_{j \in C} \|\mathbf{R}_{[i,j]}\| + \rho_i - \sum_{j \in F} \frac{1}{\rho_j} \|\mathbf{R}_{[i,j]}\| \right) \\ & \succcurlyeq \mathbf{I}_r \sum_{j \in C} \|\mathbf{R}_{[i,j]}\| \end{aligned}$$

We can compute  $\mathbf{S}$  and all  $\mathbf{r}_i^{(j)}$  in  $O(m\varepsilon^{-4})$  time and  $O(\log n)$  depth, since this is an upper bound to the number of nonzero blocks in  $\mathbf{R}$ . Let  $n_{\mathbf{r}^{(i)}}$  denote the number of nonzero blocks in  $\mathbf{r}^{(i)}$ . By Lemma 13.11.6, we may call CLIQUE SPARSIFICATION( $\mathbf{r}^{(i)}, \varepsilon/2$ ) in  $O(n_{\mathbf{r}^{(i)}} \varepsilon^{-4})$  time to sparsify  $\mathbf{L}_{G(\mathbf{r}^{(i)})}$ , producing a bDD matrix  $\tilde{\mathbf{L}}_{G(\mathbf{r}^{(i)})} \approx_{\varepsilon/2} \mathbf{L}_{G(\mathbf{r}^{(i)})}$  with  $O(n_{\mathbf{r}^{(i)}} \varepsilon^{-4})$  nonzero blocks.

The total number of nonzero blocks in

$$\mathbf{S} + \sum_i \tilde{\mathbf{L}}_{G(\mathbf{r}^{(i)})}$$

is also  $O(\varepsilon^{-8}m)$ , and  $\mathbf{S} + \sum_i \tilde{\mathbf{L}}_{G(\mathbf{r}^{(i)})} \approx_{\varepsilon/2} \text{Sc}(\mathbf{R}, F)$ . So

$$\tilde{\mathbf{M}}_{SC} = \mathbf{S} + \sum_i \tilde{\mathbf{L}}_{G(\mathbf{r}^{(i)})} \approx_{\varepsilon+2/\alpha} \text{Sc}(\mathbf{M}, F)$$

The total amount of work to compute these  $\tilde{\mathbf{L}}_{G(\mathbf{r}^{(i)})}$  is  $O(\varepsilon^{-4}(\sum_i n_{\mathbf{r}^{(i)}})) = O(\varepsilon^{-8}m)$ . The depth for this computation is  $O(\log n)$ .

Suppose  $\mathbf{T}$  is a matrix that is only nonzero inside the submatrix  $\mathbf{T}_{[C,C]}$ , and  $\mathbf{M} + \mathbf{T}$  is bDD. We can show that  $\text{LASTSTEP}(\mathbf{M} + \mathbf{T}, F, \alpha, \varepsilon) = \text{LASTSTEP}(\mathbf{M}, F, \alpha, \varepsilon) + \mathbf{T}_{[C,C]}$ , by first noting that

this type of property holds for  $\mathbf{Y}$  by Lemma 13.11.10, and from this concluding that similarly if we consider  $\mathbf{R}$  as a function of  $\mathbf{M}$  then  $\mathbf{R}_{[C,C]}(\mathbf{M} + \mathbf{T}) = \mathbf{R}_{[C,C]}(\mathbf{M}) + \mathbf{T}_{[C,C]}$ , and finally considering  $\widetilde{\mathbf{M}}_{SC}$  as a function of  $\mathbf{M}$ , we can then easily show that  $(\mathbf{M} + \mathbf{T})_{SC} = \widetilde{\mathbf{M}}_{SC} + \mathbf{T}_{[C,C]}$ .  $\square$

### ■ 13.11.3 Deferred Proofs from Section 13.11.2

To help us prove Lemma 13.11.9, we first prove the next lemma.

**Lemma 13.11.11.** *If  $\mathbf{M}_{[F,F]} = \mathbf{X} + \mathbf{D} - \mathbf{A}$  be a  $\alpha$ -bDD matrix for some  $\alpha \geq 4$ , then the operator  $\mathbf{Z}^{(last)}$  as defined in Equation 13.23 satisfies:*

$$\mathbf{M}_{[F,F]} \preceq \left( \mathbf{Z}^{(last)} \right)^{-1} \preceq \mathbf{M}_{[F,F]} + \frac{2}{\alpha} (\mathbf{D} - \mathbf{A}).$$

*Proof.* Composing both sides by  $\mathbf{X}^{-1/2}$  and substituting in  $\mathcal{L} = \mathbf{X}^{-1/2} (\mathbf{D} - \mathbf{A}) \mathbf{X}^{-1/2}$  means it suffices to show

$$\mathbf{I} + \mathcal{L} \preceq \left( \frac{1}{2} \mathbf{I} + \frac{1}{2} (\mathbf{I} - \mathcal{L})^2 \right)^{-1} \preceq \mathbf{I} + \mathcal{L} + \frac{2}{\alpha} \mathcal{L}.$$

We can use the fact that  $\mathbf{M}_{[F,F]}$  is  $\alpha$ -strongly diagonally dominant to show  $0 \preceq \mathbf{L} \preceq \frac{2}{\alpha} \mathbf{X}$ , and equivalently  $0 \preceq \mathcal{L} \preceq \frac{2}{\alpha} \mathbf{I}$ , as follows:

Firstly,  $\mathbf{D} - \mathbf{A} \succcurlyeq 0$  as  $\mathbf{D} - \mathbf{A}$  is bDD, similarly,  $\mathbf{D} + \mathbf{A} \succcurlyeq 0$  as  $\mathbf{D} + \mathbf{A}$  is also bDD. From the latter  $\mathbf{D} \succcurlyeq -\mathbf{A}$ , so  $2\mathbf{D} \succcurlyeq \mathbf{D} - \mathbf{A}$ . Finally,  $\mathbf{X} \succcurlyeq \alpha \mathbf{D} \succcurlyeq \frac{\alpha}{2} (\mathbf{D} - \mathbf{A})$ .

As  $\mathcal{L}$  and  $\mathbf{I}$  commute, the spectral theorem means it suffices to show this for any scalar  $0 \leq t \leq \frac{2}{\alpha}$ . Note that

$$\frac{1}{2} + \frac{1}{2} (1 - t)^2 = 1 - t + \frac{1}{2} t^2$$

Taking the difference between the inverse of this and the ‘true’ value of  $1 + t$  gives:

$$\left( 1 - t + \frac{1}{2} t^2 \right)^{-1} - (1 + t) = \frac{1 - (1 + t) (1 - t + \frac{1}{2} t^2)}{1 - t + \frac{1}{2} t^2} = \frac{\frac{1}{2} t^2 (1 - t)}{1 - t + \frac{1}{2} t^2}$$

Incorporating the assumption that  $0 \leq t \leq \frac{2}{\alpha}$  and  $\alpha \geq 4$  gives that the denominator is at least  $1 - \frac{2}{\alpha} \geq \frac{1}{2}$ , and the numerator term can be bounded by  $0 \leq \frac{t^2}{2} (1 - t) \leq \frac{t}{\alpha}$ . Combining these two bounds then gives the result.  $\square$

Lemma 13.11.9 allows us to extend the approximation of  $\mathbf{M}_{[F,F]}$  by the inverse of  $\mathbf{Z}^{(last)}$  to the entire matrix  $\mathbf{M}$ .

*Proof. (of Lemma 13.11.9)* Recall that when a matrix  $\mathbf{T}$  is PSD,

$$\begin{pmatrix} \mathbf{T} & 0 \\ 0 & 0 \end{pmatrix} \succcurlyeq 0. \quad (13.32)$$

The left-hand inequality of our lemma follows immediately from Eq. 13.32 and the left-hand side of the guarantee of Lemma 13.11.11. To prove the right-hand inequality we apply Eq. 13.32 and the right-hand side of the guarantee of Lemma 13.11.11. to conclude

$$\begin{pmatrix} \left( \mathbf{Z}^{(last)} \right)^{-1} & \mathbf{M}_{[F,C]} \\ \mathbf{M}_{[F,C]} & \mathbf{M}_{[C,C]} \end{pmatrix} \preceq \begin{pmatrix} \mathbf{M}_{[F,F]} + \frac{2}{\alpha} \mathbf{L} & \mathbf{M}_{[F,C]} \\ \mathbf{M}_{[F,C]} & \mathbf{M}_{[C,C]} \end{pmatrix} = \mathbf{M} + \frac{2}{\alpha} \begin{pmatrix} \mathbf{L} & 0 \\ 0 & 0 \end{pmatrix}.$$

The matrix  $\begin{bmatrix} \mathbf{X} & \mathbf{M}_{[F,C]} \\ \mathbf{M}_{[C,F]} & \mathbf{M}_{[C,C]} \end{bmatrix}$  is bDD and hence PSD. It follows that  $\begin{pmatrix} \mathbf{D} - \mathbf{A} & 0 \\ 0 & 0 \end{pmatrix} \preceq \mathbf{M}$ , by which we may conclude that  $\mathbf{M} + \frac{2}{\alpha} \begin{pmatrix} \mathbf{D} - \mathbf{A} & 0 \\ 0 & 0 \end{pmatrix} \preceq \mathbf{M} + \frac{2}{\alpha} \mathbf{M}$ .  $\square$

*Proof. (of Lemma 13.11.10)*

Each product demand clique  $\mathbf{L}_{G(g^{(i)})}$  and bipartite product demand clique  $\mathbf{L}_{G(d^{(i),F})}$  is bDD.

We now have to find an expression for  $\mathbf{Y}$  and show that  $\mathbf{Y}$  is bDD. Let us write  $\mathbf{Y}$  in terms of its blocks From

$$\mathbf{Y} = \mathbf{M} - \left( \frac{1}{2} \sum_{j \in F} \mathbf{L}_{G(g^{(i)})} + \frac{1}{2} \sum_{j \in F} \mathbf{L}_{G(d^{(i),F})} \right)$$

it then follows by a simple check that

$$\mathbf{Y}_{[C,F]}^* = \mathbf{Y}_{[F,C]} = \frac{1}{2} (\mathbf{X} - \mathbf{D}) \mathbf{X}^{-1} \mathbf{M}_{[F,C]} = \frac{1}{2} \left( 1 - \frac{1}{\alpha} \right) \mathbf{M}_{[F,C]} \quad (13.33)$$

and that

$$\mathbf{Y}_{[F,F]} = \frac{1}{2} \mathbf{X} - \frac{1}{2} \text{Diag}_{i \in F} \left( \mathbf{I}_r \sum_{j \in F} \sum_{k \in C} \|\mathbf{A}_{[i,j]} \mathbf{X}_{[j,j]}^{-1/2}\| \|\mathbf{M}_{[k,j]} \mathbf{X}_{[j,j]}^{-1/2}\| \right) \quad (13.34)$$

and

$$\begin{aligned} \mathbf{Y}_{[C,C]} &= \mathbf{M}_{[C,C]} - \frac{1}{2} \text{Diag}_{i \in C} (\mathbf{M}_{[i,F]} \mathbf{X}^{-1} \mathbf{M}_{[F,i]}) \\ &\quad - \frac{1}{2} \text{Diag}_{i \in C} \left( \mathbf{I}_r \sum_{j \in F} \sum_{k \in C \setminus \{i\}} \|\mathbf{M}_{[i,j]} \mathbf{X}_{[j,j]}^{-1/2}\| \|\mathbf{M}_{[k,j]} \mathbf{X}_{[j,j]}^{-1/2}\| \right) \\ &\quad - \frac{1}{2} \text{Diag}_{i \in C} \left( \mathbf{I}_r \sum_{j \in F} \sum_{k \in F} \|\mathbf{M}_{[i,j]} \mathbf{X}_{[j,j]}^{-1/2}\| \|\mathbf{A}_{[k,j]} \mathbf{X}_{[j,j]}^{-1/2}\| \right). \end{aligned} \quad (13.35)$$

Next, we check that  $\mathbf{Y}$  is bDD. First, let us define, for each block row  $i \in F$ ,  $\rho_i = \sum_{j \in F} \|\mathbf{A}_{[i,j]}\|$ . From  $\mathbf{M}_{[F,F]}$  being  $\alpha$ -bDD, we then conclude  $\mathbf{M}_{[i,i]} \succeq \mathbf{I}_r (1 + \alpha) \rho_i$ . And from  $\mathbf{M}$  being bDD, we find that  $\sum_{j \in C} \|\mathbf{M}_{[i,j]}\| \leq \alpha \rho_i$ . We now check that each block row  $i \in F$  of  $\mathbf{Y}$  is bDD. The off-diagonal block norm sum is at most  $\frac{1}{2} (1 - \frac{1}{\alpha}) \alpha \rho_i$ .

The diagonal satisfies

$$\begin{aligned} \mathbf{Y}_{[i,i]} &\succeq \mathbf{I}_r \left( \frac{1}{2} \alpha \rho_i - \frac{1}{2} \left( \sum_{j \in F} \sum_{k \in C} \|\mathbf{A}_{[i,j]} \mathbf{X}_{[j,j]}^{-1/2}\| \|\mathbf{M}_{[k,j]} \mathbf{X}_{[j,j]}^{-1/2}\| \right) \right) \\ &\succeq \mathbf{I}_r \left( \frac{1}{2} \alpha \rho_i - \frac{1}{2} \left( \sum_{j \in F} \sum_{k \in C} \frac{1}{\alpha \rho_j} \|\mathbf{A}_{[i,j]}\| \|\mathbf{M}_{[k,j]}\| \right) \right) \\ &\succeq \mathbf{I}_r \left( \frac{1}{2} \alpha \rho_i - \frac{1}{2} \left( \sum_{j \in F} \|\mathbf{A}_{[i,j]}\| \right) \right) = \mathbf{I}_r \left( \frac{1}{2} \left( 1 - \frac{1}{\alpha} \right) \alpha \rho_i \right) \end{aligned}$$

So the block rows with  $i \in F$  are bDD.

Next we check the block rows for  $i \in C$ . Let us define for each  $i \in C$ ,  $\rho_i = \sum_{j \in F} \|\mathbf{M}_{[i,j]}\|$ . Thus for  $\mathbf{Y}$ , the sum of block norms of row  $i$  over columns  $j \in F$

$$\sum_{j \in F} \|\mathbf{Y}_{[i,j]}\| = \frac{1}{2} \left(1 - \frac{1}{\alpha}\right) \sum_{j \in F} \|\mathbf{M}_{[i,j]}\| = \frac{1}{2} \left(1 - \frac{1}{\alpha}\right) \rho_i.$$

For  $\mathbf{Y}$ , the sum of block norms of row  $i$  over columns  $j \in C$  is  $\sum_{j \in C \setminus \{i\}} \|\mathbf{M}_{[i,j]}\|$ . So the total sum of the blocks norms of off-diagonals is

$$\frac{1}{2} \rho_i \left(1 - \frac{1}{\alpha}\right) + \sum_{j \in C \setminus \{i\}} \|\mathbf{M}_{[i,j]}\|$$

The diagonal block satisfies

$$\begin{aligned} \mathbf{Y}_{[i,i]} &\succcurlyeq \mathbf{I}_r \left( \rho_i - \frac{1}{2} \|\mathbf{M}_{[i,F]} \mathbf{X}^{-1} \mathbf{M}_{[F,i]}\| \right. \\ &\quad \left. - \frac{1}{2} \left( \sum_{j \in F} \sum_{k \in C \setminus \{i\}} \|\mathbf{M}_{[i,j]} \mathbf{X}_{[j,j]}^{-1/2}\| \|\mathbf{M}_{[k,j]} \mathbf{X}_{[j,j]}^{-1/2}\| \right) \right. \\ &\quad \left. - \frac{1}{2} \left( \sum_{j \in F} \sum_{k \in F} \|\mathbf{M}_{[i,j]} \mathbf{X}_{[j,j]}^{-1/2}\| \|\mathbf{A}_{[k,j]} \mathbf{X}_{[j,j]}^{-1/2}\| \right) \right) \\ &\succcurlyeq \mathbf{I}_r \left( \rho_i + \sum_{j \in C \setminus \{i\}} \|\mathbf{M}_{[i,j]}\| \right. \\ &\quad \left. - \frac{1}{2} \left( \sum_{j \in F} \sum_{k \in C} \frac{1}{\alpha \rho_j} \|\mathbf{M}_{[i,j]}\| \|\mathbf{M}_{[k,j]}\| \right) \right. \\ &\quad \left. - \frac{1}{2} \left( \sum_{j \in F} \sum_{k \in F} \frac{1}{\alpha \rho_j} \|\mathbf{M}_{[i,j]}\| \|\mathbf{A}_{[k,j]}\| \right) \right) \\ &\succcurlyeq \mathbf{I}_r \left( \rho_i + \sum_{j \in C \setminus \{i\}} \|\mathbf{M}_{[i,j]}\| - \frac{1}{2} \rho_i - \frac{1}{2} \frac{1}{\alpha} \rho_i \right) \\ &= \mathbf{I}_r \left( \frac{1}{2} \rho_i \left(1 - \frac{1}{\alpha}\right) + \sum_{j \in C \setminus \{i\}} \|\mathbf{M}_{[i,j]}\| \right). \end{aligned}$$

For these block rows are also bDD, and hence  $\mathbf{Y}$  is bDD. It is clear from the definitions that  $\mathbf{Y}$  as well as  $\mathbf{d}^{(i)}$  and  $\mathbf{g}^{(i)}$  for all  $i$  can be computed in  $O(m)$  time and  $O(\log n)$  depth.

It is easy to verify that if  $\mathbf{T}$  is a matrix that is only nonzero inside the submatrix  $\mathbf{T}_{[C,C]}$ , then if we apply the transformation of Eq. 13.24, to  $\mathbf{M} + \mathbf{T}$  instead of  $\mathbf{M}$ , we find  $(\mathbf{M} + \mathbf{T})_2^{(last)} = \mathbf{M}_2^{(last)} + \mathbf{T}_{[C,C]}$ , and in particular  $\mathbf{Y}(\mathbf{M} + \mathbf{T}) = \mathbf{Y}(\mathbf{M}) + \mathbf{T}_{[C,C]}$ .  $\square$

#### ■ 13.11.4 Sparsifying Product Demand Block-Laplacians

In this section, we show how to efficiently sparsify product demand block-Laplacians and their bipartite analogs. We prove the following key lemma later in this section that allows us to transfer results on



graph sparsification to sparsifying these product block-Laplacians.

**Lemma 13.11.12.** *Suppose we have two graphs on  $n$  vertices,  $H^{(1)}$  and  $H^{(2)}$  such that  $H^{(1)} \approx_\varepsilon H^{(2)}$ . Given  $\mathbf{d} \in (\mathbb{C}^{r \times r})^n$ , define the bDD matrix  $\mathbf{L}^{(\ell)} \in (\mathbb{C}^{r \times r})^{n \times n}$  for each  $\ell = 1, 2$ , as*

$$\mathbf{L}^{(\ell)}_{[i,j]} = \begin{cases} -\frac{h_{i,j}^{(\ell)}}{\|\mathbf{d}_{[i]}\| \|\mathbf{d}_{[j]}\|} \mathbf{d}_{[i]} \mathbf{d}_{[j]}^* & i \neq j, \\ \mathbf{I}_r \cdot \sum_{k:k \neq i} h_{i,k}^{(\ell)} & \text{otherwise,} \end{cases}$$

where  $h_{i,j}^{(\ell)}$  denotes the weight of the edge  $i, j$  in  $H^{(\ell)}$ . Then,  $\mathbf{L}^{(1)} \approx_\varepsilon \mathbf{L}^{(2)}$ .

We now introduce scalar versions of product block-Laplacian matrices that will be useful.

**Definition 13.11.13.** The product demand graph of a vector  $\mathbf{d} \in (\mathbb{R}_{>0})^n$ ,  $G(\mathbf{d})$ , is a complete weighted graph on  $n$  vertices whose weight between vertices  $i$  and  $j$  is given by  $\mathbf{w}_{ij} = \mathbf{d}_i \mathbf{d}_j$ .

The Laplacian of  $G(\mathbf{d})$ , denoted  $\mathbf{L}_{G(\mathbf{d})}$  is called the product demand Laplacian of  $\mathbf{d}$ .

**Definition 13.11.14.** The bipartite product demand graph of two vectors  $\mathbf{d}^A \in (\mathbb{R}_{>0})^{|A|}$ ,  $\mathbf{d}^B \in (\mathbb{R}_{>0})^{|B|}$ ,  $G(\mathbf{d}^A, \mathbf{d}^B)$ , is a weighted bipartite graph on vertices  $A \cup B$ , whose weight between vertices  $i \in A$  and  $j \in B$  is given by  $\mathbf{w}_{ij} = \mathbf{d}_i^A \mathbf{d}_j^B$ .

The Laplacian of  $G(\mathbf{d}^A, \mathbf{d}^B)$ , denoted  $\mathbf{L}_{G(\mathbf{d}^A, \mathbf{d}^B)}$  is called the bipartite product demand Laplacian of  $(\mathbf{d}^A, \mathbf{d}^B)$ .

In Section 13.14, we give results on efficiently constructing approximations to product demand Laplacians that can be summarized as follows:

*Lemma 13.11.1.* There is a routine  $\text{WEIGHTEDEXPANDER}(\mathbf{d}, \varepsilon)$  such that for any  $\varepsilon > 0$ , and a demand vector  $\mathbf{d} \in (\mathbb{R}_{>0})^n$ ,  $\text{WEIGHTEDEXPANDER}(\mathbf{d}, \varepsilon)$  returns in  $O(n\varepsilon^{-4})$  work and  $O(\log n)$  depth a graph  $H$  with  $O(n\varepsilon^{-4})$  edges such that

$$\mathbf{L}_H \approx_\varepsilon \mathbf{L}_{G(\mathbf{d})}.$$

*Lemma 13.11.2.* There is a routine  $\text{WEIGHTEDBIPARTITEEXPANDER}(\mathbf{d}^A, \mathbf{d}^B, \varepsilon)$  such that for any demand vectors  $\mathbf{d}^A$  and  $\mathbf{d}^B$  of total length  $n$ , and a parameter  $\varepsilon$ , it returns in  $O(n\varepsilon^{-4})$  work and  $O(\log n)$  depth a bipartite graph  $H$  between  $A$  and  $B$  with  $O(n\varepsilon^{-4})$  edges such that

$$\mathbf{L}_H \approx_\varepsilon \mathbf{L}_{G(\mathbf{d}^A, \mathbf{d}^B)}.$$

Finally, we need to define an operation on graphs: Given a graph  $G$ , define  $G^{(K^2)}$  to be the graph obtained by duplicating each vertex in  $G$ , and for each edge  $(i, j)$  in  $G$ , add a  $2 \times 2$  bipartite clique between the two copies of  $i$  and  $j$ .

We now combine the above construction of sparsifiers for product demand graphs with Lemma 13.11.12 to efficiently construct sparse approximations to product demand block-Laplacians.

*Proof. (of Lemma 13.11.6)* The procedure  $\text{CLIQUESPARSIFICATION}(\mathbf{d}, \varepsilon)$  returns the matrix  $\mathbf{L}$  where

$$\mathbf{L}_{[i,j]} = \begin{cases} -\frac{h_{i,j}}{\|\mathbf{d}_{[i]}\| \|\mathbf{d}_{[j]}\|} \mathbf{d}_{[i]} \mathbf{d}_{[j]}^* & i \neq j, \\ \mathbf{I}_r \cdot \sum_{k:k \neq i} h_{i,k} & \text{otherwise.} \end{cases}$$

By construction  $\mathbf{L}_H \approx_\varepsilon \mathbf{L}_{G(\mathbf{w})}$ . Thus, applying Lemma 13.11.12, with  $H^{(1)} = H$  and  $H^{(2)} = G(\mathbf{w})$ , we know that  $\mathbf{L} \approx_\varepsilon \mathbf{L}^{(2)}$ , where  $\mathbf{L}^{(2)}$  is given by

$$\mathbf{L}^{(2)}_{[i,j]} = \begin{cases} -\frac{\mathbf{w}_i \mathbf{w}_j}{\|\mathbf{d}_{[i]}\| \|\mathbf{d}_{[j]}\|} \mathbf{d}_{[i]} \mathbf{d}_{[j]}^* & i \neq j, \\ \mathbf{I}_r \cdot \mathbf{w}_i \sum_{k:k \neq i} \mathbf{w}_k & \text{otherwise.} \end{cases}$$

**Algorithm 42:**  $L = \text{CLIQUESPARSIFICATION}(\mathbf{d}, \varepsilon)$ 


---

```

Initialize  $\mathbf{L} \in (\mathbb{C}^{r \times r})^{n \times n}$  to 0.
Construct  $\mathbf{w}$  where  $w_i = \|\mathbf{d}_{[i]}\|$ .
 $H \leftarrow \text{WEIGHTEDEXPANDER}(\mathbf{w}, \varepsilon)$ .
for edge  $(i, j) \in H$ , with weight  $h_{i,j}$ , do
     $\mathbf{L}_{[i,i]} \leftarrow \mathbf{L}_{[i,i]} + h_{i,j} \mathbf{I}_r$ .
     $\mathbf{L}_{[j,j]} \leftarrow \mathbf{L}_{[j,j]} + h_{i,j} \mathbf{I}_r$ .
     $\mathbf{L}_{[i,j]} \leftarrow \mathbf{L}_{[i,j]} - \frac{h_{i,j}}{w_i w_j} \mathbf{d}_{[i]} \mathbf{d}_{[j]}^*$ .
     $\mathbf{L}_{[j,i]} \leftarrow \mathbf{L}_{[j,i]} - \frac{h_{i,j}}{w_i w_j} \mathbf{d}_{[j]} \mathbf{d}_{[i]}^*$ .
end
Output  $\mathbf{L}$ .

```

---

**Figure 13-8:** Pseudocode for Sparsifying Product Demand Block-Laplacians**Algorithm 43:**  $L = \text{BIPARTITECLIQUESPARSIFICATION}(\mathbf{b}, F, \varepsilon)$ 


---

```

Initialize  $\mathbf{L} \in (\mathbb{C}^{r \times r})^{n \times n}$  to 0.
Construct  $\mathbf{w}$  where  $w_i = \|\mathbf{d}_{[i]}\|$ . Let  $C = V \setminus F$ .
 $H \leftarrow \text{WEIGHTEDBIPARTITEEXPANDER}(\mathbf{w}|_F, \mathbf{w}|_C, \varepsilon)$ .
for edge  $(i, j) \in H$ , with weight  $h_{i,j}$ , do
     $\mathbf{L}_{[i,i]} \leftarrow \mathbf{L}_{[i,i]} + h_{i,j} \mathbf{I}_r$ .
     $\mathbf{L}_{[j,j]} \leftarrow \mathbf{L}_{[j,j]} + h_{i,j} \mathbf{I}_r$ .
     $\mathbf{L}_{[i,j]} \leftarrow \mathbf{L}_{[i,j]} - \frac{h_{i,j}}{w_i w_j} \mathbf{d}_{[i]} \mathbf{d}_{[j]}^*$ .
     $\mathbf{L}_{[j,i]} \leftarrow \mathbf{L}_{[j,i]} - \frac{h_{i,j}}{w_i w_j} \mathbf{d}_{[j]} \mathbf{d}_{[i]}^*$ .
end
Output  $\mathbf{L}$ .

```

---

**Figure 13-9:** Pseudocode for Sparsifying Bipartite Product Demand Block-Laplacians

Since  $w_i = \|\mathbf{d}_{[i]}\|$ , we have  $\mathbf{L}^{(2)} = \mathbf{L}_G(\mathbf{d})$ , proving our claim.  $\square$

*Proof. (of Lemma 13.11.7)* The procedure  $\text{BIPARTITECLIQUESPARSIFICATION}(\mathbf{d}, F, \varepsilon)$  returns the matrix  $\mathbf{L}$  where

$$\mathbf{L}_{[i,j]} = \begin{cases} -\frac{h_{i,j}}{\|\mathbf{d}_{[i]}\| \|\mathbf{d}_{[j]}\|} \mathbf{d}_{[i]} \mathbf{d}_{[j]}^* & i \neq j, \\ \mathbf{I}_r \cdot \sum_{k:k \neq i} h_{i,k} & \text{otherwise.} \end{cases}$$

Since  $H$  returned by  $\text{WEIGHTEDBIPARTITEEXPANDER}$  is guaranteed to be bipartite using Lemma 13.11.2, we obtain that  $(\mathbf{L}_H)_{[F,F]}$ , and  $(\mathbf{L}_H)_{[C,C]}$  are block-diagonal.

By construction  $\mathbf{L}_H \approx_\varepsilon \mathbf{L}_{G(\mathbf{w}|_F, \mathbf{w}|_C)}$ . Thus, applying Lemma 13.11.12, with  $H^{(1)} = H$  and  $H^{(2)} = G(\mathbf{w}|_F, \mathbf{w}|_C)$ , we know that  $\mathbf{L} \approx_\varepsilon \mathbf{L}^{(2)}$ , where  $\mathbf{L}^{(2)}$  is given by

$$\mathbf{L}^{(2)}_{[i,j]} = \begin{cases} \mathbf{I}_r \cdot w_i \sum_{k \in C} w_k & i = j, \text{ and } i \in F, \\ \mathbf{I}_r \cdot w_i \sum_{k \in F} w_k & i = j, \text{ and } i \in C, \\ 0 & i \neq j, \text{ and } (i, j \in F \text{ or } i, j \in C), \\ -\frac{w_i w_j}{\|\mathbf{d}_{[i]}\| \|\mathbf{d}_{[j]}\|} \mathbf{d}_{[i]} \mathbf{d}_{[j]}^* & i \neq j, \text{ otherwise.} \end{cases}$$

Since  $\mathbf{w}_i = \|\mathbf{d}_{[i]}\|$ , we have  $\mathbf{L}^{(2)} = \mathbf{L}_G(\mathbf{d}|_F, \mathbf{d}|_C)$ , proving our claim.  $\square$

Finally, we give a proof of Lemma 13.11.12.

*Proof.* (of Lemma 13.11.12) Using Fact 13.5.2, we can write each  $\mathbf{d}_{[i]}$  as  $\frac{1}{2}\|\mathbf{d}_{[i]}\|(\mathbf{Q}_i^{(1)} + \mathbf{Q}_i^{(2)})$ , where  $\mathbf{Q}_i^{(1)}(\mathbf{Q}_i^{(1)})^* = \mathbf{Q}_i^{(2)}(\mathbf{Q}_i^{(2)})^* = \mathbf{I}_r$ . Construct the matrix  $\mathbf{F} \in (\mathbb{C}^{r \times r})^{2n \times 2n}$  as follows:

$$\mathbf{F} = \begin{bmatrix} \mathbf{Q}_1^{(1)} & \mathbf{Q}_1^{(2)} & 0 & 0 & \cdots \\ 0 & 0 & \mathbf{Q}_2^{(1)} & \mathbf{Q}_2^{(2)} & \cdots \\ \vdots & \vdots & & & \ddots \end{bmatrix}.$$

Now, observe that for each  $\ell = 1, 2$ ,

$$\mathbf{L}^{(\ell)}_{[i,i]} = \frac{1}{2} \left( \sum_{k:k \neq i} h_{i,k}^{(\ell)} \right) \cdot \left( \mathbf{Q}_i^{(1)}(\mathbf{Q}_i^{(1)})^* + \mathbf{Q}_i^{(2)}(\mathbf{Q}_i^{(2)})^* \right) = \frac{1}{4} \left( \mathbf{F}(\mathbf{L}_{(H^{(\ell)})^{(K2)}} \otimes \mathbf{I}_r)\mathbf{F}^* \right)_{[i,i]},$$

where  $\mathbf{L}_{(H^{(\ell)})^{(K2)}}$  is the Laplacian of the graph  $(H^{(\ell)})^{(K2)}$  defined above, and  $\otimes$  tensor product. Also, for  $i \neq j$ ,

$$\mathbf{L}^{(\ell)}_{[i,j]} = -\frac{h_{i,j}^{(\ell)}}{\mathbf{w}_i \mathbf{w}_j} \mathbf{d}_{[i]} \mathbf{d}_{[j]}^* = -\frac{h_{i,j}^{(\ell)}}{4} \cdot \left( \mathbf{Q}_i^{(1)} + \mathbf{Q}_i^{(2)} \right) \left( \mathbf{Q}_j^{(1)} + \mathbf{Q}_j^{(2)} \right)^* = \frac{1}{4} \left( \mathbf{F}(\mathbf{L}_{(H^{(\ell)})^{(K2)}} \otimes \mathbf{I}_r)\mathbf{F}^* \right)_{[i,j]},$$

Thus,

$$\mathbf{L}^{(\ell)} = \frac{1}{4} \mathbf{F}(\mathbf{L}_{(H^{(\ell)})^{(K2)}} \otimes \mathbf{I}_r)\mathbf{F}^*.$$

By assumption, we have  $\mathbf{L}_{H^{(1)}} \approx_\varepsilon \mathbf{L}_{H^{(2)}}$ . Using Lemma 13.11.15, we get that  $\mathbf{L}_{H^{(1)}(K2)} \approx_\varepsilon \mathbf{L}_{H^{(2)}(K2)}$ . This implies

$$\mathbf{L}_{(H^{(1)})^{(K2)}} \otimes \mathbf{I}_r \approx_\varepsilon \mathbf{L}_{(H^{(2)})^{(K2)}} \otimes \mathbf{I}_r,$$

and thus,  $\mathbf{L}^{(1)} \approx_\varepsilon \mathbf{L}^{(2)}$ .  $\square$

### ■ 13.11.5 Constructing sparsifiers for lifts of graphs

Given a graph  $G$ , define  $G^{(K2)}$  to be the graph obtained by duplicating each vertex in  $G$ , and for each edge  $(i, j)$  in  $G$ , add a  $2 \times 2$  bipartite clique between the two copies of  $i$  and  $j$ .

**Lemma 13.11.15.** *If  $H$  is a sparsifier for  $G$ , i.e.,  $L_G \approx_\varepsilon L_H$ , then  $H^{(K2)}$  is a sparsifier for  $G^{(K2)}$ , i.e.  $L_{G^{(K2)}} \approx_\varepsilon L_{H^{(K2)}}$ .*

*Proof.* Since  $L_G \approx_\varepsilon L_H$ , we have  $e_i^\top L_G e_i \approx_\varepsilon e_i^\top L_H e_i$ . Thus, if  $D_G$  denotes the diagonal matrix of degrees of  $G$ , we have  $D_G \approx_\varepsilon D_H$ . The Laplacian for  $G^{(K2)}$  is

$$L_{G^{(K2)}} = L_G \otimes \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + D_G \otimes \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}.$$

Similarly,

$$L_{H^{(K2)}} = L_H \otimes \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + D_H \otimes \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}.$$

Observe that we can write

$$L_G \otimes \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 & 0 & \dots \\ 0 & 0 & 1 & 1 & \dots \\ \vdots & \vdots & & \ddots & \end{bmatrix} L_G \begin{bmatrix} 1 & 1 & 0 & 0 & \dots \\ 0 & 0 & 1 & 1 & \dots \\ \vdots & \vdots & & \ddots & \end{bmatrix}^\top.$$

Since  $L_G \approx_\varepsilon L_H$ , this implies

$$L_G \otimes \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \approx_\varepsilon L_H \otimes \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}.$$

Similarly, we get,

$$D_G \otimes \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \approx_\varepsilon D_H \otimes \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}.$$

Adding the above two, we get  $L_{G(\kappa_2)} \approx_\varepsilon L_{H(\kappa_2)}$ .  $\square$

## ■ 13.12 Estimating Leverage Scores by Undersampling

We will control the densities of all intermediate bDD matrices using the uniform sampling technique introduced in Chapter 3. It relies on sampling columns<sup>2</sup> of matrices by upper bounds of their true leverage scores,

$$\tau_i(\mathbf{A}) = \mathbf{a}_i^T (\mathbf{A}\mathbf{A}^*)^{-1} \mathbf{a}_i.$$

These upper bounds are measured w.r.t. a different matrix, giving generalized leverage scores of the form:

$$\tau_i^{\mathbf{B}}(\mathbf{A}) = \begin{cases} \mathbf{a}_i^* (\mathbf{B}\mathbf{B}^*)^{-1} \mathbf{a}_i & \text{if } \mathbf{a}_i \perp \ker(\mathbf{B}), \\ 1 & \text{otherwise.} \end{cases}$$

We introduced *unitary edge-vertex transfer matrices* in Definition 13.6.4. Sparsifying bDD matrices can be transformed into the more general setting described above via a unitary edge-vertex transfer matrix, which is analogous to the edge-vertex incidence matrix.

Lemma 13.6.5 proves that every bDD matrix  $\mathbf{M} \in (\mathbb{C}^{r \times r})^{n \times n}$  with  $m$  nonzero off-diagonal blocks can be written as  $\mathbf{X} + \mathbf{B}\mathbf{B}^*$  where  $\mathbf{B} \in (\mathbb{C}^{r \times r})^{n \times 2m}$  is a unitary edge-vertex transfer matrix and  $\mathbf{X}$  is a block diagonal PSD matrix. Additionally, for every block diagonal matrix  $\mathbf{Y}$  s.t.  $\mathbf{M} - \mathbf{Y}$  is bDD, we have  $\mathbf{X} \succcurlyeq \mathbf{Y}$ . This decomposition can be found in  $O(m)$  time and  $O(\log n)$  depth. We rely on this  $\mathbf{X}$  to detect some cases of high leverage scores as samples of  $\mathbf{B}$  may have lower rank.

We will reduce the number of nonzero blocks in  $\mathbf{M}$  by sampling columns blocks from this matrix. This is more restrictive than sampling individual columns. Nonetheless, it can be checked via matrix concentration bounds [3, 251] that it suffices to sample the block by analogs of leverage scores:

$$\tau_{[i]} = \text{tr} \left( \mathbf{B}_{[i]}^* (\mathbf{X} + \mathbf{B}\mathbf{B}^*)^{-1} \mathbf{B}_{[i]} \right) \quad (13.36)$$

As in Chapter 3, we recursively estimate upper bounds for these scores, leading to the pseudocode given in Figure 13-10.

**Lemma 13.12.1.** *Given a positive definite bDD matrix  $\mathbf{M} \in (\mathbb{C}^{r \times r})^{n \times n}$  with  $m$  nonzero blocks.*

*Assume that for any positive definite bDD matrix  $\mathbf{M}' \in (\mathbb{C}^{r \times r})^{n \times n}$  with  $\hat{m}$  nonzero blocks, we can find an implicit representation of a matrix  $\mathbf{W}$  such that  $\mathbf{W} \approx_1 (\mathbf{M}')^{-1}$  in depth  $d_{\text{con}}(\hat{m}, n)$  and work  $w_{\text{con}}(\hat{m}, n)$  and for any vector  $\mathbf{b}$ , we can evaluate  $\mathbf{W}\mathbf{b}$  in depth  $d_{\text{eval}}(\hat{m}, n)$  and work  $w_{\text{eval}}(\hat{m}, n)$ .*

<sup>2</sup>In Chapter 3, we sample rows of matrices. Here, we sample columns instead in order to use a more natural set of notations.

---

**Algorithm 44:**  $\widetilde{\mathbf{M}} = \text{SPARSIFY}(\mathbf{M}, \varepsilon, K)$

---

Write  $\mathbf{M}$  into the form  $\mathbf{X} + \mathbf{B}\mathbf{B}^*$  where  $\mathbf{X}$  is a block diagonal positive definite matrix  $\mathbf{X}$  and  $\mathbf{B}$  is an unitary edge-vertex transfer matrix.

Uniformly sample  $\frac{m}{K}$  column blocks of  $\mathbf{B}$  to form  $\mathbf{C}$ .

Find an implicit representation of  $(\mathbf{X} + \mathbf{C}\mathbf{C}^*)^{-1}$ ,  $\mathbf{W}$ .

Approximate  $\text{tr}\left(\mathbf{B}_{[i]}^*(\mathbf{X} + \mathbf{C}\mathbf{C}^*)^{-1}\mathbf{B}_{[i]}\right)$  using Johnson-Lindenstrauss projections.

This involves applying  $\mathbf{W}$  to  $O(\log n)$  random vectors.

Use these estimates to sample blocks of  $\mathbf{B}$  to form  $\widetilde{\mathbf{B}}$ .

Return  $\mathbf{X} + \widetilde{\mathbf{B}}\widetilde{\mathbf{B}}^*$ .

---

**Figure 13-10:** Sparsification via. Uniform Sampling

For any  $K \geq 1$ ,  $1 \geq \varepsilon > 0$ , the algorithm  $\text{SPARSIFY}(\mathbf{M}, \varepsilon, K)$  outputs an explicit positive definite bDD matrix  $\widetilde{\mathbf{M}}$  with  $O(Kn \log n / \varepsilon^2)$  nonzero blocks and  $\widetilde{\mathbf{M}} \approx_{\varepsilon} \mathbf{M}$ .

Also, this algorithm runs in

$$d_{\text{con}}\left(\frac{m}{K}, n\right) + O\left(d_{\text{eval}}\left(\frac{m}{K}, n\right) + \log n\right)$$

depth and

$$w_{\text{con}}\left(\frac{m}{K}, n\right) + O\left(w_{\text{eval}}\left(\frac{m}{K}, n\right) \log n + m \log n\right)$$

work.

The guarantees of this process can be obtained from the main result from Chapter 3:

**Theorem 13.12.2.** Let  $\mathbf{A}$  be a  $n$  by  $m$  matrix, and  $\alpha \in (0, 1]$  a density parameter. Consider the matrix  $\mathbf{A}'$  consisting of a random  $\alpha$  fraction of the columns of  $\mathbf{A}$ . Then, with high probability we have  $\mathbf{u}_i = \tau_i^{\mathbf{A}'}(\mathbf{A})$  is a vector of leverage score overestimates, i.e.  $\tau_i(\mathbf{A}) \leq \mathbf{u}_i$ , and

$$\sum_{i=1}^m \mathbf{u}_i \leq \frac{3n}{\alpha}.$$

*Proof.* (Sketch of Lemma 13.12.1) Instead of sampling  $m/K$  column blocks, consider sampling  $m/K$  individual columns to form  $\widehat{\mathbf{C}}$ . Theorem 13.12.2 gives that the total leverage scores of all the columns of  $\mathbf{B}$  w.r.t.  $\widehat{\mathbf{C}}$  is at most  $O(nKr^2)$ .

As  $\mathbf{C}$  is formed by taking blocks instead of columns, we have

$$\mathbf{X} + \mathbf{C}\mathbf{C}^* \preceq \mathbf{X} + \widehat{\mathbf{C}}\widehat{\mathbf{C}}^*,$$

so the individual leverage scores computed w.r.t.  $\mathbf{C}$  (and in turn  $\mathbf{W}$ ) sums up to less than the ones computed w.r.t.  $\widehat{\mathbf{C}}$ .

Let us use  $\mathbf{b}_{[i],j} \in \mathbb{C}^n$  to denote the  $j^{\text{th}}$  column of the  $\mathbf{B}_{[i]}$  block column. Note that

$$\mathbf{B}_{[i]}\mathbf{B}_{[i]}^* = \sum_j \mathbf{b}_{[i],j}\mathbf{b}_{[i],j}^*.$$

Define  $\tau_{[i],j}^{\mathbf{X}+\mathbf{C}\mathbf{C}^*} = \text{tr}\left(\mathbf{b}_{[i],j}^*(\mathbf{X} + \mathbf{C}\mathbf{C}^*)^{-1}\mathbf{b}_{[i],j}\right)$ .

We then have

$$\text{tr}\left(\mathbf{B}_{[i]}^*(\mathbf{X} + \mathbf{C}\mathbf{C}^*)^{-1}\mathbf{B}_{[i]}\right) = \sum_j \tau_{[i],j}^{\mathbf{X}+\mathbf{C}\mathbf{C}^*},$$

so the total sum of the estimates we obtain is at most  $O(nr^2K)$ , which gives  $\tilde{\mathbf{B}}$  with  $O(Kn \log n \varepsilon^{-2})$  blocks.

To approximate  $\sum_j \tau_{[i],j}^{\mathbf{X}+\mathbf{C}\mathbf{C}^*}$ , we apply a standard technique given by [242, 23, 170]. The rough idea is to write

$$\tau_{[i],j}^{\mathbf{X}+\mathbf{C}\mathbf{C}^*} = \|\mathbf{C}^*(\mathbf{X} + \mathbf{C}\mathbf{C}^*)^{-1}\mathbf{b}_{[i],j}\|^2 + \|\sqrt{\mathbf{X}}(\mathbf{X} + \mathbf{C}\mathbf{C}^*)^{-1}\mathbf{b}_{[i],j}\|^2$$

By Johnson-Lindenstrauss Lemma, for a random Gaussian matrix  $\mathbf{G}$  with  $\Theta(\log(n))$  rows, we know that

$$\tau_{[i],j}^{\mathbf{X}+\mathbf{C}\mathbf{C}^*} \approx_2 \|\mathbf{G}\mathbf{C}^*(\mathbf{X} + \mathbf{C}\mathbf{C}^*)^{-1}\mathbf{b}_{[i],j}\|^2 + \|\mathbf{G}\sqrt{\mathbf{X}}(\mathbf{X} + \mathbf{C}\mathbf{C}^*)^{-1}\mathbf{b}_{[i],j}\|^2$$

for high probability. Since  $\mathbf{G}$  has  $O(\log(n))$  rows, this can be approximated by applying the approximate inverse  $\mathbf{W}$  to  $O(\log(n))$  vectors. Hence, step 3 runs in  $O(d_{\text{eval}}(\frac{m}{K}, n) + \log n)$  depth and  $O(w_{\text{eval}}(\frac{m}{K}, n) \log n + m \log n)$  work.  $\square$

We remark that the extra factor of  $r^2$  can likely be improved by modifying the proof of Theorem 13.12.2 to work with blocks. We omit this improvement for simplicity, especially since we're already treating  $r$  as a constant.

## ■ 13.13 Recursive Construction of Schur Complement Chains

We now give the full details for the recursive algorithm that proves the running times as stated in Theorem 13.1.2:

**Theorem 13.1.2** (Sparsified Multigrid). *There is an algorithm that, when given a bDD matrix  $\mathbf{M}$  with  $n$  block rows and  $m$  nonzero blocks, produces a solver for  $\mathbf{M}$  in  $O(m \log n + n \log^{2+o(1)} n)$  work and  $O(n^{o(1)})$  depth so that the solver finds  $\varepsilon$ -approximate solutions to systems of equations in  $\mathbf{M}$  in  $O((m + n \log^{1+o(1)} n) \log(1/\varepsilon))$  work and  $O(\log^2 n \log \log n \log(1/\varepsilon))$  depth.*

This argument can be viewed as a more sophisticated version of the one in Section 13.3.6. The main idea is to invoke routines for reducing edges and vertices in a slightly unbalanced recursion, where several steps of vertex reductions take place before a single edge reduction. The routines that we will call are:

1. BDDSUBSET given by Lemma 13.3.2 proven in Section 13.8 for finding a large set of  $\alpha$ -bDD subset.
2. APPROXSCHUR given by Lemma 13.11.1 proven in Section 13.11 that gives sparse approximations to Schur complements.
3. SPARSIFY given by Lemma 13.12.1 proven in Section 13.12 that allows us to sparsify a bDD matrix by recursing on a uniform subsample of its non-zero blocks.

An additional level of complication comes from the fact that the approximation guarantees of our constructions rely on gradually decreasing errors down the Schur complement chain. This means that the density increases faster and larger reduction factors  $K$  are required as the iteration goes on. Pseudocode of our algorithm is given in Figure 13-11.

Note that since  $\mathbf{M}^{(0)}$  may be initially dense, we first make a recursive call to SPARSIFY before computing the approximate Schur complements.

In our analysis, we will use  $n^{(i)}$  to denote the number of non-zero column/row blocks in  $\mathbf{M}^{(i)}$ , and  $m^{(i)}$  to denote the number of non-zero blocks in  $\mathbf{M}^{(i)}$ . These are analogous to dimension and

---

**Algorithm 45:**  $(\mathbf{M}^{(1)}, \mathbf{M}^{(2)}, \dots; F_1, F_2, \dots) = \text{RECURSIVECONSTRUCT}(\mathbf{M}^{(0)})$

---

Set  $i \leftarrow 0$ , and  $k$  and  $c$  are parameters to be decided.

**while**  $\mathbf{M}^{(i)}$  has more than  $\Theta(1)$  blocks **do**

If  $i \bmod k = 0$ , Then  $\mathbf{M}^{(i)} \leftarrow \text{SPARSIFY}(\mathbf{M}^{(i)}, (i+8)^{-2}, 2^{2ck \log^2(i+1)})$ .

$F_{i+1} \leftarrow \text{BDDSUBSET}(\mathbf{M}^{(i)}, 4)$ .

$\mathbf{M}^{(i+1)} \leftarrow \text{APPROXSCHUR}(\mathbf{M}^{(i)}, F_{i+1}, 4, (i+8)^{-2})$ .

$i \leftarrow i + 1$ .

**end**

---

**Figure 13-11:** Pseudocode for Recursively Constructing Schur Complement Chains, the recursive calls happen through SPARSIFY, which constructs its own Schur Complement Chains to perform the sparsification

number of non-zeros in the matrix. It is also useful to refer to the steps between calls to SPARSIFY as phases. Specifically, phase  $j$  consists of iterations  $(j-1)k$  to  $jk-1$ . We will also use  $K_j$  to denote the reduction factor used to perform the sparsification at the start of phase  $j$ , aka.

$$K_j = 2^{2ck \log^2((j-1)k+1)}.$$

A further technicality is that Lemma 13.12.1 require a strictly positive definite block-diagonal part to facilitate the detection of vectors in the null space of the sample. We do so by padding  $\mathbf{M}$  with a small copy of the identity, and will check below that this copy stays throughout the course of this algorithm.

**Lemma 13.13.1.** *For any bDD matrix  $\mathbf{M}^{(0)}$  expressible as  $\xi \mathbf{I} + \mathbf{B}^{(0)}(\mathbf{B}^{(0)})^*$  for some  $\xi > 0$  and unitary edge-vertex transfer matrix  $\mathbf{B}$ , all intermediate matrices  $\mathbf{M}'$  generated during the algorithm  $\text{RECURSIVECONSTRUCT}(\mathbf{M}^{(0)})$  can be expressed as  $\xi \mathbf{I} + \mathbf{B}'(\mathbf{B}')^*$ .*

*Proof.* There are two mechanisms by which this recursive algorithm generates new matrices: through uniform sampling within SPARSIFY and via APPROXSCHUR. We will show inductively down the algorithmic calls that all matrices satisfy this property.

Lemma 13.12.1 gives that this  $\mathbf{X}$  is preserved in the sample.

For calls to  $\text{APPROXSCHUR}(\mathbf{M}, C)$ , the inductive hypothesis gives that  $\mathbf{M}_{[C,C]}$  is expressible as  $\widehat{\mathbf{M}} + \xi \mathbf{I}_C$ . The last condition in Lemma 13.3.4 then gives that the result also has a  $\xi \mathbf{I}_C$  part, which gives the inductive hypothesis for it as well.  $\square$

As the interaction between this padding and our recursion is minimal, we will simply state the input conditions as  $\mathbf{M} = \xi \mathbf{I} + \mathbf{B}\mathbf{B}^*$  for some  $\xi > 0$ . We will first bound the sizes of the matrices in the Schur complement chain produced by RECURSIVECONSTRUCT.

**Lemma 13.13.2.** *For any bDD matrix  $\mathbf{M}^{(0)} = \xi \mathbf{I} + \mathbf{B}^{(0)}(\mathbf{B}^{(0)})^*$  for some  $\xi > 0$ , the algorithm  $\text{RECURSIVECONSTRUCT}(\mathbf{M}^{(0)})$  returns a Schur complement chain  $(\mathbf{M}^{(1)}, \mathbf{M}^{(2)}, \dots; F_1, F_2, \dots)$  such that:*

1.  $n^{(i)} \leq \beta^i n^{(0)}$  for some absolute constant  $\beta < 1$ .
2. The number of non-zero blocks in any iteration  $i$  of phase  $j$  is at most  $2^{3ck \log^2(jk)} n^{((j-1)k)} \log n^{((j-1)k)}$ .

*Proof.* Lemma 13.3.2 and the choice of  $\alpha = 4$  ensures  $|F_k| = \Omega(n^{(i)})$ , which means there is constant  $\beta < 1$  such that  $n^{(i)} \leq \beta^k n$ . Furthermore, we do not increase vertex count at any point in this recursion, and all recursive calls are made to smaller graphs. Therefore, the recursion terminates.

Lemma 13.11.1 shows that after computing each approximate Schur complement,

$$\begin{aligned} m^{(i+1)} &= O(m^{(i)}(i^2 \log(i+8))^{O(\log(i+8))}) \\ &\leq 2^{O(\log^2(i+1))} m^{(i)}. \end{aligned}$$

Hence, by picking  $c$  appropriately we can guarantee that the density increases by a factor of most  $2^{c \log^2(i+1)}$  during each iteration. This size increase is controlled by calls to SPARSIFY. Specifically, Lemma 13.12.1 gives that at the start of phase  $j$  we have:

$$\frac{m^{((j-1)k)}}{n^{((j-1)k)} \log n^{((j-1)k)}} \leq 2^{2ck \log^2((j-1)k)}.$$

Then, since we go at most  $k$  steps without calling SPARSIFY, this increase in density can be bounded by:

$$2^{2ck \log^2((j-1)k)} 2^{2c \log^2(jk-1)} \dots 2^{2c \log^2((j-1)k+1)} \leq 2^{3ck \log^2(jk)}.$$

□

**Lemma 13.13.3.** *For any bDD matrix  $\mathbf{M}^{(0)} = \xi \mathbf{I} + \mathbf{B}^{(0)}(\mathbf{B}^{(0)})^*$  for some  $\xi > 0$  that has  $n$  blocks, the algorithm RECURSIVECONSTRUCT( $\mathbf{M}^{(0)}$ ) returns a Schur Complement Chain  $(\mathbf{M}^{(1)}, \dots; F_1, \dots)$  whose corresponding the linear operator  $\mathbf{W}$  satisfies*

$$\mathbf{W} \approx_{O(1)} \left( \mathbf{M}^{(0)} \right)^{-1}.$$

*Also, we can evaluate  $\mathbf{W}\mathbf{b}$  in  $O(\log^2 n \log \log n)$  depth and  $2^{O(k \log^2 k)} n \log n$  work for any vector  $\mathbf{b}$ .*

*Proof.* We first bound the quality of approximation between  $\mathbf{W}$  and  $\left( \mathbf{M}^{(0)} \right)^{-1}$ . The approximate Schur complement  $\mathbf{M}^{(i+1)}$  was constructed so that  $\mathbf{M}^{(i+1)} \approx_{(i+8)^{-2}} \text{Sc} \left( \mathbf{M}^{(i)}, F_i \right)$ . The other source of error, SPARSIFY, is called only for some  $i$ . In those iterations, Lemma 13.12.1 guarantee that  $\mathbf{M}^{(i)}$  changes only by  $(i+8)^{-2}$  factor. This means overall we have  $\mathbf{M}^{(i+1)} \approx_{2(i+8)^{-2}} \text{Sc} \left( \mathbf{M}^{(i)}, F_i \right)$ . By Lemma 13.3.1, we have:

$$\mathbf{W} \approx_{1/2+4 \sum_i (i+8)^{-2}} \left( \mathbf{M}^{(0)} \right)^{-1},$$

and it can be checked that  $\sum_i (i+8)^{-2}$  is a constant.

The cost of APPLYCHAIN is dominated by the sequence of calls to JACOBI. As each  $F_i$  is chosen to be 4-bDD, the number of iterations required is  $O(\log(\varepsilon_i)) = O(\log i)$ . As matrix-vector multiplications take  $O(\log n)$  depth, the total depth can be bounded by.

$$O \left( \sum_i \log i \log n^{(i)} \right) = O(\log^2 n \log \log n)$$

The total work of these steps depend on the number of non-zero blocks,  $m^{(i)}$ . Substituting the bounds from Lemma 13.13.2 into JACOBI gives a total of:

$$O \left( \sum_i 2^{3ck \log^2(i+k)} n^{(i)} \log n^{(i)} \log i \right) \leq 2^{O(k \log^2 k)} n \log n$$



where the inequality follows from the fact that the  $n^{(i)}$ s are geometrically decreasing.  $\square$

This allows us to view the additional overhead of SPARSIFY as a black box, and analyze the total cost incurred by the non-recursive parts of RECURSIVECONSTRUCT during each phase.

**Lemma 13.13.4.** *The total cost of RECURSIVECONSTRUCT during phase  $j$ , including the linear system solves made by SPARSIFY at iteration  $(j-1)k$  (but not its recursive invocation to RECURSIVECONSTRUCT) is*

$$2^{O(k \log^2(jk))} n^{((j-1)k)} \log^2 n^{((j-1)k)}$$

in work and

$$O\left(k \log(jk) \log^2 n^{((j-1)k)}\right) + O\left(\log^2 n^{((j-1)k)} \log \log n^{((j-1)k)}\right)$$

in depth.

*Proof.* Let  $m^{(i)}$  and  $n^{(i)}$  be the number of non zeros and variables in  $\mathbf{M}^{(i)}$  before the SPARSIFY call if there is. Lemmas 13.3.2 and 13.11.1 show that the  $i^{\text{th}}$  iteration takes  $O(m^{(i)} + m^{(i+1)})$  work and  $O(\log i \log n^{(i)})$  depth. By Lemma 13.13.2 the cost during these iterations excluding the call to SPARSIFY is:

$$\begin{aligned} \sum_{i=(j-1)k}^{jk-1} O\left(m^{(i)} + m^{(i+1)}\right) &\leq \sum_{i=(j-1)k}^{jk-1} 2^{O(k \log^2 i)} n^{(i)} \log n^{(i)} \\ &\leq 2^{O(k \log^2(jk))} n^{((j-1)k)} \log n^{((j-1)k)} \end{aligned}$$

work and

$$\sum_{i=(j-1)k}^{jk-1} O\left(\log i \log n^{(i)}\right) \leq O(k \log((j-1)k) \log n^{((j-1)k)})$$

depth.

We now consider the call to SPARSIFY made during iteration  $(j-1)k$ . Access to a fast solver for the sampled bDD matrix is obtained via recursive calls to RECURSIVECONSTRUCT. The guarantees of the chain given by Lemma 13.13.3 above means each solve takes depth

$$d_{eval} = \log^2 n^{((j-1)k)} \log \log n^{((j-1)k)},$$

and work

$$w_{eval} = n^{((j-1)k)} \log n^{((j-1)k)}$$

Incorporating these parameters into Lemma 13.12.1 allows us to bound the overhead from these solves by

$$2^{O(k \log^2(jk))} n^{((j-1)k)} \log^2 n^{((j-1)k)}$$

work and

$$O\left(k \log(jk) \log n^{(jk)}\right) + O\left(\log^2 n^{((j-1)k)} \log \log \left(n^{((j-1)k)}\right)\right)$$

depth.  $\square$

Note that at the end of the  $j^{\text{th}}$  phase, the time required to construct an extra Schur complement chain for the SPARSIFY call is less than the remaining cost after the  $j^{\text{th}}$  phase. This is the reason why we use  $2^{2ck \log^2(i+1)}$  as the reduction factor for the SPARSIFY call. The following theorem takes account for the recursive call and show the total running time for the algorithm.

**Lemma 13.13.5.** *With high probability,  $\text{RECURSIVECONSTRUCT}(\mathbf{M}^{(0)})$  takes  $2^{O(\log n/k)}$  depth and  $m \log n + 2^{O(k \log^2 k)} n \log^2 n$  work.*

*Proof.* Lemma 13.13.2 shows that the call to SPARSIFY at the start of each phase  $j$  requires the construction of an extra Schur complement chain on a matrix with  $n^{((j-1)k)}$  row/column blocks and at most  $2^{ck \log^2((j-1)k)} n^{((j-2)k)} \log n^{((j-2)k)}$  non-zeros blocks. The guarantees of Lemma 13.12.1 gives that the number of non-zero block in the sparsified matrix is at most

$$C 2^{2ck \log^2((j-1)k)} n^{((j-1)k)} \log n^{((j-1)k)}$$

for some absolute constant  $C$ . Therefore the cost of this additional call is less than the cost of constructing the rest of the phases during the construction process. Therefore, every recursive call except the first one can be viewed as an extra branching factor of 2 at each subsequent phase.

Depth can be bounded by the total number of recursive invocations to  $\text{RECURSIVECONSTRUCT}$ . The fact that  $n^{(i)}$  is geometrically decreasing means there are  $O(\log n/k)$  phases. Choosing  $k$  so that  $k \log^2 k = o(\log n)$  gives a depth of:

$$\begin{aligned} & \sum_{j=1}^{O(\log n/k)} 2^j \left( k \log(jk) \log n^{(jk)} + \log^2 n^{(jk)} \log \log n^{(jk)} \right) \\ &= 2^{O(\log n/k)} O \left( k \log \log n \log n^{(last)} + \log^2 n^{(last)} \log \log n^{(last)} \right) \\ &= 2^{O(\log n/k)} O \left( k^2 \log \log n + k^2 \log k \right) \\ &= 2^{O(\log n/k)} k^2 \log \log(n) \\ &= 2^{O(\log n/k)}. \end{aligned}$$

For bounding work, we start with the sparse case since all intermediate matrices generated during the construction process have density at most  $2^{O(k \log^2 k)}$ . In this case, the extra branching factor of 2 at each phase can be accounted for by weighting the cost of iteration  $j$  by  $2^j$ , giving:

$$\begin{aligned} & \sum_{j=1}^{O(\log n/k)} 2^j \left( 2^{O(k \log^2(jk))} n^{(jk)} \log^2 n^{(jk)} \right) \\ &= 2^{O(k \log^2 k)} n \log^2 n. \end{aligned}$$

For the dense case, the first recursive call to  $\text{SPARSIFY}(\mathbf{M}^{(0)}, O(1), 2^{2ck})$  is made to a graph whose edge count is much less. This leads to a geometric series, and an overhead of  $O(m \log n)$  work at each step. As this can happen at most  $O(\log n)$  times, it gives an additional factor of  $O(\log n)$  in depth, which is still  $2^{O(\log n/k)}$ . The work obeys the recurrence:

$$W(m)_{dense} \leq \begin{cases} 2^{O(k \log^2 k)} n \log^2 n & \text{if } m \leq 2^{O(k \log^2 k)} n \log n, \text{ and} \\ W_{dense}(m/2) + m \log n + 2^{O(k \log^2 k)} n \log^2 n & \text{otherwise.} \end{cases}$$

which solves to:

$$W_{dense}(m) \leq O(m \log n) + 2^{O(k \log^2 k)} n \log^2 n.$$

□

To obtain Theorem 13.1.2, we simply choose an appropriate initial padding and set the parameter  $k$ .

*Proof.* (of Theorem 13.1.2) Lemma 13.2.2 allows us to solve the linear system  $M + \varepsilon\mu I$  instead. Invoking Lemma 13.13.5 with  $k = \log \log \log n$  gives a depth of  $2^{O(\log n / \log \log \log n)}$ , and work of  $m \log n + 2^{O(\log \log \log n \log \log \log \log^4 n)} n \log^2 n$ . The depth term can be simplified to  $n^{o(1)}$  while the work term can be simplified to  $2^{O(\log \log \log^2 n)} = \log^{o(1)} n$ .  $\square$

## ■ 13.14 Weighted Expander Constructions

In this section, we give a linear time algorithm for computing linear sized spectral sparsifiers of complete and bipartite product demand graphs. These routines give Lemmas 13.11.2 and 13.11.1, which are crucial for controlling the densities of intermediate matrices in the spectral vertex sparsification algorithm from Section 13.11. Recall that the *product demand graph* with vertex set  $V$  and demands  $\mathbf{d} : V \rightarrow \mathbb{R}_{>0}$  is the complete graph in which the weight of edge  $(u, v)$  is the product  $d_u d_v$ . Similarly, the *bipartite demand graph* with vertex set  $U \cup V$  and demands  $\mathbf{d} : U \cup V \rightarrow \mathbb{R}_{>0}$  is the complete bipartite graph on which the weight of the edge  $(u, v)$  is the product  $d_u d_v$ . Our routines are based on reductions to the unweighted, uniform case. In particular, we

1. Split all of the high demand vertices into many vertices that all have the same demand. This demand will still be the highest.
2. Given a graph in which almost all of the vertices have the same highest demand, we
  - (a) drop all of the edges between vertices of lower demand,
  - (b) replace the complete graph between the vertices of highest demand with an expander, and
  - (c) replace the bipartite graph between the high and low demand vertices with a union of stars.
3. To finish, we merge back together the vertices that split off from each original vertex.

We start by showing how to construct the expanders that we need for step (2b). We state formally and analyze the rest of the algorithm for the complete case in the following two sections. We explain how to handle the bipartite case in Section 13.14.3.

Expanders give good approximations to unweighted complete graphs, and our constructions will use the spectrally best expanders—Ramanujan graphs. These are defined in terms of the eigenvalues of their adjacency matrices. We recall that the adjacency matrix of every  $d$ -regular graph has eigenvalue  $d$  with multiplicity 1 corresponding to the constant eigenvector. If the graph is bipartite, then it also has an eigenvalue of  $-d$  corresponding to an eigenvector that takes value 1 on one side of the bipartition and  $-1$  on the other side. These are called the *trivial* eigenvalues. A  $d$ -regular graph is called a Ramanujan graph if all of its non-trivial eigenvalues have absolute value at most  $2\sqrt{d-1}$ . Ramanujan graphs were constructed independently by Margulis [188] and Lubotzky, Phillips, and Sarnak [181]. The following theorem and proposition summarizes part of their results.

**Theorem 13.14.1.** *Let  $p$  and  $q$  be unequal primes congruent to 1 modulo 4. If  $p$  is a quadratic residue modulo  $q$ , then there is a non-bipartite Ramanujan graph of degree  $p+1$  with  $q^2(q-1)/2$  vertices. If  $p$  is not a quadratic residue modulo  $q$ , then there is a bipartite Ramanujan graph of degree  $p+1$  with  $q^2(q-1)$  vertices.*

The construction is explicit.

**Proposition 13.14.2.** *If  $p < q$ , then the graph guaranteed to exist by Theorem 13.14.1 can be constructed in parallel depth  $O(\log n)$  and work  $O(n)$ , where  $n$  is its number of vertices.*

*Proof.* When  $p$  is a quadratic residue modulo  $q$ , the graph is a Cayley graph of  $PSL(2, Z/qZ)$ . In the other case, it is a Cayley graph of  $PGL(2, Z/qZ)$ . In both cases, the generators are determined by the  $p + 1$  solutions to the equation  $p = a_0^2 + a_1^2 + a_2^2 + a_3^2$  where  $a_0 > 0$  is odd and  $a_1, a_2$ , and  $a_3$  are even. Clearly, all of the numbers  $a_0, a_1, a_2$  and  $a_3$  must be at most  $\sqrt{p}$ . So, we can compute a list of all sums  $a_0^2 + a_1^2$  and all of the sums  $a_2^2 + a_3^2$  with work  $O(p)$ , and thus a list of all  $p + 1$  solutions with work  $O(p^2) < O(n)$ .

As the construction requires arithmetic modulo  $q$ , it is convenient to compute the entire multiplication table modulo  $q$ . This takes time  $O(q^2) < O(n)$ . The construction also requires the computation of a square root of  $-1$  modulo  $q$ , which may be computed from the multiplication table. Given this data, the list of edges attached to each vertex of the graph may be produced using linear work and logarithmic depth.  $\square$

For our purposes, there are three obstacles to using these graphs:

1. They do not come in every degree.
2. They do not come in every number of vertices.
3. Some are bipartite and some are not.

We handle the first two issues by observing that the primes congruent to 1 modulo 4 are sufficiently dense. To address the third issue, we give a procedure to convert a non-bipartite expander into a bipartite expander, and *vice versa*.

An upper bound on the gaps between consecutive primes congruent to 1 modulo 4 can be obtained from the following theorem of Tchudakoff.

**Theorem 13.14.3.** *For two integers  $a$  and  $b$ , let  $p_i$  be the  $i$ th prime congruent to  $a$  modulo  $b$ . For every  $\varepsilon > 0$ ,*

$$p_{i+1} - p_i \leq O(p_i^{3/4+\varepsilon}).$$

**Corollary 13.14.4.** *There exists an  $n_0$  so that for all  $n \geq n_0$  there is a prime congruent to 1 modulo 4 between  $n$  and  $2n$ .*

We now explain how we convert between bipartite and non-bipartite expander graphs. To convert a non-bipartite expander into a bipartite expander, we take its double-cover. We recall that if  $G = (V, E)$  is a graph with adjacency matrix  $\mathbf{A}$ , then its double-cover is the graph with adjacency matrix

$$\begin{pmatrix} 0 & \mathbf{A} \\ \mathbf{A}^T & 0 \end{pmatrix}.$$

It is immediate from this construction that the eigenvalues of the adjacency matrix of the double-cover are the union of the eigenvalues of  $\mathbf{A}$  with the eigenvalues of  $-\mathbf{A}$ .

**Proposition 13.14.5.** *Let  $G$  be a connected,  $d$ -regular graph in which all matrix eigenvalues other than  $d$  are bounded in absolute value by  $\lambda$ . Then, all non-trivial adjacency matrix eigenvalues of the double-cover of  $G$  are also bounded in absolute value by  $\lambda$ .*

To convert a bipartite expander into a non-bipartite expander, we will simply collapse the two vertex sets onto one another. If  $G = (U \cup V, E)$  is a bipartite graph, we specify how the vertices of  $V$  are mapped onto  $U$  by a permutation  $\pi : V \rightarrow U$ . We then define the *collapse* of  $G$  induced by  $\pi$  to be the graph with vertex set  $U$  and edge set

$$\{(u, \pi(v)) : (u, v) \in E\}.$$

Note that the collapse will have self-loops at vertices  $u$  for which  $(u, v) \in E$  and  $u = \pi(v)$ . We assign a weight of 2 to every self loop. When a double-edge would be created, that is when  $(\pi(v), \pi^{-1}(u))$  is also an edge in the graph, we give the edge a weight of 2. Thus, the collapse can be a weighted graph.

**Proposition 13.14.6.** *Let  $G$  be a  $d$ -regular bipartite graph with all non-trivial adjacency matrix eigenvalues bounded by  $\lambda$ , and let  $H$  be a collapse of  $G$ . Then, every vertex in  $H$  has weighted degree  $2d$  and all adjacency matrix eigenvalues of  $H$  other than  $d$  are bounded in absolute value by  $2\lambda$ .*

*Proof.* To prove the bound on the eigenvalues, let  $G$  have adjacency matrix

$$\begin{pmatrix} 0 & \mathbf{A} \\ \mathbf{A}^T & 0 \end{pmatrix}.$$

After possibly rearranging rows and columns, we may assume that the adjacency matrix of the collapse is given by

$$\mathbf{A} + \mathbf{A}^T.$$

Note that the self-loops, if they exist, correspond to diagonal entries of value 2. Now, let  $\mathbf{x}$  be a unit vector orthogonal to the all-1s vector. We have

$$\mathbf{x}^T(\mathbf{A} + \mathbf{A}^T)\mathbf{x} = \begin{pmatrix} \mathbf{x} \\ \mathbf{x} \end{pmatrix}^T \begin{pmatrix} 0 & \mathbf{A} \\ \mathbf{A}^T & 0 \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{x} \end{pmatrix} \leq \lambda \left\| \begin{pmatrix} \mathbf{x} \\ \mathbf{x} \end{pmatrix} \right\|^2 \leq 2\lambda,$$

as the vector  $[\mathbf{x}; \mathbf{x}]$  is orthogonal to the eigenvectors of the trivial eigenvalues of the adjacency matrix of  $G$ .  $\square$

We now state how bounds on the eigenvalues of the adjacency matrices of graphs lead to approximations of complete graphs and complete bipartite graphs.

**Proposition 13.14.7.** *Let  $G$  be a graph with  $n$  vertices, possibly with self-loops and weighted edges, such that every vertex of  $G$  has weighted degree  $d$  and such that all non-trivial eigenvalues of the adjacency matrix of  $G$  have absolute value at most  $\lambda \leq d/2$ . If  $G$  is not bipartite, then  $(n/d)\mathbf{L}_G$  is an  $\varepsilon$ -approximation of  $K_n$  for  $\varepsilon = (2 \ln 2)(\lambda)/d$ . If  $G$  is bipartite, then  $(n/d)\mathbf{L}_G$  is an  $\varepsilon$ -approximation of  $K_{n,n}$  for  $\varepsilon = (2 \ln 2)(\lambda)/d$ .*

*Proof.* Let  $\mathbf{A}$  be the adjacency matrix of  $G$ . Then,

$$\mathbf{L}_G = d\mathbf{I} - \mathbf{A}.$$

In the non-bipartite case, we observe that all of the non-zero eigenvalues of  $\mathbf{L}_{K_n}$  are  $n$ , so for all vectors  $x$  orthogonal to the constant vector,

$$x^T \mathbf{L}_{K_n} x = nx^T x.$$

As all of the non-zero eigenvalues of  $\mathbf{L}_G$  are between  $d - \lambda$  and  $d + \lambda$ , for all vectors  $x$  orthogonal to the constant vector

$$n \left(1 - \frac{\lambda}{d}\right) x^T x \leq x^T (n/d)\mathbf{L}_G x \leq n \left(1 + \frac{\lambda}{d}\right) x^T x.$$

Thus,

$$\left(1 - \frac{\lambda}{d}\right) \mathbf{L}_{K_n} \preceq \mathbf{L}_G \preceq \left(1 + \frac{\lambda}{d}\right) \mathbf{L}_{K_n}.$$

In the bipartite case, we naturally assume that the bipartition is the same in both  $G$  and  $K_{n,n}$ . Now, let  $\mathbf{x}$  be any vector on the vertex set of  $G$ . Both the graphs  $K_{n,n}$  and  $(n/d)G$  have Laplacian

matrix eigenvalue 0 with the constant eigenvector, and eigenvalue  $2n$  with eigenvector  $[1; -1]$ . The other eigenvalues of the Laplacian of  $K_{n,n}$  are  $n$ , while the other eigenvalues of the Laplacian of  $(n/d)G$  are between

$$n \left(1 - \frac{\lambda}{d}\right) \quad \text{and} \quad n \left(1 + \frac{\lambda}{d}\right).$$

Thus,

$$\left(1 - \frac{\lambda}{d}\right) \mathbf{L}_{K_{n,n}} \preceq \mathbf{L}_G \preceq \left(1 + \frac{\lambda}{d}\right) \mathbf{L}_{K_{n,n}}.$$

The proposition now follows from our choice of  $\varepsilon$ , which guarantees that

$$e^{-\varepsilon} \leq 1 - \lambda/d \quad \text{and} \quad 1 + \lambda/d \leq e^\varepsilon,$$

provided that  $\lambda/d \leq 1/2$ . □

**Lemma 13.14.8.** *There are algorithms that on input  $n$  and  $\varepsilon > n^{-1/6}$  produce a graph having  $O(n/\varepsilon^2)$  edges that is an  $O(\varepsilon)$  approximation of  $K_{n'}$  or  $K_{n',n'}$  for some  $n \leq n' \leq 8n$ . These algorithms run in  $O(\log n)$  depth and  $O(n/\varepsilon^2)$  work.*

*Proof.* We first consider the problem of constructing an approximation of  $K_{n',n'}$ . By Corollary 13.14.4 there is a constant  $n_0$  so that if  $n > n_0$ , then there is a prime  $q$  that is equivalent to 1 modulo 4 so that  $q^2(q-1)$  is between  $n$  and  $8n$ . Let  $q$  be such a prime and let  $n' = q^2(q-1)$ . Similarly, for  $\varepsilon$  sufficiently small, there is a prime  $p$  equivalent to 1 modulo 4 that is between  $\varepsilon^{-2}/2$  and  $\varepsilon^{-2}$ . Our algorithm should construct the corresponding Ramanujan graph, as described in Theorem 13.14.1 and Proposition 13.14.2. If the graph is bipartite, then Proposition 13.14.7 tells us that it provides the desired approximation of  $K_{n',n'}$ . If the graph is not bipartite, then we form its double cover to obtain a bipartite graph and use Proposition 13.14.5 and Proposition 13.14.7 to see that it provides the desired approximation of  $K_{n',n'}$ .

The non-bipartite case is similar, except that we require a prime  $q$  so that  $q^2(q-1)/2$  is between  $n$  and  $8n$ , and we use a collapse to convert a bipartite expander to a non-bipartite one, as analyzed in Proposition 13.14.6. □

For the existence results in Section 13.10, we just need to know that there exist graphs of low degree that are good approximations of complete graphs. We may obtain them from the recent theorem of Marcus, Spielman and Srivastava that there exist bipartite Ramanujan graphs of every degree and number of vertices [187].

**Lemma 13.14.9.** *For every integer  $n$  and even integer  $d$ , there is a weighted graph on  $n$  vertices of degree at most  $d$  that is a  $4/\sqrt{d}$  approximation of  $K_n$ .*

*Proof.* The main theorem of [187] tells us that there is a bipartite Ramanujan graph on  $2n$  vertices of degree  $k$  for every  $k \leq n$ . By Propositions 13.14.6 and 13.14.7, a collapse of this graph is a weighted graph of degree at most  $2k$  that is a  $(4 \ln 2)/\sqrt{k}$  approximation of  $K_{n,n}$ . The result now follows by setting  $d = 2k$ . □

### ■ 13.14.1 Sparsifying Complete Product Demand Graphs

In the rest of the section, we will adapt these expander constructions to weighted settings. We start with product demand graphs.

---

**Algorithm 46:**  $G' = \text{WEIGHTEDEXPANDER}(\mathbf{d}, \varepsilon)$

---

1. Let  $\hat{n}$  be the least integer greater than  $2n/\varepsilon^2$  such that the algorithm described in Lemma 13.14.8 produces an  $\varepsilon$ -approximation of  $K_{\hat{n}}$ .
2. Let  $t = \frac{\sum_k d_k}{\hat{n}}$ .
3. Create a new product demand graph  $\hat{G}$  with demand vector  $\hat{\mathbf{d}}$  by splitting each vertex  $i$  into a set of  $\lceil d_i/t \rceil$  vertices,  $S_i$ :
  - (a)  $\lceil d_i/t \rceil$  vertices with demand  $t$ .
  - (b) one vertex with demand  $d_i - t \lceil d_i/t \rceil$ .
4. Let  $H$  be a set of  $\hat{n}$  vertices in  $\hat{G}$  with demand  $t$ , and let  $L$  contain the other vertices. Set  $k = |L|$ .
5. Partition  $H$  arbitrarily into sets  $V_1, \dots, V_k$ , so that  $|V_i| \geq \lfloor \hat{n}/k \rfloor$  for all  $1 \leq i \leq k$ .
6. Use the algorithm described in Lemma 13.14.8 to produce  $\tilde{K}_{HH}$ , an  $\varepsilon$ -approximation of the complete graph on  $H$ . Set

$$\tilde{G} = t^2 \tilde{K}_{HH} + \sum_{l \in L} \frac{|H|}{|V_l|} \sum_{h \in V_l} \hat{d}_l \hat{d}_h(l, h).$$

7. Let  $G'$  be the graph obtained by collapsing together all vertices in each set  $S_i$ .
-

*Lemma 13.11.1.* There is a routine  $\text{WEIGHTEDEXPANDER}(\mathbf{d}, \varepsilon)$  such that for any  $\varepsilon > 0$ , and a demand vector  $\mathbf{d} \in (\mathbb{R}_{>0})^n$ ,  $\text{WEIGHTEDEXPANDER}(\mathbf{d}, \varepsilon)$  returns in  $O(n\varepsilon^{-4})$  work and  $O(\log n)$  depth a graph  $H$  with  $O(n\varepsilon^{-4})$  edges such that

$$\mathbf{L}_H \approx_\varepsilon \mathbf{L}_{G(\mathbf{d})}.$$

Our algorithm for sparsifying complete product demand graphs begins by splitting the vertices of highest demands into many vertices. By *splitting* a vertex, we mean replacing it by many vertices whose demands sum to its original demand. In this way, we obtain a larger product demand graph. We observe that we can obtain a sparsifier of the original graph by sparsifying the larger graph, and then collapsing back together the vertices that were split.

**Proposition 13.14.10.** *Let  $G$  be a product demand graph with vertex set  $\{1, \dots, n\}$  and demands  $\mathbf{d}$ , and let  $\widehat{G} = (\widehat{V}, \widehat{E})$  be a product demand graph with demands  $\widehat{\mathbf{d}}$ . If there is a partition of  $\widehat{V}$  into sets  $S_1, \dots, S_n$  so that for all  $i \in V$ ,  $\sum_{j \in S_i} \widehat{d}_j = d_i$ , then  $\widehat{G}$  is a splitting of  $G$  and there is a matrix  $\mathbf{M}$  so that*

$$\mathbf{L}_G = \mathbf{M} \mathbf{L}_{\widehat{G}} \mathbf{M}^T.$$

*Proof.* The  $(i, j)$  entry of matrix  $\mathbf{M}$  is 1 if and only if  $j \in S_i$ . Otherwise, it is zero.  $\square$

We now show that we can sparsify  $G$  by sparsifying  $\widehat{G}$ .

**Proposition 13.14.11.** *Let  $\widehat{G}_1$  and  $\widehat{G}_2$  be graphs on the same vertex set  $\widehat{V}$  such that  $\widehat{G}_1 \approx_\varepsilon \widehat{G}_2$  for some  $\varepsilon$ . Let  $S_1, \dots, S_n$  be a partition of  $\widehat{V}$ , and let  $G_1$  and  $G_2$  be the graphs obtained by collapsing together all the vertices in each set  $S_i$  and eliminating any self loops that are created. Then*

$$G_1 \approx_\varepsilon G_2.$$

*Proof.* Let  $\mathbf{M}$  be the matrix introduced in Proposition 13.14.10. The proof follows from the fact that

$$\mathbf{L}_{G_1} = \mathbf{M} \mathbf{L}_{\widehat{G}_1} \mathbf{M}^T \quad \text{and} \quad \mathbf{L}_{G_2} = \mathbf{M} \mathbf{L}_{\widehat{G}_2} \mathbf{M}^T.$$

$\square$

For distinct vertices  $i$  and  $j$ , we let  $(i, j)$  denote the graph with an edge of weight 1 between vertex  $i$  and vertex  $j$ . If  $i = j$ , we let  $(i, j)$  be the empty graph. With this notation, we can express the product demand graph as

$$\sum_{i < j} d_i d_j (i, j) = \frac{1}{2} \sum_{i, j \in V} d_i d_j (i, j).$$

This notation also allows us to precisely express our algorithm for sparsifying product demand graphs. This section and the next are devoted to the analysis of this algorithm. Given Proposition 13.14.11, we just need to show that  $\widetilde{G}$  is a good approximation to  $\widehat{G}$ .

**Proposition 13.14.12.** *The number of vertices in  $\widehat{G}$  is at most  $n + \widehat{n}$ .*

*Proof.* The number of vertices in  $\widehat{G}$  is

$$\sum_{i \in V} \lceil d_i/t \rceil \leq n + \sum_{i \in V} d_i/t = n + \widehat{n}.$$

$\square$

So,  $k \leq n$  and  $\widehat{n} \geq 2k/\varepsilon^2$ . That is,  $|H| \geq 2|L|/\varepsilon^2$ . In the next section, we prove the lemmas that show that for these special product demand graphs  $\widehat{G}$  in which almost all weights are the maximum, our algorithm produces a graph  $\widetilde{G}$  that is a good approximation of  $\widehat{G}$ .



*Proof.* (of Lemma 13.11.1) The number of vertices in the graph  $\widehat{G}$  will be between  $n + 2n/\varepsilon^2$  and  $n + 16n/\varepsilon^2$ . So, the algorithm described in Lemma 13.14.8 will take  $O(\log n)$  depth and  $O(n/\varepsilon^4)$  work to produce an  $\varepsilon$  approximation of the complete graph on  $\widehat{n}$  vertices. This dominates the computational cost of the algorithm.

Proposition 13.14.11 tells us that  $G'$  approximates  $G$  at least as well as  $\widetilde{G}$  approximates  $\widehat{G}$ . To bound how well  $\widetilde{G}$  approximates  $\widehat{G}$ , we use two lemmas that are stated in the next section. Lemma 13.14.14 shows that

$$\widehat{G}_{HH} + \widehat{G}_{LH} \approx_{O(\varepsilon^2)} \widehat{G}.$$

Lemma 13.14.16 shows that

$$\widehat{G}_{HH} + \widehat{G}_{LH} \approx_{4\varepsilon} \widehat{G}_{HH} + \sum_{l \in L} \frac{|H|}{|V_l|} \sum_{h \in V_l} \widehat{d}_l \widehat{d}_h(l, h).$$

And, we already know that  $t^2 \widetilde{K}$  is an  $\varepsilon$ -approximation of  $\widehat{G}_{HH}$ . Fact 13.2.1 says that we can combine these three approximations to conclude that  $\widetilde{G}$  is an  $O(\varepsilon)$ -approximation of  $\widehat{G}$ .  $\square$

### ■ 13.14.2 Product demand graphs with most weights maximal

In this section, we consider product demand graphs in which almost all weights are the maximum. For simplicity, we make a slight change of notation from the previous section. We drop the hats, we let  $n$  be the number of vertices in the product demand graph, and we order the demands so that

$$d_1 \leq d_2 \leq \cdots \leq d_k \leq d_{k+1} = \cdots = d_n = 1.$$

We let  $L = \{1, \dots, k\}$  and  $H = \{k+1, \dots, n\}$  be the set of low and high demand vertices, respectively. Let  $G$  be the product demand graph corresponding to  $\mathbf{d}$ , and let  $G_{LL}$ ,  $G_{HH}$  and  $G_{LH}$  be the subgraphs containing the low-low, high-high and low-high edges respectively. We now show that little is lost by dropping the edges in  $G_{LL}$  when  $k$  is small.

Our analysis will make frequent use of the following Poincare inequality:

**Lemma 13.14.13.** *Let  $c(u, v)$  be an edge of weight  $c$  and let  $P$  be a path from  $u$  to  $v$  consisting of edges of weights  $c_1, c_2, \dots, c_k$ . Then*

$$c(u, v) \preceq c \left( \sum c_i^{-1} \right) P.$$

As the weights of the edges we consider in this section are determined by the demands of their vertices, we introduce the notation

$$(i, j)_2 = d_i d_j (i, j).$$

With this notation, we can express the product demand graph as

$$\sum_{i < j} (i, j)_2 = \frac{1}{2} \sum_{i, j \in V} (i, j)_2.$$

**Lemma 13.14.14.** *If  $|L| \leq |H|$ , then*

$$G_{HH} + G_{LH} \approx_{3 \frac{|L|}{|H|}} G.$$

*Proof.* The lower bound  $G_{HH} + G_{LH} \preceq G_{HH} + G_{LH} + G_{LL}$  follows from  $G_{LL} \succeq 0$ .

Using lemma 13.14.13 and the assumptions  $d_l \leq 1$  for  $l \in L$  and  $d_h = 1$  for  $h \in H$ , we derive for every  $l_1, l_2 \in L$ ,

$$(l_1, l_2)_2 = \frac{1}{|H|^2} \sum_{h_1, h_2 \in H} (l_1, l_2)_2$$

(by Lemma 13.14.13)

$$\begin{aligned} &\preceq \frac{1}{|H|^2} \sum_{h_1, h_2 \in H} d_{l_1} d_{l_2} \left( \frac{1}{d_{l_1} d_{h_1}} + \frac{1}{d_{h_1} d_{h_2}} + \frac{1}{d_{h_2} d_{l_2}} \right) ((l_1, h_1)_2 + (h_1, h_2)_2 + (h_2, l_2)_2) \\ &\preceq \frac{3}{|H|^2} \sum_{h_1, h_2 \in H} ((l_1, h_1)_2 + (h_1, h_2)_2 + (h_2, l_2)_2) \\ &= \frac{3}{|H|} \sum_{h \in H} ((l_1, h)_2 + (l_2, h)_2) + \frac{6}{|H|^2} G_{HH}. \end{aligned}$$

So,

$$\begin{aligned} G_{LL} &= \frac{1}{2} \sum_{l_1, l_2 \in L} (l_1, l_2)_2 \\ &\preceq \frac{1}{2} \sum_{l_1, l_2} \left( \frac{3}{|H|} \sum_{h \in H} ((l_1, h)_2 + (l_2, h)_2) + \frac{6}{|H|^2} G_{HH} \right) \\ &= \frac{3|L|}{|H|} G_{LH} + \frac{3|L|^2}{|H|^2} G_{HH}. \end{aligned}$$

The assumption  $|L| \leq |H|$  then allows us to conclude

$$G_{HH} + G_{LH} + G_{LL} \preceq \left( 1 + 3 \frac{|L|}{|H|} \right) (G_{HH} + G_{LH}).$$

□

Using a similar technique, we will show that the edges between  $L$  and  $H$  can be replaced by the union of a small number of stars. In particular, we will partition the vertices of  $H$  into  $k$  sets, and for each of these sets we will create one star connecting the vertices in that set to a corresponding vertex in  $L$ .

We employ the following consequence of the Poincare inequality in Lemma 13.14.13.

**Lemma 13.14.15.** *For any  $\varepsilon \leq 1$ ,  $l \in L$  and  $h_1, h_2 \in H$ ,*

$$\varepsilon(h_1, l)_2 + (1/2)(h_1, h_2)_2 \approx_{4\sqrt{\varepsilon}} \varepsilon(h_2, l)_2 + (1/2)(h_1, h_2)_2.$$

*Proof.* By applying Lemma 13.14.13 and recalling that  $d_{h_1} = d_{h_2} = 1$  and  $d_l \leq 1$ , we compute

$$\begin{aligned} (h_1, l)_2 &\preceq d_{h_1} d_l \left( \frac{\sqrt{\varepsilon}}{d_{h_1} d_{h_2}} + \frac{1}{d_{h_2} d_l} \right) \left( \frac{1}{\sqrt{\varepsilon}} (h_1, h_2)_2 + (h_2, l)_2 \right) \\ &\preceq \frac{1 + \sqrt{\varepsilon}}{\sqrt{\varepsilon}} (h_1, h_2)_2 + (1 + \sqrt{\varepsilon}) (h_2, l)_2 \\ &\preceq (1 + \sqrt{\varepsilon}) (h_2, l)_2 + \frac{2}{\sqrt{\varepsilon}} (h_1, h_2)_2. \end{aligned}$$

Multiplying both sides by  $\varepsilon$  and adding  $(1/2)(h_1, h_2)_2$  then gives

$$\begin{aligned} \varepsilon(h_1, l)_2 + (1/2)(h_1, h_2)_2 &\preceq (1 + \sqrt{\varepsilon})\varepsilon(h_2, l)_2 + (2\sqrt{\varepsilon} + 1/2)(h_1, h_2)_2 \\ &\preceq (1 + 4\sqrt{\varepsilon})(\varepsilon(h_2, l)_2 + (1/2)(h_1, h_2)_2) \\ &\preceq e^{4\sqrt{\varepsilon}}(\varepsilon(h_2, l)_2 + (1/2)(h_1, h_2)_2). \end{aligned}$$

By symmetry, we also have

$$\varepsilon(h_2, l)_2 + (1/2)(h_1, h_2)_2 \preceq e^{4\sqrt{\varepsilon}}(\varepsilon(h_1, l)_2 + (1/2)(h_1, h_2)_2).$$

□

**Lemma 13.14.16.** *Recall that  $L = \{1, \dots, k\}$  and let  $V_1, \dots, V_k$  be a partition of  $H = \{k+1, \dots, n\}$  so that  $|V_l| \geq s$  for all  $l$ . Then,*

$$G_{HH} + G_{LH} \approx_{4/\sqrt{s}} G_{HH} + \sum_{l \in L} \frac{|H|}{|V_l|} \sum_{h \in V_l} (l, h)_2.$$

*Proof.* Observe that

$$G_{LH} = \sum_{l \in L} \sum_{h \in H} (l, h)_2.$$

For each  $l \in L$ ,  $h_1 \in H$  and  $h_2 \in V_l$  we apply Lemma 13.14.15 to show that

$$\frac{1}{|V_l|}(l, h_1)_2 + \frac{1}{2}(h_1, h_2)_2 \approx_{4/\sqrt{s}} \frac{1}{|V_l|}(l, h_2)_2 + \frac{1}{2}(h_1, h_2)_2.$$

Summing this approximation over all  $h_2 \in V_l$  gives

$$(l, h_1)_2 + \sum_{h_2 \in V_l} \frac{1}{2}(h_1, h_2)_2 \approx_{4/\sqrt{s}} \sum_{h_2 \in V_l} \left( \frac{1}{|V_l|}(l, h_2)_2 + \frac{1}{2}(h_1, h_2)_2 \right).$$

Summing the left-hand side of this approximation over all  $l \in L$  and  $h_1 \in H$  gives

$$\sum_{l \in L, h_1 \in H} (l, h_1)_2 + \sum_{h_2 \in V_l} \frac{1}{2}(h_1, h_2)_2 = \sum_{l \in L, h_1 \in H} (l, h_1)_2 + \frac{1}{2} \sum_{h_1 \in H, l \in L} \sum_{h_2 \in V_l} (h_1, h_2)_2 = G_{LH} + G_{HH}.$$

On the other hand, the sum of the right-hand terms gives

$$G_{HH} + \sum_{l \in L, h_1 \in H} \sum_{h_2 \in V_l} \frac{1}{|V_l|}(l, h_2)_2 = G_{HH} + \sum_{l \in L} \sum_{h_2 \in V_l} \frac{|H|}{|V_l|}(l, h_2)_2.$$

□

### ■ 13.14.3 Weighted Bipartite Expanders

This construction extends analogously to bipartite product graphs. The bipartite product demand graph of vectors  $(\mathbf{d}^A, \mathbf{d}^B)$  is a complete bipartite graph whose weight between vertices  $i \in A$  and  $j \in B$  is given by  $w_{ij} = d_i^A d_j^B$ . The main result that we will prove is:

*Lemma 13.11.2.* There is a routine `WEIGHTEDBIPARTITEEXPANDER` $(\mathbf{d}^A, \mathbf{d}^B, \varepsilon)$  such that for any demand vectors  $\mathbf{d}^A$  and  $\mathbf{d}^B$  of total length  $n$ , and a parameter  $\varepsilon$ , it returns in  $O(n\varepsilon^{-4})$  work and  $O(\log n)$  depth a bipartite graph  $H$  between  $A$  and  $B$  with  $O(n\varepsilon^{-4})$  edges such that

$$\mathbf{L}_H \approx_\varepsilon \mathbf{L}_{G(\mathbf{d}^A, \mathbf{d}^B)}.$$

Without loss of generality, we will assume  $d_1^A \geq d_2^A \geq \dots \geq d_{n^A}^A$  and  $d_1^B \geq d_2^B \geq \dots \geq d_{n^B}^B$ . As the weights of the edges we consider in this section are determined by the demands of their vertices, we introduce the notation

$$(i, j)_2 = d_i^A d_j^B(i, j).$$

Our construction is based on a similar observation that if most vertices on  $A$  side have  $d_i^A$  equaling to  $d_1^A$  and most vertices on  $B$  side have  $d_i^B$  equaling to  $d_1^B$ , then the uniform demand graph on these vertices dominates the graph.

---

**Algorithm 47:**  $G' = \text{WEIGHTEDBIPARTITEEXPANDER}(\mathbf{d}^A, \mathbf{d}^B, \varepsilon)$

---

1. Let  $n' = \max(n^A, n^B)$  and  $\hat{n}$  be the least integer greater than  $2n'/\varepsilon^2$  such that the algorithm described in Lemma 13.14.8 produces an  $\varepsilon$ -approximation of  $K_{\hat{n}, \hat{n}}$ .
2. Let  $t^A = \frac{\sum_k d_k^A}{\hat{n}}$  and  $t^B = \frac{\sum_k d_k^B}{\hat{n}}$ .
3. Create a new bipartite demand graph  $\hat{G}$  with demands  $\hat{\mathbf{d}}^A$  and  $\hat{\mathbf{d}}^B$  follows:
  - (a) On the side  $A$  of the graph, for each vertex  $i$ , create a subset  $S_i$  consisting of  $\lceil d_i^A/t^A \rceil$  vertices:
    - i.  $\lfloor d_i^A/t^A \rfloor$  with demand  $t^A$ .
    - ii. one vertex with demand  $d_i^A - t^A \lfloor d_i^A/t^A \rfloor$ .
  - (b) Let  $H^A$  contain  $\hat{n}$  vertices of  $A$  of with demand  $t^A$ , and let  $L^A$  contain the rest. Set  $k^A = |L^A|$ .
  - (c) Create the side  $B$  of the graph with partition  $H^B, L^B$  and demand vector  $\hat{\mathbf{d}}^B$  similarly.
4. Partition  $H^A$  into sets of size  $|V_i^A| \geq \lfloor \hat{n}/k^A \rfloor$ , one corresponding to each vertex  $l \in L^A$ . Partition  $V_B$  similarly.
5. Let  $\tilde{K}_{H^A H^B}$  be a bipartite expander produced by Lemma 13.14.8 that  $\varepsilon$ -approximates  $K_{\hat{n}, \hat{n}}$ , identified with the vertices  $H^A$  and  $H^B$ .
6. Set

$$\tilde{G} = t^A t^B \tilde{K} + \sum_{l \in L^A} \frac{|H^B|}{|V_l^B|} \sum_{h \in V_l^B} \hat{d}_l^A \hat{d}_h^B(l, h) + \sum_{l \in L^B} \frac{|H^A|}{|V_l^A|} \sum_{h \in V_l^A} \hat{d}_l^B \hat{d}_h^A(l, h).$$

7. Let  $G'$  be the graph obtained by collapsing together all vertices in each set  $S_i^A$  and  $S_i^B$ .
- 

Similarly to the non-bipartite case, the Poincare inequality show that the edges between low demand vertices can be completely omitted if there are many high demand vertices which allows the demand routes through high demand vertices.

**Lemma 13.14.17.** *Let  $G$  be the bipartite product demand graph of the demand  $(\mathbf{d}_i^A, \mathbf{d}_j^B)$ . Let  $H^A$  a subset of vertices on  $A$  side with demand higher than the set of remaining vertices  $L^A$  on  $A$  side. Define  $H^B, L^B$  similarly. Assume that  $|L^A| \leq |H^A|$  and  $|L^B| \leq |H^B|$ , then*

$$G_{H^A H^B} + G_{H^A L^B} + G_{L^A H^B} \approx_{3 \max\left(\frac{|L^A|}{|H^A|}, \frac{|L^B|}{|H^B|}\right)} G.$$

*Proof.* The proof is analogous to Lemma 13.14.14, but with the upper bound modified for bipartite

graphs.

For every edge  $l_A, l_B$ , we embed it evenly into paths of the form  $l_A, h_B, h_A, l_B$  over all choices of  $h_A$  and  $h_B$ . The support of this embedding can be calculated using Lemma 13.14.13, and the overall accounting follows in the same manner as Lemma 13.14.14.  $\square$

It remains to show that the edges between low demand and high demand vertices can be compressed into a few edges. The proof here is also analogous to Lemma 13.14.15: we use the Poincare inequality to show that all demands can routes through high demand vertices. The structure of the bipartite graph makes it helpful to further abstract these inequalities via the following Lemma for four edges.

**Lemma 13.14.18.** *Let  $G$  be the bipartite product demand graph of the demand  $(d_i^A, d_j^B)$ . Given  $h_A, l_A \in A$  and  $h_{B,1}, h_{B,2} \in B$ . Assume that  $d_{h_A}^A = d_{h_{B,1}}^B = d_{h_{B,2}}^B \geq d_{l_A}^A$ . For any  $\varepsilon < 1$ , we have*

$$\varepsilon(l_A, h_{B,1})_2 + (h_A, h_{B,2})_2 + (h_A, h_{B,1})_2 \approx_{3\sqrt{\varepsilon}} \varepsilon(l_A, h_{B,2})_2 + (h_A, h_{B,2})_2 + (h_A, h_{B,1})_2.$$

*Proof.* Using Lemma 13.14.13 and  $d_{h_A}^A = d_{h_{B,1}}^B = d_{h_{B,2}}^B \geq d_{l_A}^A$ , we have

$$\begin{aligned} & (l_A, h_{B,1})_2 \\ & \preceq d_{l_A}^A d_{h_{B,1}}^B \left( \frac{1}{d_{l_A}^A d_{h_{B,2}}^B} + \frac{\sqrt{\varepsilon}}{d_{h_A}^A d_{h_{B,2}}^B} + \frac{\sqrt{\varepsilon}}{d_{h_A}^A d_{h_{B,1}}^B} \right) \left( (l_A, h_{B,2})_2 + \frac{1}{\sqrt{\varepsilon}}(h_A, h_{B,2})_2 + \frac{1}{\sqrt{\varepsilon}}(h_A, h_{B,1})_2 \right) \\ & \preceq (1 + 2\sqrt{\varepsilon})(l_A, h_{B,2})_2 + \frac{1 + 2\sqrt{\varepsilon}}{\sqrt{\varepsilon}}(h_A, h_{B,2})_2 + \frac{1 + 2\sqrt{\varepsilon}}{\sqrt{\varepsilon}}(h_A, h_{B,1})_2. \end{aligned}$$

Therefore,

$$\varepsilon(l_A, h_{B,1})_2 + (h_A, h_{B,2})_2 + (h_A, h_{B,1})_2 \preceq (1 + 3\sqrt{\varepsilon}) (\varepsilon(l_A, h_{B,2})_2 + (h_A, h_{B,2})_2 + (h_A, h_{B,1})_2).$$

The other side is similar due to the symmetry.  $\square$

(of Lemma 13.11.2) The proof is analogous to Lemma 13.11.1. After the splitting, the demands in  $H^A$  are higher than the demands in  $L^A$  and so is  $H^B$  to  $L^B$ . Therefore, Lemma 13.14.17 shows that that

*Proof.*

$$\widehat{G}_{H^A H^B} + \widehat{G}_{H^A L^B} + \widehat{G}_{L^A H^B} \approx_{3\varepsilon^2/2} \widehat{G}.$$

By a proof analogous to Lemma 13.14.16, one can use Lemma 13.14.18 to show that

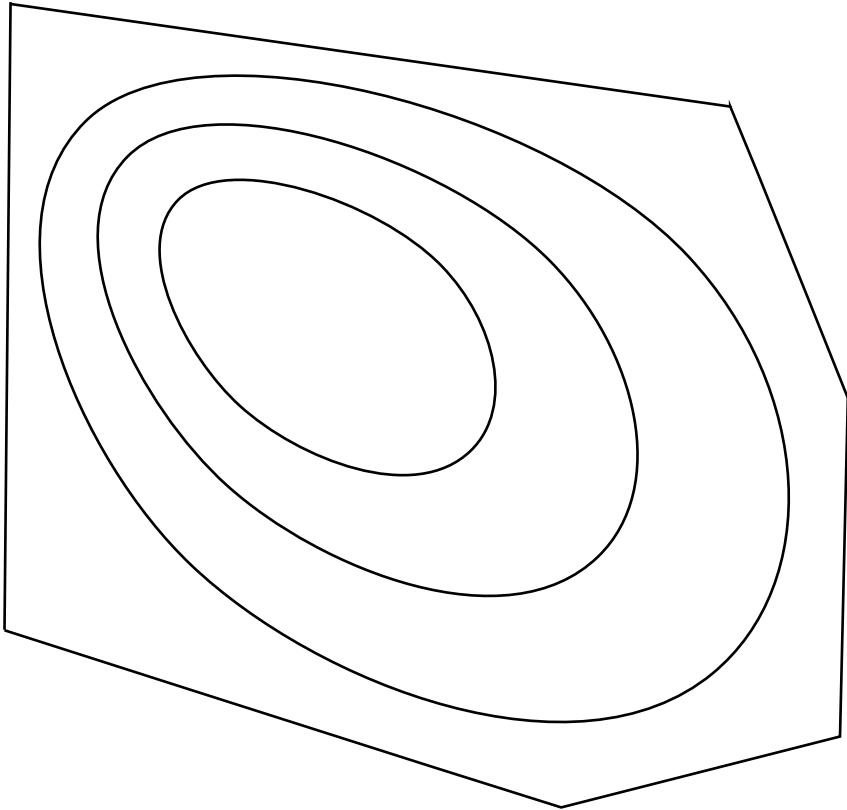
$$\widehat{G}_{H^A H^B} + \widehat{G}_{H^A L^B} + \widehat{G}_{L^A H^B} \approx_{O(\varepsilon)} \widehat{G}_{H^A H^B} + \frac{|H^B|}{|V_l^B|} \sum_{h \in V_l^B} \hat{d}_l^A \hat{d}_h^B(l, h) + \sum_{l \in L^B} \frac{|H^A|}{|V_l^A|} \sum_{h \in V_l^A} \hat{d}_l^B \hat{d}_h^A(l, h).$$

And, we already know that  $t^A t^B \tilde{K}$  is an  $\varepsilon$ -approximation of  $\widehat{G}_{H^A H^B}$ . Fact 13.2.1 says that we can combine these three approximations to conclude that  $\tilde{G}$  is an  $O(\varepsilon)$ -approximation of  $\widehat{G}$ .  $\square$

# Part IV

---

# Geometry



---

---

## Sampling In Sub-Quadratic Steps

### ■ 14.1 Introduction

In this chapter, we study the problem of sampling a high-dimensional polytope, a fundamental algorithmic problem with many applications. The problem can be solved in randomized polynomial time. Progress on the more general problem of sampling a convex body given only by a membership oracle [78, 79, 175, 176, 177, 131, 179, 178, 261] has lead to a set of general-purpose techniques, both for algorithms and for analysis in high dimension. All known algorithms are based on sampling by discrete-time Markov chains. These include the ball walk [174], hit-and-run [241, 179] and the Dikin walk [132], the last requiring stronger access than a membership oracle. In each case, the main challenge is analyzing the mixing time of the Markov chain. For a polytope defined by  $m$  inequalities in  $\mathbb{R}^n$ , the current best complexity of sampling is roughly the minimum of  $n^3 \cdot mn$  and  $mn \cdot mn^{\omega-1}$  where the first factor in each term is the mixing time and the second factor is the time to implement one step. In fact, the bound of  $n^3$  on the mixing time (achieved by the ball walk and hit-and-run) holds for arbitrary convex bodies, and  $O(mn)$  is just the time to implement a membership oracle. The second term is for the Dikin walk, for which Kannan and Narayanan showed a mixing time of  $O(mn)$  for the Dikin walk [132], with each step implementable in roughly matrix multiplication time. For general convex bodies given by membership oracles,  $\Omega(n^2)$  is a lower bound on the number of oracle calls for all known walks. A quadratic upper bound would essentially follow from a positive resolution of the KLS hyperplane conjecture (we mention that [59] show a mixing bound of  $\tilde{O}(n^2)$  for the ball walk for sampling from a Gaussian distribution restricted to a convex body). The quadratic barrier seems inherent for sampling convex bodies given by membership oracles, holding even for cubes and cylinders for the known walks based on membership oracles. It has not been surpassed thus far for explicitly described polytopes by any means.

For a polytope in  $\mathbb{R}^n$ , the Euclidean perspective is natural and predominant. The game has been to define a process on the points of the polytope so that the distribution of the current point quickly approaches the uniform (or a desired stationary distribution). The difficulty is that for points near the boundary of a body, steps are necessarily small due to the nature of volume distribution in high dimension. The Dikin walk departs from the standard perspective by making the distribution of the next step depend on the distance(s) of the current point to the boundary of the polytope. At each step, the process picks a random point from a suitable ellipsoid that is guaranteed to almost lie inside. This process adapts to the boundary, but runs into the similar difficulties — the ellipsoid has to shrink as the point approaches the boundary in order to ensure that (a) the stationary distribution is close to uniform and (b) the 1-step distribution has to be *smooth*, both necessary properties for rapidly converging to the uniform distribution. While the walk has the appealing property of being affine-invariant, and thus avoiding having to explicitly *round* the polytope, the current best upper bound is still quadratic, even for cylinders.

An alternative approach for sampling is the simulation of Brownian motion with boundary reflection [114, 67, 64, 43]. While there has been much study of this process, several difficulties arise in

turning it into an efficient algorithm. In particular, if the current point is close to the boundary of the polytope, extra care is needed in simulation and the process effectively slows down. However, if it is deep inside the polytope, we should expect that Brownian motion locally looks like a Gaussian distribution and hence it is easier to simulate. This suggests that the standard Euclidean metric, which does not take into account distance to the boundary, is perhaps not the right notion for getting a fast sampling algorithm.

In this chapter, we combine the use of stochastic differential equations (SDE) with non-Euclidean geometries (Riemannian manifolds) to break the quadratic barrier for mixing in polytopes. As a result we obtain significantly faster sampling algorithms.

Roughly speaking, our work is based on three key conceptual ideas. The first is the use of a Riemannian metric rather than the Euclidean metric. This allows us to scale space as we get closer to the boundary and incorporate boundary information much more smoothly. This idea was already used by Narayanan [198] to extend the Dikin walk to more general classes of convex bodies. The relevant metrics are induced by Hessians of convex barrier functions, objects that have been put to remarkable use for efficient linear and convex optimization [211]. The second idea is to simulate a Stochastic Differential Equation (SDE) on the manifold corresponding to the metric, via its geodesics (shortest-path curves, to be defined presently). Unlike straight lines, geodesics bend away from the boundary and this allows us to take larger steps while staying inside the polytope. The third idea is to use a modification of standard Brownian motion via a *drift* term, i.e., rather than centering the next step at the current point, we first shift the current point deterministically, then take a random step. This drift term compensates the changes of the step size and this makes the process closer to symmetric. Taken together, these ideas allow us to simulate an SDE by a discrete series of ordinary differential equations (ODE), which we are able to solve efficiently to the required accuracy. In order to state our contributions and results more precisely, we introduce some background, under three headings.

**Riemannian Geometry.** A manifold can be viewed as a surface embedded in a Euclidean space. Each point in the manifold (on the surface), has a tangent space (the linear approximate of the surface at that point) and a local metric. For a point  $x$  in a manifold  $M$ , the metric at  $x$  is defined by a positive definite matrix  $g(x)$  and the length of a vector  $u$  in the tangent space  $T_x M$  is defined as  $\|u\|_x \stackrel{\text{def}}{=} u^T g(x) u$ . By integration, the length of any curve on the manifold is defined as  $\int \|\frac{dc}{dt}\|_{c(t)}$ . A basic fact about Riemannian manifolds is that for any point in the manifold, in any direction (from the tangent space), there is locally a shortest path (*geodesic*) starting in that direction. In Euclidean space, this is just a straight line in the starting direction. Previous random walks involve generating a random direction and going along a straight line in that direction. However such straight lines do not take into account the local geometry, while geodesics do. We give formal definitions in Section 14.2.1.

**Hessian Manifolds.** In this chapter, we are concerned primarily with Riemannian manifolds induced by the Hessians of smooth (infinitely differentiable) convex functions. More precisely, for any such function  $\phi$ , the metric (of the manifold induced by  $\phi$ ) at a point  $x \in M$  is given by the Hessian of  $\phi$  at  $x$ , i.e.,  $\nabla^2 \phi(x)$ . Since  $\phi$  is convex, the Hessian is positive definite and hence the Riemannian manifold induced by  $\phi$  is well-defined and is called a Hessian manifold. In the context of convex optimization, we are interested in a class of convex functions called self-concordant barriers. Such convex function is smooth in a certain sense and blows up on the boundary of a certain convex set. The class of Hessian manifolds corresponding to self-concordant barriers has been studied and used to study interior-point methods (IPM) [137, 214, 212].

The barrier we are particularly interested is the logarithmic barrier. For a polytope  $Ax > b$ , with



$A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$ , for any  $x$  in the polytope, the logarithmic barrier is given by

$$\phi(x) = - \sum_i \ln(Ax - b)_i.$$

In this chapter, we study the logarithmic barrier and show how to use this to develop a faster algorithm for sampling polytopes.

**Stochastic Differential Equations.** Given a self-concordant barrier  $\phi$  on a convex set  $K$ , there is an “unique” Brownian motion with drift on the Hessian manifold  $M$  induced by  $\phi$  that has uniform stationary distribution. In the Euclidean coordinate system, the SDE is given by

$$dx_t = \mu(x_t)dt + (\nabla^2\phi(x_t))^{-1/2} dW_t \quad (14.1)$$

where the first term, called *drift*, is given by:

$$\mu_i(x_t) = \frac{1}{2} \sum_{j=1}^n \frac{\partial}{\partial x_j} \left( (\nabla^2\phi(x_t))^{-1} \right)_{ij}. \quad (14.2)$$

This suggests an approach for generating a random point in a polytope, namely to simulate the SDE. The running time of such an algorithm depends on the convergence rate of the SDE and the cost of simulating the SDE in discrete time steps.

Since the SDE is defined on the Riemannian manifold  $M$ , it is natural to consider the following geodesic walk:

$$x^{(j+1)} = \exp_{X^{(j)}}(\sqrt{h}w + \frac{h}{2}\mu(x^{(j)})) \quad (14.3)$$

where  $\exp_{x^{(j)}}$  is a map from  $T_{x^{(j)}}M$  back to the manifold,  $w$  is a random Gaussian vector on  $T_{x^{(j)}}M$ ,  $\mu(x^{(j)}) \in T_{x^{(j)}}M$  is the drift term and  $h$  is the step size. The coefficient of the drift term depends on the coordinate system we use; as we show in Lemma 14.3.1, for the coordinate system induced by a geodesic, the drift term is  $\mu/2$  instead of  $\mu$  as in (14.1). The Gaussian vector  $w$  has mean 0 and variance 1 in the metric at  $x$ , i.e.  $\mathbf{E}_w \|w\|_x^2 = 1$ . We write it as  $w \sim N_x(0, I)$ .

It can be shown that this discrete walk converges to (14.1) as  $h \rightarrow 0$  and it converges in a rate faster than the walk suggested by Euclidean coordinates (Theorem 14.3.3), namely,  $x^{(j+1)} = x^{(j)} + \sqrt{h}w + h\mu(x^{(j)})$ . (Note the drift here is proportional to  $h$  and not  $h/2$ .) This is the reason we study the geodesic walk.

### ■ 14.1.1 Algorithm

The algorithm is a discrete-time simulation of the geodesic process (14.3). For step-size  $h$  chosen in advance, let  $p(x \xrightarrow{w} y)$  be the probability density (in Euclidean coordinates) of going from  $x$  to  $y$  using the local step  $w$ . In general, the stationary distribution of the geodesic process is not uniform and it is difficult to analyze the stationary distribution unless  $h$  is very small, which would lead to a high number of steps. To get around this issue, we use the standard method of rejection sampling to get a uniform stationary distribution. We call this geodesic walk.

To implement the geodesic walk, we need to compute the exponential map  $\exp_x$  and  $p(x \xrightarrow{w} y)$  efficiently. We show how to implement this in Section 14.4.6. Our goal is to make sure that each iteration of geodesic walks only uses matrix multiplication and matrix inverse for  $O(\log^{O(1)} m)$  many  $O(m) \times O(m)$ -size matrices, and the rejection probability is small.

It turns out that the problems of computing the exponential map and  $p(x \xrightarrow{w} y)$  are similar; both involve solving some ordinary differential equations (ODEs) to accuracy  $1/n^{\Theta(1)}$ . Hence, one can view

---

**Algorithm 48:** Geodesic Walk (See Algo 49 for details)

---

Pick a Gaussian random vector  $w \sim N_x(0, 1)$ , i.e.  $\mathbf{E}_w \|w\|_x^2 = 1$ .

Compute  $y = \exp_x(\sqrt{h}w + \frac{h}{2}\mu(x))$  where  $\mu(x)$  is given by (14.2).

Let  $p(x \xrightarrow{w} y)$  is the probability density of going from  $x$  to  $y$  using the above step  $w$ .

Compute a corresponding  $w'$  s.t.  $x = \exp_y(\sqrt{h}w' + \frac{h}{2}\mu(y))$ .

With probability  $\min\left(1, \frac{p(y \xrightarrow{w'} x)}{p(x \xrightarrow{w} y)}\right)$ , go to  $y$ ;

Otherwise, stay at  $x$ .

---

the geodesic walk as reducing the problem of simulating an SDE 14.1 to solving a sequence of ODEs. Although solving ODEs is well-studied, existing literature seems quite implicit about the dependence on the dimension and hence it is difficult to apply it directly. In Section 14.6, we rederive some existing result about solving ODEs efficiently, but with quantitative estimates of the dependence on the dimension and desired error.

### ■ 14.1.2 Main results

In this chapter, we analyze the geodesic walk for the logarithmic barrier. The convergence analysis will need tools from Riemannian geometry and stochastic calculus, while the implementation uses efficient (approximate) solution of ODEs. Both aspects appear to be of independent interest. For the reader unfamiliar with these topics, we include an exposition of the relevant background.

We analyze geodesic walk in general and give a bound on the mixing time in terms of a set of manifold parameters (Theorem 14.4.2). Applying this on the logarithmic barrier, we obtain a faster sampling algorithm, going below the  $mn$  mixing time of the Dikin walk, while maintaining the same per-step complexity.

**Theorem 14.1.1** (Sampling with logarithmic barrier). *Given a polytope  $\{Ax \geq b\}$  with an uniformly random initial point in the polytope. Using the logarithmic barrier, we can sample another uniformly random point in the polytope via geodesic walk using  $\tilde{O}\left(mn^{\frac{3}{4}}\right)$  steps, with each step taking  $\tilde{O}\left(mn^{\omega-1}\right)$  time to implement.*

The implementation of each step of sampling is based on the efficient solution of high-dimensional ODEs (Theorem 14.6.4); this might have other applications.

### ■ 14.1.3 Discussion and open problems

At a high level, our algorithm is based on the following sequence of natural choices. First, we consider an Brownian motion that gives uniform stationary distribution and such Brownian motion is unique for a given metric (Fokker-Planck equation). Since the set we sample is a polytope, we use the metric given by the Hessian of a self-concordant barrier, a well-studied class of metrics in convex optimization. This allows us to reduce the sampling problem to the problem of simulating an SDE<sup>1</sup>. To simulate an SDE, we apply the Milstein method, well-known in that field. To implement the Milstein method, we perform a change of coordinates to make the metric locally constant and this makes most of the terms in the Milstein approximation method vanish. These coordinates are called normal coordinates

---

<sup>1</sup>Coincidentally, when we use the best known self-concordant barrier, the canonical barrier, our SDE becomes a Brownian motion with drift on an Einstein manifold. This is similar to how physical particles mix under general relativity, a sampling algorithm that has been executed for over 10 billion years!

and can be calculated by geodesics (Lemma 14.3.1). After all of this, we obtain the step of our walk (14.3).

There are two choices which are not natural. First, it is unclear that Hessians of self-concordant barriers are the best metrics for the sampling problems; after all, these barriers are designed for solving linear programs. Second, it is unclear that Milstein method is the best choice. There are other numerical method for SDE with better convergence rates, such as higher-order Runge-Kutta schemes. However, these methods are very complex and it is not clear how to implement the steps of these methods in  $\tilde{O}(mn^{\omega-1})$  time.

Clearly, we are not simultaneously experts on Riemannian geometry, numerical SDE, numerical ODE and convex geometry; we view this chapter as a sampling algorithm that connects different areas while improving the state of the art. Although the sampling problem is a harder problem than solving linear programs, the step size of geodesic walk we use is larger than that of the short-step interior point method. We hope that the relations revealed in this chapter might be useful for further development of both samplers and linear programming solvers.

There are several avenues for improving the sampling complexity further. One is to take longer steps and use a higher-order simulation of the SDE that is accurate up to a larger distance. Another part that is not tight in the analysis is the isoperimetry. At the moment, we incur a linear factor in the mixing time due to the isoperimetry ( $m$  for the log barrier,  $\tilde{O}(n)$  for the LS barrier). As far as we know, this factor might in fact be as small as  $O(1)$ . We make this precise via the following generalization of the KLS hyperplane conjecture. If true, it would directly improve our mixing time bound to  $\tilde{O}\left(n^{\frac{3}{4}}\right)$ .

**Conjecture 14.1.2.** *For any Hessian manifold  $M$  with metric  $d$ , let the isoperimetric ratio of a subset  $S$  w.r.t.  $d$  be defined as*

$$\psi_d(S) = \inf_{\varepsilon > 0} \frac{\text{vol}(\{x \in K \setminus S, d(x, y) \leq \varepsilon\})}{\varepsilon \cdot \min\{\text{vol}(S), \text{vol}(K \setminus S)\}}$$

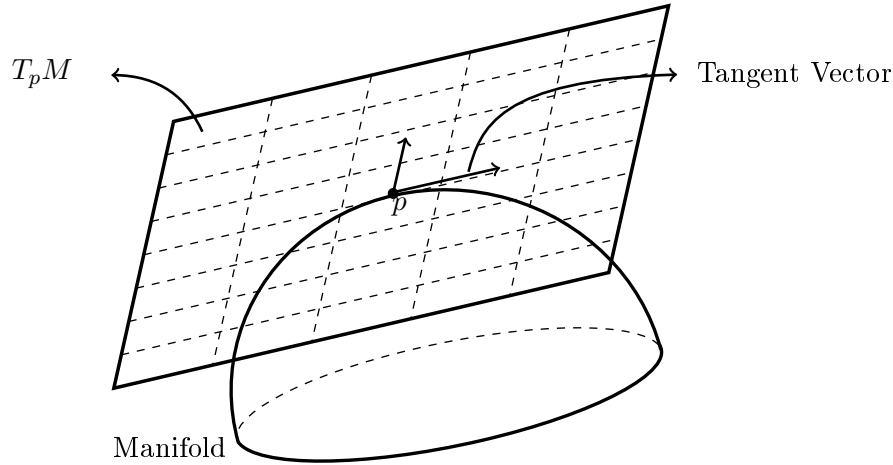
and the isoperimetric ratio of  $d$  as  $\psi = \inf_{S \subset K} \psi_d(S)$ . Then there is a subset  $S$  defined as a halfspace intersected with  $K$  with  $\psi_d(S) = O(\psi)$ .

#### ■ 14.1.4 Outline

In the next section, we recall the basic definitions and properties of Riemannian geometry and Hessian manifolds. Following that, in Section 14.3, we derive the discrete-time geodesic walk, showing how the magic formula naturally arises. In Section 14.4, we prove the main mixing bound for the geodesic walk in terms of a set of manifold parameters. In subsequent sections, we bound these parameters for the logarithmic barrier, thus proving Theorem 14.1.1. The algorithm for solving ODEs (collocation method) is presented and analyzed in Section 14.6.

### ■ 14.2 Preliminaries

Throughout the chapter, we use lowercase letter for vectors and vector fields and uppercase letter for matrices and tensors. We use  $e_k$  to denote coordinate vectors. We use  $\frac{d}{dt}$  for the usual derivative, e.g.  $\frac{df(c(t))}{dt}$  is the derivative of some function along a curve parametrized by  $t$ ,  $\nabla$  for the connection (manifold derivative, defined below which takes into account the local metric),  $D_v$  for the directional derivative wrt to the vector (or vector field)  $v$ , and  $D_t$  if the parametrization is clear from the context. We use  $g$  for the local metric. Given a point  $x \in M$ ,  $g$  is a matrix with entries  $g_{ij}$ . Its inverse has



**Figure 14-1:** Riemannian Manifold and Tangent Space

entries  $g^{ij}$ . Also  $n$  is the dimension,  $m$  the number of inequalities,  $\gamma$  is a geodesic, and  $\phi$  is a smooth convex function.

### ■ 14.2.1 Basic Definitions of Riemannian geometry

Here we first introduce basic notions of Riemannian geometry. In this chapter, we use lowercase letter for vectors and vector fields and uppercase letter for matrices (this is not the convention used in Riemannian geometry). One can think of a manifold  $M$  as a  $n$ -dimensional “surface” in  $\mathbb{R}^m$  for some  $m \geq n$ .

1. Tangent space  $T_p M$ : For any point  $p$ , the tangent space  $T_p M$  of  $M$  at point  $p$  is a linear subspace of  $\mathbb{R}^m$ . Intuitively,  $T_p M$  is the vector space of possible directions that are tangential to the manifold at  $x$ . Equivalently, it can be thought as the first-order linear approximation of the manifold  $M$  at  $p$ . For any curve  $c$  on  $M$ , the direction  $\frac{d}{dt}c(t)$  is tangent to  $M$  and hence lies in  $T_{c(t)}M$ . When it is clear from context, we define  $c'(t) = \frac{dc}{dt}(t)$ . For any open subset  $M$  of  $\mathbb{R}^n$ , we can identify  $T_p M$  with  $\mathbb{R}^n$  because all directions can be realized by derivatives of some curves in  $\mathbb{R}^n$ .
2. Riemannian metric: For any  $v, u \in T_p M$ , the inner product (Riemannian metric) at  $p$  is given by  $\langle v, u \rangle_p$  and this allows us to define the norm of a vector  $\|v\|_p = \sqrt{\langle v, v \rangle_p}$ . We call a manifold a Riemannian manifold if it is equipped with a Riemannian metric. When it is clear from context, we define  $\langle v, u \rangle = \langle v, u \rangle_p$ . In  $\mathbb{R}^n$ ,  $\langle v, u \rangle_p$  is the usual  $\ell_2$  inner product.
3. Pushforward  $d$ : Given a function  $f$  from a manifold  $M$  to a manifold  $N$ , we define  $df(x)$  as the linear map from  $T_x M$  to  $T_{f(x)}N$  such that

$$df(x)(c'(0)) = (f \circ c)'(0)$$

for any curve  $c$  on  $M$  starting at  $x = c(0)$ . When  $M$  and  $N$  are Euclidean spaces,  $df(x)$  is the Jacobian of  $f$  at  $x$ .

4. Hessian manifold: We call  $M$  a Hessian manifold if  $M$  is an open subset of  $\mathbb{R}^n$  with the Riemannian

nian metric at a point  $p \in M$  defined by

$$\langle v, u \rangle_p = v^T \nabla^2 \phi(p) u$$

where  $v, u \in T_p M$  and  $\phi$  is a smooth convex function on  $M$ .

5. Length: For any curve  $c : [0, 1] \rightarrow M$ , we define its length by

$$L(c) = \int_0^1 \left\| \frac{d}{dt} c(t) \right\|_{c(t)} dt.$$

6. Distance: For any  $x, y \in M$ , we define  $d(x, y)$  be the infimum of the lengths of all paths connecting  $x$  and  $y$ . In  $\mathbb{R}^n$ ,  $d(x, y) = \|x - y\|_2$ .

7. Geodesic: We call a curve  $\gamma(t) : [a, b] \rightarrow M$  a geodesic if

- (a) The curve  $\gamma(t)$  is parameterized with constant velocity. Namely, there is some  $c$  such that  $\left\| \frac{d}{dt} \gamma(t) \right\|_{\gamma(t)} = c$  for all  $a \leq t \leq b$ .
- (b) The curve is the locally shortest length curve between  $\gamma(a)$  and  $\gamma(b)$ . Namely, for any family of curve  $c(t, s)$  with  $c(t, 0) = \gamma(t)$  and  $c(0, a) = \gamma(a)$  and  $c(0, b) = \gamma(b)$ , we have that  $\frac{d}{ds} \Big|_{s=0} \int_a^b \left\| \frac{d}{dt} c(t, s) \right\|_{c(t,s)} dt = 0$ .

Note that, if  $\gamma(t)$  is a geodesic, then  $\gamma(\alpha t)$  is a geodesic for any  $\alpha$ . Intuitively, geodesics are local shortest-path curves. In  $\mathbb{R}^n$ , geodesics are straight lines.

8. Exponential map: The map  $\exp_p : T_p M \rightarrow M$  is defined as

$$\exp_p(v) = \gamma_v(1)$$

where  $\gamma_v$  is the unique geodesic starting at  $p$  with initial velocity  $\gamma'_v(0)$  equals to  $v$ . The exponential map takes a straight line  $tv \in T_p M$  to a geodesic  $\gamma_{tv}(1) = \gamma_v(t) \in M$ . Intuitively, the exponential map can be thought as vector addition in a manifold. In  $\mathbb{R}^n$ , we have  $\exp_p(v) = p + v$ .

9. Parallel transport: Given any geodesic  $c(t)$  and a vector  $v$  such that  $\langle v, c'(0) \rangle_{c(0)} = 0$ , we define the parallel transport of  $v$  along  $c(t)$  by the following process: Take  $h$  to be infinitesimally small and  $v_0 = v$ . For  $i = 1, 2, \dots, 1/h$ , we let  $v_{ih}$  be the vector orthogonal to  $c'(ih)$  that minimizes the distance between  $\exp_{c(ih)}(h v_{ih})$  and  $\exp_{c((i-1)h)}(h v_{(i-1)h})$ . Intuitively, the parallel transport finds the vectors on the curve such that their end points are closest to the end points of  $v$ . For general vector  $v \in T_{c'(0)}$ , we write  $v = \alpha c'(0) + w$  and we define the parallel transport of  $v$  along  $c(t)$  is the sum of  $\alpha c'(t)$  and the parallel transport of  $w$  along  $c(t)$ . For non-geodesic curve, see the definition in Fact 14.2.1.

10. Orthonormal frame: Given a vector fields  $v_1, v_2, \dots, v_n$  on a subset of  $M$ , we call  $\{v_i\}_{i=1}^n$  is an orthonormal frame if  $\langle v_i, v_j \rangle_x = \delta_{ij}$  for all  $x$ . Given a curve  $c(t)$  and an orthonormal frame at  $c(0)$ , we can extend it on the whole curve by parallel transport and it remains orthonormal on the whole curve.

11. Directional derivatives and Levi-Civita connection: Given any vector  $v \in T_p M$  and a vector field  $u$  in a neighborhood of  $p$ . Let  $\gamma_v$  is the unique geodesic starting at  $p$  with initial velocity  $\gamma'_v(0) = v$ , we define

$$\nabla_v u = \lim_{h \rightarrow 0} \frac{u(h) - u(0)}{h}$$

where  $u(h)$  is the parallel transport of  $u(\gamma(h))$  from  $\gamma(h)$  to  $\gamma(0)$ . Intuitively, Levi-Civita connection is the directional derivative of  $u$  along direction  $v$ , *taking the metric into account*. In particular, for  $\mathbb{R}^n$ , we have  $\nabla_v u(x) = \frac{d}{dt}u(x + tv)$ . When  $u$  is defined on a curve  $c$ , we define  $D_t u = \nabla_{c'(t)} u$ . In  $\mathbb{R}^n$ , we have  $D_t u(\gamma(t)) = \frac{d}{dt}u(\gamma(t))$ . We reserve  $\frac{d}{dt}$  for the usual derivative with Euclidean coordinates.

We list some basic facts about the definitions introduced above that are useful for computation and intuition.

**Fact 14.2.1.** *Given a manifold  $M$ , a curve  $c(t) \in M$ , a vector  $v$  and vector fields  $u, w$  on  $M$ , we have the following:*

1. (alternative definition of parallel transport)  $v(t)$  is the parallel transport of  $v$  along  $c(t)$  if and only if  $\nabla_{c'(t)} v(t) = 0$ .
2. (alternative definition of geodesic)  $c$  is a geodesic if and only if  $\nabla_{c'(t)} c'(t) = 0$ .
3. (linearity)  $\nabla_v(u + w) = \nabla_v u + \nabla_v w$ .
4. (product rule) For any scalar-valued function  $f$ ,  $\nabla_v(f \cdot u) = \frac{\partial f}{\partial v} u + f \cdot \nabla_v u$ .
5. (metric preserving)  $\frac{d}{dt} \langle u, w \rangle_{c(t)} = \langle D_t u, w \rangle + \langle u, D_t w \rangle$ .
6. (torsion free-ness) For any map  $c(t, s)$  from a subset of  $\mathbb{R}^2$  to  $M$ , we have that  $D_s \frac{\partial c}{\partial t} = D_t \frac{\partial c}{\partial s}$  where  $D_s = \nabla_{\frac{\partial c}{\partial s}}$  and  $D_t = \nabla_{\frac{\partial c}{\partial t}}$ .
7. (alternative definition of Levi-Civita connection)  $\nabla_v u$  is the unique linear mapping from the product of vector and vector field to vector field that satisfies (3), (4), (5) and (6).

### 14.2.1.1 Curvature

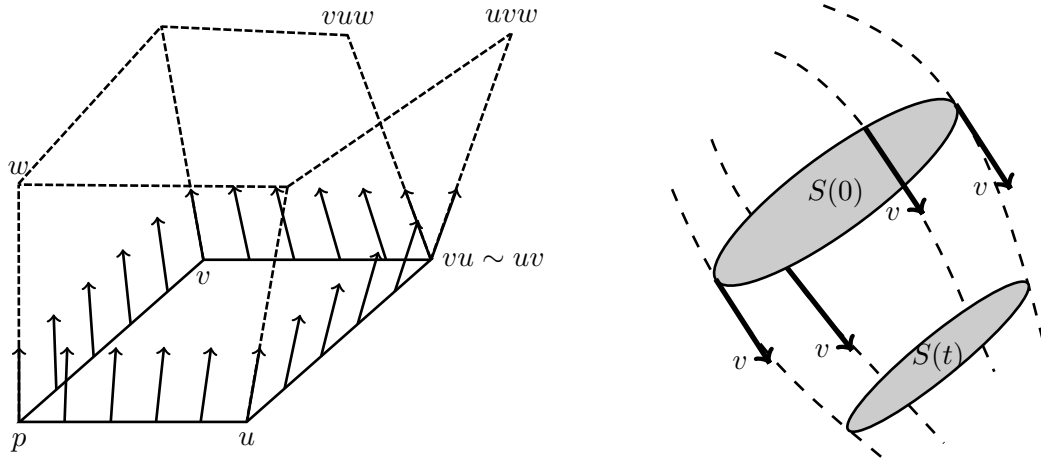
Here, we define various notions of curvature. Roughly speaking, they measure the amount by which a manifold deviates from Euclidean space.

Given vector  $u, v \in T_p M$ , in this section, we define  $uv$  be the point obtained from moving from  $x$  along direction  $u$  with distance  $\|u\|_p$  (using geodesic), then moving along direction “ $v$ ” with distance  $\|v\|_p$  where “ $v$ ” is the parallel transport of  $v$  along the path  $u$ . In  $\mathbb{R}^n$ ,  $uv$  is exactly  $p + u + v$  and hence  $uv = vu$ , namely, parallelograms close up. For a manifold, parallelograms almost close up, namely,  $d(uv, vu) = o(\|u\|\|v\|)$ . This property is called being *torsion-free*.

1. Riemann curvature tensor: Three-dimensional parallelepipeds might not close up, and the curvature tensor measures how far they are from closing up. Given vector  $u, v, w \in T_p M$ , we define  $uvw$  as the point obtained from moving from  $uv$  along direction “ $w$ ” for distance  $\|w\|_p$  where “ $w$ ” is the parallel transport of  $w$  along the path  $uv$ . In a manifold, parallelepipeds do not close up and the Riemann curvature tensor how much  $uvw$  deviates from  $vuw$ . Formally, for vector fields  $v, w$ , we define  $\tau_v w$  be the parallel transport of  $w$  along the vector field  $v$  for one unit of time. Given vector field  $v, w, u$ , we define the Riemann curvature tensor by

$$R(u, v)w = \frac{d}{ds} \frac{d}{dt} \tau_{su}^{-1} \tau_{tv}^{-1} \tau_{su} \tau_{tv} w \Big|_{t,s=0}. \quad (14.4)$$

Riemann curvature tensor is a tensor, namely,  $R(u, v)w$  at point  $p$  depends only on  $u(p)$ ,  $v(p)$  and  $w(p)$ .



**Figure 14-2:** Riemann curvature tensor measures the deviation of parallelepipeds (14.4) and Ricci curvature measures the change of volumes (14.5).

2. Ricci curvature: Given a vector  $v \in T_pM$ , the Ricci curvature  $\text{Ric}(v)$  measures if the geodesics starting around  $p$  with direction  $v$  converge together. Positive Ricci curvature indicates the geodesics converge while negative curvature indicates they diverge. For example, let  $S(0)$  be a small shape around  $p$  and  $S(t)$  be the set of point obtained from moving  $S(0)$  along direction  $v$  with  $t$  unit of time. Then,

$$\text{vol}S(t) = \text{vol}S(0)\left(1 - \frac{t^2}{2}\text{Ric}(v) + \text{smaller terms}\right). \tag{14.5}$$

Formally, we define

$$\text{Ric}(v) = \sum_{u_i} \langle R(v, u_i)u_i, v \rangle$$

where  $u_i$  is an orthonormal basis of  $T_pM$ . Equivalently, we have  $\text{Ric}(v) = \mathbf{E}_{u \sim N(0,I)} \langle R(v, u)u, v \rangle$ . For  $\mathbb{R}^n$ ,  $\text{Ric}(v) = 0$ . For a sphere in  $n + 1$  dimension with radius  $r$ ,  $\text{Ric}(v) = \frac{n-1}{r^2} \|v\|^2$ .

**Fact 14.2.2** (Alternative definition of Riemann curvature tensor). *Given vector fields  $x, y, z$  on  $M$ . Then, we have*

$$R(x, y)z = \nabla_x \nabla_y z - \nabla_y \nabla_x z - \nabla_{[x, y]}z$$

where  $[x, y]$  is the vector field such that  $\nabla_{[x, y]}f = \nabla_x \nabla_y f - \nabla_y \nabla_x f$  for any smooth scalar function  $f$ .

Given any  $M$ -valued function  $c(t, s)$ , we have vector fields  $\frac{\partial c}{\partial t}$  and  $\frac{\partial c}{\partial s}$  on  $M$ . Then, we have that

$$R\left(\frac{\partial c}{\partial t}, \frac{\partial c}{\partial s}\right)z = \nabla_{\frac{\partial c}{\partial t}} \nabla_{\frac{\partial c}{\partial s}} z - \nabla_{\frac{\partial c}{\partial s}} \nabla_{\frac{\partial c}{\partial t}} z.$$

Sometimes, we simply write it as  $R(\partial_t c, \partial_s c)z = D_t D_s z - D_s D_t z$ .

**Fact 14.2.3.** *Given vector fields  $v, u, w, z$  on  $M$ . Then, we have*

$$\langle R(v, u)w, z \rangle = \langle R(w, z)v, u \rangle = -\langle R(u, v)w, z \rangle = -\langle R(v, u)z, w \rangle.$$

### 14.2.1.2 Jacobi field

Let  $c : [0, \ell] \rightarrow M$  be a geodesic and  $c(t, s)$  be a variation of  $c(t)$  (i.e.  $c(t, 0) = c(t)$ ) such that  $c_s(t) \stackrel{\text{def}}{=} c(t, s)$  is a geodesic for all  $s$ . Then,  $x(t) \stackrel{\text{def}}{=} \frac{\partial}{\partial s} c(t, s)|_{s=0}$  satisfies the following equation

$$D_t D_t x + R(x, \frac{dc}{dt}) \frac{dc}{dt} = 0$$

where  $R(\cdot, \cdot)$  is Riemann curvature tensor defined before. We call any vector field satisfying the equation above a *Jacobi field*. In  $\mathbb{R}^n$ , geodesics are straight lines and Jacobi field is linear, namely,  $x(t) = x(0) + x'(0)t$ .

**Fact 14.2.4.** *Every Jacobi field  $x$  can be split into the tangential part  $x_1$  and the normal part  $x_2$  such that*

1.  $x = x_1 + x_2$ ,
2.  $x_1$  and  $x_2$  are Jacobi fields, namely,  $D_t D_t x_1 + R(x_1, \frac{dc}{dt}) \frac{dc}{dt} = 0$  and  $D_t D_t x_2 + R(x_2, \frac{dc}{dt}) \frac{dc}{dt} = 0$ ,
3.  $x_1$  is parallel to  $\frac{dc}{dt}$  and is linear, namely,  $x_1(t) = \left( \langle x(0), \frac{dc}{dt}(0) \rangle_{c(0)} + \langle D_t x(0), \frac{dc}{dt}(0) \rangle_{c(0)} t \right) \frac{dc}{dt}(t)$ ,
4.  $x_2$  is orthogonal to  $\frac{dc}{dt}$ , namely,  $\langle x_2(t), \frac{dc}{dt}(t) \rangle_{c(t)} = 0$ .

### 14.2.1.3 Hessian manifolds

Recall that a manifold is called Hessian if it is a subset of  $\mathbb{R}^n$  and its metric is given by  $g_{ij} = \frac{\partial^2 \phi}{\partial x^i \partial x^j}$  for some smooth convex function  $\phi$ . We let  $g^{ij}$  be entries of the inverse matrix of  $g_{ij}$ . For example, we have  $\sum_j g^{ij} g_{jk} = \delta_{ik}$ . We use  $\phi_{ij}$  to denote  $\frac{\partial^2 \phi}{\partial x^i \partial x^j}$  and  $\phi_{ijk}$  to denote  $\frac{\partial^3 \phi}{\partial x^i \partial x^j \partial x^k}$ .

Since Hessian manifold is a subset of Euclidean space, we identify tangent spaces  $T_p M$  by Euclidean coordinates. The following lemma give formulas for Levi-Civita connection and curvature under the Euclidean coordinates.

**Lemma 14.2.5** ([250]). *Given a Hessian manifold  $M$ , vector fields  $v, u, w, z$  on  $M$ , we have the following:*

1. (Levi-Civita connection)  $\nabla_v u = \sum_{ik} v_i \frac{\partial u_k}{\partial x^i} e_k + \sum_{ijk} v_i u_j \Gamma_{ij}^k e_k$  where  $e_k$  are coordinate vectors and the Christoffel symbol

$$\Gamma_{ij}^k = \frac{1}{2} \sum_l g^{kl} \phi_{ijl}$$

2. (Riemann curvature tensor)  $\langle R(u, v)w, z \rangle = \sum_{ijkl} R_{klij} u_i v_j w_l z_k$  where

$$R_{klij} = \frac{1}{4} \sum_{pq} g^{pq} (\phi_{jkp} \phi_{ilq} - \phi_{ikp} \phi_{jlq}) .$$

3. (Ricci curvature)  $Ric(v) = \frac{1}{4} \sum_{ijklpq} g^{pq} g^{jl} (\phi_{jkp} \phi_{ilq} - \phi_{ikp} \phi_{jlq}) v^i v^k$ .

In this chapter, geodesic is used everywhere and we will implicitly using the following lemma in all of our calculations.



**Lemma 14.2.6** ([214, Cor 3.1]). *If  $\phi$  is a self-concordant function, namely,  $|D^3f(x)[h, h, h]| \leq 2|D^2f(x)[h, h]|^{3/2}$ , then the corresponding Hessian manifold  $M$  is geodesically complete, namely, for any  $p \in M$ , the exponential map is defined on the entire tangent space  $T_pM$  and for any two points  $p, q \in M$ , there is a length minimizing geodesic connecting them.*

*In particular, for a polytope  $M = \{x : Ax > b\}$ , the Hessian manifold induced by the function  $\phi(x) = -\sum_i \log(a_i^T x - b_i)$  is geodesically complete.*

#### 14.2.1.4 Normal coordinates

For any manifold  $M$ , and any  $p \in M$ , the exponential map  $\exp_x$  maps from  $T_xM$  to  $M$ . Since  $T_xM$  is isomorphic to  $\mathbb{R}^n$ ,  $\exp_x^{-1}$  gives a local coordinate system of  $M$  around  $x$ . We call this system the normal coordinates at  $x$ . In a normal coordinate system, the metric is locally constant.

**Lemma 14.2.7.** *In normal coordinates, we have*

$$g_{ij} = \delta_{ij} - \frac{1}{3} \sum_{kl} R_{ikjl} x^k x^l + O(|x|^3).$$

For a Hessian manifold, one can do a linear transformation on the normal coordinate to make this coincide with Euclidean coordinate up to the first order.

**Lemma 14.2.8.** *Given a Hessian manifold  $M$  and any point  $x \in M$ . We pick a basis of  $T_xM$  such that*

$$\exp_x(tv) = x + tv + O(t^2).$$

*Let a local coordinate system  $F : M \rightarrow \mathbb{R}^n$  defined by  $F(y) = \exp_x^{-1}$ . Then, we have*

$$DF[h] = h \text{ and } D^2F_k(x)[h, h] = h^T \Gamma^k h$$

*where  $F_k$  is the  $k^{\text{th}}$  coordinate of  $F$  and  $\Gamma^k$  is the matrix with entries  $\Gamma_{ij}^k$  defined in Lemma 14.2.5.*

### ■ 14.2.2 Stochastic calculus

A stochastic differential equation (SDE) describes a stochastic process over a domain  $\Omega$ . It has the form  $dx_t = \mu(x_t, t)dt + \sigma(x_t, t)dW_t$  where  $x_t$  is the current point at time  $t$ ,  $W_t$  is a standard Brownian motion, and  $\mu(x_t, t), \sigma(x_t, t)$  are the mean and covariance of the next infinitesimal step at time  $t$ .

**Lemma 14.2.9** (Itô's lemma). *Given a SDE  $dx_t = \mu(x_t)dt + \sigma(x_t)dW_t$  and any smooth function  $f$ , we have*

$$df(t, x_t) = \left\{ \frac{\partial f}{\partial t} + \langle \nabla f, \mu \rangle + \frac{1}{2} \text{Tr} [\sigma^T (\nabla^2 f) \sigma] \right\} dt + (\nabla f)^T \sigma dW_t.$$

SDEs are closely related to diffusion equations:

$$\frac{\partial}{\partial t} p(x, t) = \frac{1}{2} \nabla \cdot (A(x, t) \nabla p(x, t))$$

where  $p(x, t)$  is the density at point  $x$  and time  $t$ ,  $\nabla \cdot$  is the usual divergence operator,  $\nabla p$  is the gradient of  $p$  and the matrix  $A(x, t)$  represents the diffusion coefficient at point  $x$  and time  $t$ . When  $A(x, t) = I$ , we get the familiar heat equation:

$$\frac{\partial}{\partial t} p(x, t) = \frac{1}{2} \Delta p(x, t).$$

In this chapter, the diffusion coefficient will be a symmetric positive definite matrix given by the Hessian  $(\nabla^2\phi(x))^{-1}$  of a convex function  $\phi(x)$ .

The Fokker-Planck equation connects an SDE to a diffusion equation.

**Theorem 14.2.10** (Fokker-Planck equation). *For any stochastic differential equation (SDE) of the form*

$$dx_t = \mu(x_t, t)dt + \sqrt{A(x_t, t)}dW_t,$$

*the probability density of the SDE is given by the diffusion equation*

$$\frac{\partial}{\partial t}p(x, t) = -\sum_{i=1}^n \frac{\partial}{\partial x_i}[\mu_i(x, t)p(x, t)] + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2}{\partial x_i \partial x_j}[A_{ij}(x, t)p(x, t)].$$

### ■ 14.2.3 Complex analysis

A complex function is said to be (*complex*) *analytic* (equivalently, *holomorphic*) if it is locally defined by a convergent power series. Hartog's theorem shows that a complex function in several variables  $f : \mathbb{C}^n \rightarrow \mathbb{C}$  is holomorphic iff it is analytic in each variable (while fixing all the other variables). For any power series expansion, we define the radius of convergence at  $x$  as the largest number  $r$  such that the series converges on the sphere with radius  $r$  centered at  $x$ . Roughly speaking, complex analytic functions behave very nicely up to the radius of convergence, and one can avoid complex and tedious computations via general convergence theorems.

**Theorem 14.2.11** (Cauchy's Estimates). *Suppose  $f$  is holomorphic on a neighborhood of the ball  $B \stackrel{\text{def}}{=} \{z \in \mathbb{C} : |z - z_0| \leq r\}$ , then we have that*

$$\left| f^{(k)}(z_0) \right| \leq \frac{k!}{r^k} \sup_{z \in B} |f(z)|.$$

*In particular, for any rational function  $f(z) = \frac{\prod_{i=1}^{\alpha}(z-a_i)}{\prod_{j=1}^{\beta}(z-b_j)}$  and any ball  $B \stackrel{\text{def}}{=} \{z \in \mathbb{C} : |z - z_0| \leq r\}$  such that  $b_j \notin B$ , we have that*

$$\left| f^{(k)}(z_0) \right| \leq \frac{k!}{r^k} \sup_{z \in B} |f(z)|.$$

A similar estimate holds for analytic functions in several variables.

## ■ 14.3 Intuition

In this section, we derive the formula of the geodesic walk from first principles.

### ■ 14.3.1 Derivation of the Geodesic walk

Given a smooth convex function  $\phi$  on the convex domain  $M$ , namely that it is convex and is infinitely differentiable at every interior point of  $M$ , we consider the corresponding diffusion equation

$$\frac{\partial}{\partial t}p(x, t) = \frac{1}{2} \nabla \cdot (\nabla^2\phi)^{-1} \nabla p.$$

We can expand it by

$$\begin{aligned}\frac{\partial}{\partial t}p(x, t) &= \frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial x_i} \left( \sum_{j=1}^n \left( (\nabla^2 \phi)^{-1} \right)_{ij} \frac{\partial}{\partial x_j} p(x, t) \right) \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2}{\partial x_i \partial x_j} \left( \left( (\nabla^2 \phi)^{-1} \right)_{ij} p(x, t) \right) - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \frac{\partial}{\partial x_i} \left( \frac{\partial}{\partial x_j} \left( (\nabla^2 \phi)^{-1} \right)_{ij} p(x, t) \right).\end{aligned}$$

The uniform distribution is the stationary distribution of this diffusion equation. Now applying the Fokker-Planck equation (Theorem 14.2.10) with  $A = (\nabla^2 \phi)^{-1}$ , the SDE for the above diffusion is given by:

$$dx_t = \mu(x_t)dt + (\nabla^2 \phi(x_t))^{-1/2} dW_t.$$

This explains the definition of (14.1). To simplify the notation, we write the SDE as

$$dx_t = \mu(x_t)dt + \sigma(x_t)dW_t \quad (14.6)$$

where  $\sigma(x_t) = (\nabla^2 \phi(x_t))^{-1/2}$ . One way to simulate this is via the Euler–Maruyama method, namely

$$x_{(t+1)h} = x_{th} + \mu(x_{th})h + \sigma(x_{th})w_{th}\sqrt{h}$$

where  $w_{th} \sim N_{x_{th}}(0, I)$ . We find the direction we are heading and take a small step along that direction. However, if we view  $M$  as a manifold, then directly adding the direction  $\mu(x_{th})h + \sigma(x_{th})w_{th}\sqrt{h}$  to  $x_{th}$  is not natural; the Euclidean coordinate is just an arbitrary coordinate system and we could pick any other coordinate systems and add the direction into  $x_{th}$ , giving a different step. Instead, we take the step in normal coordinates (Section 14.2.1.4).

In particular, given an initial point  $x_0$ , we define  $F = \exp_{x_0}^{-1}$  and we note that  $F(x_t)$  is another SDE. To see the defining equation of this transformed SDE, we use Itô's lemma to show that the transformed SDE looks the same but with half the drift term. This explains the formulation of geodesic walk:  $x^{(j+1)} = \exp_{x^{(j)}}(\sqrt{h}w + \frac{h}{2}\mu(x^{(j)}))$ .

**Lemma 14.3.1.** *Let  $F = \exp_{x_0}^{-1}$  and  $x_t$  satisfies the SDE (14.6) Then we have*

$$dF(x_0) = \frac{1}{2}\mu(x_0)dt + \sigma(x_0)dW_0. \quad (14.7)$$

*Proof.* Itô's lemma (Lemma 14.2.9) shows that

$$dF_k(x_t) = \left\{ \langle \nabla F_k, \mu \rangle + \frac{1}{2} \text{Tr} [\sigma^T (\nabla^2 F_k) \sigma] \right\} dt + (\nabla F_k)^T \sigma dW_t$$

where  $F_k$  indicates the  $k^{\text{th}}$  coordinate of  $F$ . From Lemma 14.2.8, we have that  $\langle \nabla F_k(x_0), \mu \rangle = \mu_k$ ,  $(\nabla F_k(x_0))^T \sigma = e_k^T \sigma$  and

$$\begin{aligned}\text{Tr} [\sigma(x_0)^T (\nabla^2 F_k(x_0)) \sigma(x_0)] &= \sum_i D^2 F_k[\sigma e_i, \sigma e_i] \\ &= \sum_i e_i^T \sigma^T \Gamma^k \sigma e_i \\ &= \text{Tr} (\sigma^T \Gamma^k \sigma) \\ &= \sum_i e_i^T \Gamma^k (\nabla^2 \phi)^{-1} e_i.\end{aligned}$$

Now, using Lemma 14.2.5, we have that

$$\mathrm{Tr} [\sigma(x_0)^T (\nabla^2 F_k(x_0)) \sigma(x_0)] = \frac{1}{2} \sum_{ijl} g^{kl} \phi_{ijl} g^{ji}$$

Hence, we have that

$$dF_k(x_0) = \left\{ \mu_k + \frac{1}{4} \sum_{ijl} g^{kl} \phi_{ijl} g^{ji} \right\} dt + e_k^T \sigma dW_t$$

Recall that the drift term (14.2) is given by

$$\begin{aligned} \mu_k &= \frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial x_i} \left( (\nabla^2 \phi)^{-1} \right)_{ki} \\ &= -\frac{1}{2} \sum_i e_k^T (\nabla^2 \phi)^{-1} \frac{\partial}{\partial x_i} \nabla^2 \phi (\nabla^2 \phi)^{-1} e_i \\ &= -\frac{1}{2} \sum_{l,i,j} g^{kl} \frac{\partial}{\partial x_i} \phi_{lj} g^{ji} \\ &= -\frac{1}{2} \sum_{ijl} g^{kl} \phi_{ijl} g^{ji}. \end{aligned} \tag{14.8}$$

Therefore, we have the result.  $\square$

### ■ 14.3.2 Discussion of Geodesic Walk

In this section, we relate our geodesic walk with the following commonly-used method for solving SDE.

**Theorem 14.3.2** (Milstein Method). *Given a SDE  $dx_t = \mu(x_t)dt + \sigma(x_t)dW_t$ , consider the algorithm*

$$\bar{x}_h = x_0 + \mu(x_0)h + \sigma(x_0)\Delta W + \sum_{j_1, j_2=1}^m L^{j_1} \sigma^{k, j_2} I_{j_1, j_2}$$

where  $L^j = \sum_{k=1}^d \sigma^{k, j} \frac{\partial}{\partial x^k}$ ,  $\Delta W_j \sim \int_0^h dW_t^j$  and  $I_{j_1, j_2} \sim \int_0^h \int_0^{t_1} dW_{t_2}^{j_2} dW_{t_1}^{j_1}$ . Then, we have

$$\mathbf{E} \|\bar{x}_h - x_h\|_2 = O(h^2)$$

where the  $O(\cdot)$  notation hides constants related to anything except  $h$ .

Under normal coordinates, the metric is locally constant (Lemma 14.2.7). Due to this, the term  $\sum_{j_1, j_2=1}^m L^{j_1} \sigma^{k, j_2} I_{j_1, j_2}$  in the Milstein method vanishes. Hence, we have the following result.

**Theorem 14.3.3** (Geodesic Walk). *Let  $X_t$  satisfies the SDE (14.6). Consider the discrete step*

$$\bar{x}_h = \exp_{x_0} \left( \frac{1}{2} \mu(x_0)h + \sigma(x_0)\Delta W \right)$$

where  $\Delta W \sim N(0, hI)$ . Then,

$$\mathbf{E} \|\bar{x}_h - x_h\|_2 = O(h^2)$$

where the  $O(\cdot)$  notation hides constants related to anything except  $h$ .

*Proof.* In normal coordinates, we consider the step

$$\overline{F(x_h)} = \frac{1}{2}\mu(x_0)h + \sigma(x_0)\Delta W.$$

Since  $F(x_h)$  satisfies 14.7 and  $\sigma$  is locally constant at  $X_0$  (Lemma 14.2.7), Theorem 14.3.2 (Milstein Method) shows that

$$\|\overline{F(x_h)} - F(x_h)\|_2 = O(h^2).$$

The result follows from  $\bar{x}_h = \exp_{X_0}(\overline{F(x_h)})$ ,  $x_h = \exp_{X_0}(F(x_h))$  and the smoothness of  $\exp_{x_0}$ .  $\square$

For the Euler-Maruyama method,  $\bar{x}_h = x_0 + \mu(x_0)h + \sigma(x_0)\Delta W$ , we only get  $O(h^{1.5})$  as the bound on the error, rather than  $O(h^2)$ . This is one of the reasons we use geodesic instead of just adding up the vector in Euclidean coordinates.

## ■ 14.4 Convergence of the Geodesic Walk

The geodesic walk is a Metropolis-filtered Markov chain, whose stationary distribution is the uniform distribution over the polytope to be sampled. We will prove that the conductance of this chain is large with an appropriate choice of the step-size parameter. Therefore, its mixing time to converge to the stationary distribution will be small. The proof of high conductance involves showing (a) the acceptance probability of the Metropolis filter is at least a constant (b) the induced metric satisfies a strong isoperimetric inequality (c) two points that are close in metric distance are also close in probabilistic distance, namely, the one-step distributions from them have large overlap. Besides bounding the number of steps, we also have to show that each step of the Markov chain can be implemented efficiently. We describe the implementation via reduction to solving ODEs in Section 14.4.6. We do this in later sections via an efficient algorithm for approximately solving ODEs.

In this section, we present the general conductance proof for Hessian manifolds. The proof will need bounds on seven parameters determined by the specific barrier function. In Section 14.5, we will bound these parameters for the logarithmic barrier.

For a Markov chain with state space  $M$ , stationary distribution  $Q$  and next step distribution  $P_u(\cdot)$  for any  $u \in M$ , the conductance of the Markov chain is

$$\phi = \inf_{S \subset M} \frac{\int P_u(M \setminus S) dQ(u)}{\min\{Q(S), Q(M \setminus S)\}}.$$

The conductance of an ergodic Markov chain allows us to bound its mixing time, i.e., the rate of convergence to its stationary distribution, e.g., via the following theorem of Lovász and Simonovits.

**Theorem 14.4.1** ([177]). *Let  $Q_t$  be the distribution of the current point after  $t$  steps of a Markov chain with stationary distribution  $Q$  and conductance at least  $\phi$ , starting from initial distribution  $Q_0$ . Then,*

$$d_{TV}(Q_t, Q) \leq \sqrt{d_0} \left(1 - \frac{\phi^2}{2}\right)^t$$

where  $d_0 = \mathbf{E}_{Q_0}(dQ_0(u)/dQ(u))$  is a measure of the distance of the starting distribution from the stationary and  $d_{TV}$  is the total variation distance.

### ■ 14.4.1 Hessian parameters

The mixing of the walk depends on the maximum values of several smoothness parameters of the manifold. For some of them, we only need to bound the parameters with high probability over the choice of steps of the geodesic walk. For  $x \in M$ , let  $\tilde{\Pi}_x$  be the set of geodesics used in one step of the geodesic walk starting at  $x$  with the parameterization  $\gamma : [0, \ell] \rightarrow M$  and  $\ell = \sqrt{n\hbar}$ . Parameters  $D_1, G_1, R_1, R_2$  below will be bounded with high probability over the choice of the next step (the others hold for all points in the domain). Let  $\Pi_x$  denote a  $1 - \exp(n^{-\Theta(1)})$  probability subset of  $\tilde{\Pi}_x$  and  $\Pi = \bigcup_{x \in M} \Pi_x$ . These will be fixed when bounding the parameters for specific manifolds.

1. The maximum norm of the drift,  $D_0 = \sup_{x \in M} \|\mu(x)\|_x$ .
2. The smoothness of drift norm,  $D_1 = \sup_{\gamma \in \Pi, 0 \leq t \leq \ell} \left\| \frac{d}{dt} \|\mu(\gamma(t))\|^2 \right\|$ .
3. The smoothness of the drift,  $D_2 = \sup_{x \in M, \|s\|_x \leq 1} \|\nabla_s \mu(x)\|_x$ .
4. The smoothness of the volume,  $G_1 = \sup_{\gamma \in \Pi, 0 \leq t \leq \ell} |\log \det(g(\gamma(t)))|$  where  $g(x)$  is the metric at  $x$ .
5. The smoothness of the metric,  $G_2 = \sup \frac{d(x,y)}{d_H(x,y)}$  where  $d_H$  is the Hilbert distance defined in Section 14.4.3 and  $d$  is the shortest path distance in  $M$ .
6. The stability of the Jacobian field,  $R_1 = \sup_{\gamma \in \Pi, 0 \leq t \leq \ell} \|R(t)\|_F$  where  $R(t)_{ij} = \langle R(X_i, \gamma'(t)) \gamma'(t), X_j \rangle$  where  $\{X_i\}$  is any orthonormal frame at  $\gamma(t)$ .
7. The smoothness of the Ricci curvature,  $R_2 = \sup_{\gamma \in \Pi} \left| \frac{d}{ds} \text{Ric}(\gamma'_s(t)) \right|$  (see Definition 14.4.16).

We refer to these as the smoothness parameters of a Hessian manifold. Our main theorem for convergence can be stated as follows.

**Theorem 14.4.2.** *On a Hessian manifold with smoothness parameters  $D_0, D_1, D_2, G_1, G_2, R_1, R_2$  and assumption 14.4.23, the geodesic walk with step size*

$$0 < h \leq \Theta(1) \min \left\{ \frac{1}{(nD_0R_1)^{2/3}}, \frac{1}{D_2}, \frac{1}{nR_1}, \frac{1}{n^{1/3}D_1^{2/3}}, \frac{1}{nG_1^{2/3}}, \frac{1}{(nR_2)^{2/3}} \right\}$$

has conductance  $\Omega(\sqrt{\hbar}/G_2)$  and its mixing time is  $O(G_2^2/h)$ .

### ■ 14.4.2 1-step distribution

We derive a formula for the drift term — it is in fact a Newton step of the volumetric barrier function  $\log \det \nabla^2 \phi(x)$ .

**Lemma 14.4.3.** *We have*

$$\mu(x) = -(\nabla^2 \phi(x))^{-1} \nabla \psi(x)$$

where  $\psi(x) = \frac{1}{2} \log \det \nabla^2 \phi(x)$ .

*Proof.* We note that

$$\begin{aligned} \frac{\partial}{\partial x_j} \log \det (\nabla^2 \phi)^{-1} &= \text{Tr} \left( (\nabla^2 \phi) \frac{\partial}{\partial x_j} (\nabla^2 \phi)^{-1} \right) \\ &= -\text{Tr} \left( (\nabla^2 \phi) (\nabla^2 \phi)^{-1} \left( \frac{\partial}{\partial x_j} \nabla^2 \phi \right) (\nabla^2 \phi)^{-1} \right) \\ &= -\sum_k e_k^T \left( \frac{\partial}{\partial x_j} \nabla^2 \phi \right) (\nabla^2 \phi)^{-1} e_k. \end{aligned}$$

Hence, we have

$$\begin{aligned} \frac{1}{2} e_i^T (\nabla^2 \phi)^{-1} \nabla \log \det (\nabla^2 \phi)^{-1} &= -\frac{1}{2} \sum_{jk} (\nabla^2 \phi)^{-1}_{ij} e_k^T \left( \frac{\partial}{\partial x_j} \nabla^2 \phi \right) (\nabla^2 \phi)^{-1} e_k \\ &= -\frac{1}{2} \sum_{jk} e_i^T (\nabla^2 \phi)^{-1} \left( \frac{\partial}{\partial x_k} \nabla^2 \phi \right) (\nabla^2 \phi)^{-1} e_k \end{aligned}$$

On the other hand, we have

$$\begin{aligned} \mu_i &= \frac{1}{2} \sum_k \frac{\partial}{\partial x_k} \left( (\nabla^2 \phi)^{-1} \right)_{ik} \\ &= -\frac{1}{2} \sum_k e_i^T (\nabla^2 \phi)^{-1} \left( \frac{\partial}{\partial x_k} \nabla^2 \phi \right) (\nabla^2 \phi)^{-1} e_k. \end{aligned}$$

□

To have a uniform stationary distribution, the geodesic walk uses a Metropolis filter. The transition probability before applying the filter is given as follows.

**Lemma 14.4.4.** *For any  $x \in M$  and  $h > 0$ , the probability density of the 1-step distribution from  $x$  (before applying the Metropolis filter) is given by*

$$p_x(y) = \sum_{v_x: \exp_x(v_x)=y} \det(d \exp_x(v_x))^{-1} \sqrt{\frac{\det(g(y))}{(2\pi h)^n}} \exp \left( -\frac{1}{2} \left\| \frac{v_x - \frac{h}{2} \mu(x)}{\sqrt{h}} \right\|_x^2 \right) \quad (14.9)$$

where  $y = \exp_x(v_x)$  and  $d \exp_x$  is the differential of the exponential map at  $x$ .

*Proof.* We prove the formula by separately considering each  $v_x \in T_x M$  s.t.  $\exp_x(v_x) = y$ , then summing up. In the tangent space  $T_x M$ , the point  $v_x$  is a Gaussian step. Therefore, the probability density of  $v_x$  in  $T_x M$  as follows.

$$p_x^{T_x M}(v_x) = \frac{1}{(2\pi h)^{n/2}} \exp \left( -\frac{1}{2} \left\| \frac{v_x - \frac{h}{2} \mu(x)}{\sqrt{h}} \right\|_x^2 \right).$$

Note that  $v_x, \mu(x) \in T_x M$ . Let  $y = \exp_x(v_x)$ . In the tangent space  $T_y M$ , we have that  $y$  maps to 0. Let  $\phi : T_x M \rightarrow K$  defined by  $F(v) = \text{id}_{M \rightarrow K} \circ \exp_x(v)$ . Here  $K$  is the same set as  $M$  but endowed with the Euclidean metric (we use  $F$  instead of  $\exp_y^{-1}(\exp_x(v))$  since the exponential map might not be 1-1). Hence, we have

$$dF(v_x) = d \text{id}_{M \rightarrow K}(y) d \exp_x(v_x).$$

The result follows from  $p_x(y) = \det(dF(v_x))^{-1} p_x^{T_x M}(v_x)$  and

$$\begin{aligned} \det dF(v_x) &= \det(\text{did}_{M \rightarrow K}(y)) \det(d \exp_x(v_x)) \\ &= \det(g(y))^{-1/2} \det(d \exp_x(v_x)). \end{aligned}$$

□

In Section 14.4.5, we bound the acceptance probability of the Metropolis filter.

**Lemma 14.4.5.** *Under the condition on  $h$  of Theorem 14.4.2,*

$$\left| \log \left( \frac{p_x(y)}{p_y(x)} \right) \right| = O \left( \sqrt{nh}^{3/2} D_1 + (nh)^{3/2} G_1 + (nh R_1)^2 \right) < \frac{1}{4}.$$

In Section 14.4.7, we bound the overlap of one-step distributions from nearby points.

**Lemma 14.4.6.** *Let  $x, y \in M$  be such that  $d(x, y) \leq 0.1\sqrt{h}$ . Then, under the condition on  $h$  of Theorem 14.4.2, the one-step distributions  $P_x, P_y$  from  $x, y$  satisfy*

$$d_{TV}(P_x, P_y) = O(nh^{3/2} R_2 + h D_2 + nh^{3/2} R_1 D_0) < 0.3.$$

### ■ 14.4.3 Isoperimetry

For a convex body  $K$ , the *cross-ratio distance* of  $x$  and  $y$  is

$$d_K(x, y) = \frac{|x - y||p - q|}{|p - x||y - q|}$$

where  $p$  and  $q$  are on the boundary of  $K$  such that  $p, x, y, q$  are on the straight line  $\overline{xy}$  and are in order. In this section, we show that if the distance  $d(x, y)$  induced by the Riemannian metric is upper bounded by the cross-ratio distance, then the body has good isoperimetric constant in terms of  $d$ . We note that although the cross-ratio distance is not a metric, the closely-related Hilbert distance is a metric:

$$d_H(x, y) = \log \left( 1 + \frac{|x - y||p - q|}{|p - x||y - q|} \right).$$

**Theorem 14.4.7.** *For a Hessian manifold  $M$  with smoothness parameters  $G_2$ , for any partition of  $M$  into three measurable subsets  $S_1, S_2, S_3$ , we have that*

$$\text{vol}(S_3) \geq \frac{d(S_1, S_2)}{G_2} \min\{\text{vol}(S_1), \text{vol}(S_2)\}.$$

The theorem follows from the following isoperimetric inequality from [173], the definition of  $G_2$  and the fact  $d_H \leq d_K$ .

**Theorem 14.4.8** ([173]). *For any convex body  $K$  and any partition of  $K$  into disjoint measurable subsets  $S_1, S_2, S_3$*

$$\text{vol}(S_3) \geq d_K(S_1, S_2) \text{vol}(S_1) \text{vol}(S_2).$$



### ■ 14.4.4 Conductance and Mixing

*Proof of Theorem 14.4.2.* The proof follows the standard outline for geometric random walks (see e.g., [261]). Let  $Q$  be the uniform distribution over  $M$  and  $S$  be any measurable subset of  $M$ . Then our goal is to show that

$$\frac{\int_S P_x(M \setminus S) dQ(x)}{\min\{Q(S), Q(M \setminus S)\}} = \Omega\left(\frac{\sqrt{h}}{G_2}\right).$$

Since the Markov chain is time-reversible (For any two subsets  $A, B$ ,  $\int_A P_x(B) dx = \int_B P_x(A) dx$ ), we can write the numerator of the LHS above as

$$\frac{1}{2} \left( \int_S P_x(M \setminus S) dQ(x) + \int_{M \setminus S} P_x(S) dQ(x) \right).$$

Define

$$\begin{aligned} S_1 &= \{x \in S : P_x(M \setminus S) < 0.05\} \\ S_2 &= \{x \in M \setminus S : P_x(S) < 0.05\} \\ S_3 &= M \setminus S_1 \setminus S_2. \end{aligned}$$

We can assume wlog that  $Q(S_1) \geq (1/2)Q(S)$  and  $Q(S_2) \geq (1/2)Q(M \setminus S)$  (if not, the conductance is  $\Omega(1)$ ).

Next, we note that for any two points  $x \in S_1$  and  $y \in S_2$ ,  $d_{TV}(P_x, P_y) > 0.9$ . Therefore, by Lemma 14.4.6, we have that  $d(x, y) \geq 0.1\sqrt{h}$  and hence  $d(S_1, S_2) \geq 0.1\sqrt{h}$ . Therefore, using Theorem 14.4.7,

$$\text{vol}(S_3) \geq \frac{0.1\sqrt{h}}{G_2} \min\{\text{vol}(S_1), \text{vol}(S_2)\}.$$

Going back to the conductance,

$$\begin{aligned} \frac{1}{2} \left( \int_S P_x(M \setminus S) dQ(x) + \int_{M \setminus S} P_x(S) dQ(x) \right) &\geq \frac{1}{2} \int_{S_3} (0.05) dQ(x) \\ &\geq \frac{\sqrt{h}}{800G_2} \min\{\text{vol}(S_1), \text{vol}(S_2)\} \frac{1}{\text{vol}(K)} \\ &= \frac{\sqrt{h}}{1600G_2} \min\left\{ \frac{\text{vol}(S)}{\text{vol}(M)}, \frac{\text{vol}(M \setminus S)}{\text{vol}(M)} \right\} \\ &= \frac{\sqrt{h}}{1600G_2} \min\{Q(S), Q(M \setminus S)\} \end{aligned}$$

Therefore,  $\phi(S) \geq (1/1600)\sqrt{h}/G_2$ . □

**Corollary 14.4.9.** *Let  $K$  be a polytope with  $m$  facets. Let  $Q$  be the uniform distribution over  $K$  and  $Q_t$  be the distribution obtained after  $t$  steps of the geodesic walk started from a distribution  $Q_0$  with  $d_0 = \sup_K \frac{dQ_0}{dQ}$ . Then after  $t > C(G_2^2/h) \log\left(\frac{d_0}{\varepsilon}\right)$  steps, with probability at least  $1 - \delta$ , we have  $d_{TV}(Q_t, Q) \leq \varepsilon$ .*

### ■ 14.4.5 Rejection Probability

The goal of this section is to prove Lemma 14.4.5, i.e., the rejection probability of the Metropolis filter is small. For a transition from  $x$  to  $y$  on the manifold, the filter is applied with respect to a randomly chosen  $v_x$ , i.e., for one geodesic from  $x$  to  $y$ . We will bound the ratio of the transition probabilities (without the filter) as follows:

$$\begin{aligned} \log \left( \frac{p(x \xrightarrow{v_x} y)}{p(y \xrightarrow{v_y} x)} \right) &= \log \left( \frac{\det(d \exp_x(v_x))^{-1}}{\det(d \exp_y(v_y))^{-1}} \right) + \frac{1}{2} \log \det(g(y)) - \frac{1}{2} \log \det(g(x)) \\ &\quad - \frac{1}{2} \left\| \frac{v_x - \frac{h}{2} \mu(x)}{\sqrt{h}} \right\|_2^2 + \left\| \frac{v_y - \frac{h}{2} \mu(y)}{\sqrt{h}} \right\|_2^2. \end{aligned}$$

Since a geodesic has constant speed, we have  $\|v_x\| = \|v_y\|$ . Therefore, we have that

$$\begin{aligned} \log \left( \frac{p(x \xrightarrow{v_x} y)}{p(y \xrightarrow{v_y} x)} \right) &= \log \left( \frac{\det(d \exp_y(v_y))}{\det(d \exp_x(v_x))} \right) + \frac{1}{2} \log \det(g(y)) - \frac{1}{2} \log \det(g(x)) \quad (14.10) \\ &\quad + \frac{1}{2} \langle v_x, \mu(x) \rangle - \frac{h}{8} \|\mu(x)\|^2 - \frac{1}{2} \langle v_y, \mu(y) \rangle + \frac{h}{8} \|\mu(y)\|^2. \end{aligned}$$

We separate the proof into three parts:

- $|\|\mu(x)\|^2 - \|\mu(y)\|^2| \leq O(\sqrt{nh}D_1)$ , immediate from the definition of  $D_1$  and the fact that the  $\|x - y\| = O(\sqrt{nh})$  whp.
- Sec 14.4.5.1:  $|\log \det(g(y)) - \log \det(g(x)) + \langle v_x, \mu(x) \rangle - \langle v_y, \mu(y) \rangle| = O((nh)^{3/2}G_1)$ ,
- Sec 14.4.5.2: Assuming that  $h < \frac{1}{2nR_1}$ , we have  $\left| \log \left( \frac{\det(d \exp_y(v_y))}{\det(d \exp_x(v_x))} \right) \right| = O((nhR_1)^2)$ .

Together, these facts imply the lemma.

#### 14.4.5.1 Trapezoidal Rule for Volumetric Barrier

Recall that the Trapezoidal rule is to approximate  $\int_0^h f(t)dt$  by  $\frac{h}{2} (f(0) + f(h))$ . The nice thing about this rule is that the error is  $O(h^3)$  instead of  $O(h^2)$  because the second order term cancels by symmetry. Our main observation here is that the geodesic walk is implicitly following a trapezoidal rule on the metric and hence it has a small error.

**Lemma 14.4.10.** *We have that*

$$\left| \int_0^\ell f'(t)dt - \frac{\ell}{2} (f'(0) + f'(\ell)) \right| \leq \frac{\ell^3}{12} \max_{0 \leq t \leq \ell} |f'''(t)|.$$

*Proof.* Note that

$$\begin{aligned}
\int_0^\ell f'(t)dt - \frac{\ell}{2}(f'(0) + f'(\ell)) &= \int_0^\ell \left( f'(0) + \int_0^t f''(s)ds \right) dt - \ell f'(0) - \frac{\ell}{2} \int_0^\ell f''(s)ds \\
&= \int_0^\ell \int_0^t f''(s)dsdt - \frac{\ell}{2} \int_0^\ell f''(s)ds \\
&= \int_0^\ell \left( \frac{\ell}{2} - s \right) f''(s)ds \\
&= \int_0^\ell \left( \frac{\ell}{2} - s \right) \left( f''(0) + \int_0^s f'''(t)dt \right) ds \\
&= \int_0^\ell \left( \frac{\ell}{2} - s \right) \int_0^s f'''(t)dt ds \\
&\leq \frac{\ell^3}{12} \max_{0 \leq t \leq \ell} |f'''(t)|.
\end{aligned}$$

□

We apply this to the logdet function.

**Lemma 14.4.11.** *Let  $f(t) = \log \det(g(\gamma(t)))$  where  $\gamma(t) = \exp_x(\frac{t}{\ell}v_x) \in \Pi$  (See Definition of  $\Pi$  in Section 14.4.1). Then,*

$$\left| \log \det(g(y)) - \log \det(g(x)) + \langle v_x, \mu(x) \rangle_x - \langle v_y, \mu(y) \rangle_y \right| = O((nh)^{3/2}G_1).$$

*Proof.* Let  $f(t) = \log \det g(\gamma(t))$  and  $z = \gamma(t)$ . By Lemma 14.4.3,  $\mu(z) = -\frac{1}{2}g(z)^{-1}\nabla f(z)$ . Using this,

$$\begin{aligned}
f'(t) &= \langle \nabla_z \log \det g(z), \gamma'(t) \rangle_z \\
&= \langle g(z)^{-1} \nabla_z \log \det g(z), \gamma'(t) \rangle_z \\
&= -\langle 2\mu(z), \gamma'(t) \rangle_z.
\end{aligned}$$

Noting that  $v_x = \ell\gamma'(0)$  and  $v_y = -\ell\gamma'(\ell)$ , and using Lemma 14.4.10, we have

$$\begin{aligned}
&\left| \log \det(g(y)) - \log \det(g(x)) + \langle v_x, \mu(x) \rangle_x - \langle v_y, \mu(y) \rangle_y \right| \\
&= \left| \log \det(g(y)) - \log \det(g(x)) + \ell \left( \langle \gamma'(0), \mu(x) \rangle_x + \langle \gamma'(\ell), \mu(y) \rangle_y \right) \right| \\
&= \left| \int_0^\ell f'(t)dt - \frac{\ell}{2}(f'(0) + f'(\ell)) \right| \\
&\leq \frac{\ell^3}{12} \max_{0 \leq t \leq \ell} |f'''(t)| = O\left((nh)^{3/2}G_1\right).
\end{aligned}$$

□

### 14.4.5.2 Smoothness of exponential map

First, we show the relation between the differential of exponential map  $d \exp_x(v_x)$  and the Jacobi field along the geodesic  $\exp_x(tv_x)$ .

**Lemma 14.4.12.** *Given a geodesic  $\gamma(t) = \exp_x(\frac{t}{\ell}v_x) \in \Pi$ , let  $\{X_i(t)\}_{i=1}^n$  be the parallel transport of some orthonormal frame along  $\gamma(t)$ . Then, for any  $w \in \mathbb{R}^n$ , we have that*

$$d \exp_x(v_x) \left( \sum w_i X_i(0) \right) = \sum_i \psi_w(\ell) X_i(\ell)$$

where  $\psi_w$  satisfies the Jacobi field along  $\gamma(t)$ , i.e.,

$$\begin{aligned} \frac{d^2}{dt^2} \psi_w(t) + R(t) \psi_w(t) &= 0, \\ \frac{d}{dt} \psi_w(0) &= w/\ell, \\ \psi_w(t) &= 0, \end{aligned} \tag{14.11}$$

and  $R(t)$  is a matrix given by  $R_{ij}(t) = \langle R(X_i(t), \gamma'(t))\gamma'(t), X_j(t) \rangle$ .

*Proof.* We want to compute  $d \exp_x(v_x)(w)$  for some  $w \in T_y M$ . By definition, we have that

$$d \exp_x(v_x)(w) = \frac{d}{ds} \gamma(t, s)|_{t=\ell, s=0}$$

where  $\gamma(t, s) = \exp_x(tv_x/\ell + sw)$ . Let  $\eta_w(t) = \frac{d}{ds} \gamma(t, s)|_{s=0}$  be a Jacobi field given by the formula

$$\begin{aligned} D_t D_t \eta_w + R(\eta_w, \gamma'(t))\gamma'(t) &= 0, \text{ for } 0 \leq t \leq \ell \\ D_t \eta_w(0) &= w/\ell, \\ \dot{\eta}_w(0) &= 0. \end{aligned}$$

Recall that the parallel transport of an orthonormal frame remains orthonormal because  $\frac{d}{dt} \langle X_i, X_j \rangle = \langle D_t X_i, X_j \rangle + \langle X_i, D_t X_j \rangle = 0$ . Since  $X_i(t)$  is an orthonormal basis at  $T_{\gamma(t)} M$ , we can write

$$\eta_w(t) = \sum_i \psi_i(t) X_i(t).$$

Since  $\eta_w$  satisfies the ODE above, we have

$$D_t D_t \sum_i \psi_i(t) X_i(t) + R \left( \sum_i \psi_i(t) X_i(t), \gamma'(t) \right) \gamma'(t) = 0.$$

Since  $X_i(t)$  is a parallel transport, we have that  $D_t D_t (\psi_i(t) X_i(t)) = \frac{d^2 \psi_i(t)}{dt^2} X_i(t)$  and hence

$$\sum_i \frac{d^2}{dt^2} \psi_i(t) X_i(t) + \sum_i \psi_i(t) R(X_i(t), \gamma'(t)) \gamma'(t) = 0$$

Let  $R_{ij}(t) = \langle R(X_i(t), \gamma'(t))\gamma'(t), X_j(t) \rangle$ . Since  $R(t)$  is a symmetric matrix (Fact 14.2.3), we have

$$\frac{d^2}{dt^2} \psi(t) + R(t) \psi(t) = 0.$$

Hence, we have

$$d \exp_x(v_x)(w) = \psi(\ell).$$

□

Next, we have a lemma about determinant.

**Lemma 14.4.13.** *Suppose that  $E$  is a matrix (not necessarily symmetric) with  $\|E\|_F \leq \frac{1}{4}$ , we have*

$$|\log \det(I + E) - \text{Tr}E| \leq \|E\|_F^2.$$

*Proof.* Let  $f(t) = \log \det(I + tE)$ . Then, by Jacobi's formula, we have

$$\begin{aligned} f'(t) &= \text{Tr}((I + tE)^{-1}E), \\ f''(t) &= -\text{Tr}((I + tE)^{-1}E(I + tE)^{-1}E). \end{aligned}$$

Therefore, we have

$$\begin{aligned} f(1) &= f(0) + f'(0) + \int_0^1 (1-s)f''(s)ds \\ &= \text{Tr}(E) + \int_0^1 (1-s)f''(s)ds. \end{aligned}$$

We note that

$$\begin{aligned} |f''(s)| &= |\text{Tr}((I + tE)^{-1}E(I + tE)^{-1}E)| \\ &\leq |\text{Tr}(E^T((I + tE)^{-1})^T(I + tE)^{-1}E)|. \end{aligned}$$

Since  $\|E\|_F \leq \frac{1}{4}$ , we have  $\|E\|_2 \leq \frac{1}{4}$  and hence  $\|(I + tE)^{-1}\|_2 \leq \frac{4}{3}$  for all  $0 \leq t \leq 1$ . Therefore, we have

$$|f''(s)| \leq 2|\text{Tr}(E^T E)| = 2\|E\|_F^2.$$

Therefore, we have

$$|f(1) - \text{Tr}E| \leq \|E\|_F^2.$$

□

Using lemma 14.4.13 and lemma 14.4.21, we have the following:

**Lemma 14.4.14.** *Given a geodesic walk  $\gamma(t) = \exp_x(\frac{t}{\ell}v_x) \in \Pi$  with step size  $h$  satisfying  $0 < h \leq \frac{1}{nR_1}$ . We have that  $d\exp_x(v_x)$  is invertible and*

$$\left| \log \det(d\exp_x(v_x)) - \int_0^\ell \frac{s(\ell-s)}{\ell} \text{Ric}(\gamma'(s))ds \right| \leq \frac{(nhR_1)^2}{6}. \quad (14.12)$$

Therefore, we have that

$$|\log \det(d\exp_x(v_x)) - \log \det(d\exp_y(v_y))| \leq \frac{(nhR_1)^2}{3}.$$

*Proof.* Let  $\Psi$  be the solution of the ODE  $\Psi''(t) + R(t)\Psi(t) = 0$ ,  $\Psi'(0) = I/\ell$  and  $\Psi(0) = 0$ . We know that  $\|R(t)\|_2 \leq \|R(t)\|_F \leq R_1$  for all  $0 \leq t \leq \ell$ . Hence, Lemma 14.4.21 shows that

$$\begin{aligned} \|\Psi(t) - \frac{t}{\ell}I\|_F &\leq \max_{0 \leq s \leq \ell} \|R(s)\|_F \left( \frac{t^3}{5} \|I/\ell\|_2 \right) \\ &\leq \frac{1}{5} R_1 n h \leq \frac{1}{4} \end{aligned} \quad (14.13)$$

By the Lemma 14.4.13, we have that

$$|\log \det(\Psi(\ell)) - \text{Tr}(\Psi(\ell) - I)| \leq \left( \frac{1}{5} R_1 n h \right)^2.$$

Now, we need to estimate  $\text{Tr}(\Psi(\ell) - I)$ . Note that

$$\begin{aligned}\Psi(\ell) &= \Psi(0) + \Psi'(0)\ell + \int_0^\ell (\ell - s)R(s)\Psi(s)ds \\ &= I + \int_0^\ell (\ell - s)R(s)\Psi(s)ds.\end{aligned}$$

Hence, we have

$$\Psi(\ell) - I - \int_0^\ell \frac{s(\ell - s)}{\ell} R(s)ds = \int_0^\ell (\ell - s)R(s) \left( \Psi(s) - \frac{s}{\ell}I \right) ds.$$

Using Lemma (14.13), we have

$$\begin{aligned}\left| \text{Tr} \left( \Psi(\ell) - I - \int_0^\ell \frac{s(\ell - s)}{\ell} R(s)ds \right) \right| &\leq \int_0^\ell (\ell - s) \left| \text{Tr} R(s) \left( \Psi(s) - \frac{s}{\ell}I \right) \right| ds \\ &\leq \int_0^\ell (\ell - s) \|R(s)\|_F \|\Psi(s) - \frac{s}{\ell}I\| ds \\ &\leq \frac{\ell^2}{2} \cdot R_1 \cdot \frac{1}{5} R_1 nh \leq \frac{(nhR_1)^2}{10}.\end{aligned}$$

Lemma 14.4.12 shows that  $d \exp_x(v_x) (\sum w_i X_i(0)) = \sum_i \psi_w(\ell)_i X(\ell) = \sum_i (\Psi(\ell)w)_i X_i(\ell)$ . Since  $\{X_i(t)\}_{i=1}^n$  are orthonormal, this shows that

$$d \exp_x(v_x) = X(\ell)\Psi(\ell)X(0)^T$$

where  $X$  is the matrix  $[X_1, X_2, \dots, X_n]$ . Since  $\|\Psi(\ell) - I\|_F \leq \frac{1}{5}$  (14.13), we have that  $\Psi(\ell)$  is invertible and so is  $d \exp_x(v_x)$ .

Since  $X(\ell)$  and  $X(0)$  are orthonormal, we have that

$$\log \det (d \exp_x(v_x)) = \log \det \Psi(\ell).$$

Therefore, this gives

$$\left| \log \det (d \exp_x(v_x)) - \int_0^\ell \frac{s(\ell - s)}{\ell} \text{Tr} R(s)ds \right| \leq \frac{(nhR_1)^2}{6}.$$

By the definition of Ricci curvature and Fact 14.2.3, we have that

$$\begin{aligned}\text{Tr} R(s) &= \sum_i \langle R(X_i(s), \gamma'(s))\gamma'(s), X_j(s) \rangle \\ &= \text{Ric}(\gamma'(s)).\end{aligned}$$

This gives the result (14.12).

Since the geodesic  $\exp_x(tv_x)$  is the same as  $\exp_y(tv_y)$  except for swapping the parameterization, and since  $\int_0^\ell \frac{s(\ell - s)}{\ell} \text{Tr} R(s)ds$  is invariant under this swap, (14.12) implies that  $\log \det (d \exp_x(v_x))$  is close to  $\log \det (d \exp_y(v_y))$ .  $\square$

## ■ 14.4.6 Implementation

Here, we explain in high level how to implement the geodesic walk in general via an efficient algorithm for approximately solving ODEs. Note that to implement the step, we need to compute the geodesic and compute the probability  $p(x \xrightarrow{w} y)$  and  $p(y \xrightarrow{w'} x)$ . From the formula (14.10), we see that they

involves the term  $d\exp_x(w)$ . Lemma 14.4.12 shows that  $d\exp_x(w)$  can be computed by Jacobi field. Also, Lemma 14.4.12 shows that Jacobi field can be computed by a ODE if it is written in the orthonormal frame systems. Therefore, to implement the geodesic walk, we need to compute geodesic, parallel transport and Jacobi field (See Algo 49). All of these are ODEs and can be solved using the collocation method. In the later sections, we will see how to collocation method can indeed solves these ODEs in matrix multiplication time.

---

**Algorithm 49: Geodesic Walk (Detailed)**


---

Pick a Gaussian random vector  $w \sim N_x(0, 1)$ , i.e.  $\mathbf{E}_w \|w\|_x^2 = 1$ .  
 /\* Compute  $y = \exp_x(\sqrt{h}w + \frac{h}{2}\mu(x))$  where  $\mu(x)$  is given by (14.2). \*/  
 /\* Compute a corresponding  $w'$  s.t.  $x = \exp_y(\sqrt{h}w' + \frac{h}{2}\mu(y))$ . \*/  
 Generate a random direction  $d = \sqrt{h}w + \frac{h}{2}\mu(x)$  where  $\mu(x)$  is given by (14.2).  
 Solve the geodesic equation  $D_{\gamma'}\gamma' = 0$  with  $\gamma(0) = x$  and  $\gamma'(0) = d$  using collocation method (Sec 14.6).  
 Set  $y = \gamma(1)$  and  $w' = -\gamma'(1)$ .  
  
 /\* Compute the probability  $p(x \xrightarrow{w} y)$  of going from  $x$  to  $y$  using the step  $w$ . \*/  
 Pick an orthonormal frame  $\{X_i\}_{i=1}^n$  at  $x$ .  
 Compute the parallel transport of  $\{X_i\}_{i=1}^n$  along  $\gamma(t)$  using collocation method.  
 Compute  $d\exp_x(w)$  via solve the Jacobi field equation (14.11) (in the coordinate systems  $\{X_i(t)\}$ ).  
 Set  $p(x \xrightarrow{w} y) = \det(d\exp_x(w))^{-1} \sqrt{\frac{\det(g(y))}{(2\pi h)^n}} \exp\left(-\frac{1}{2}\left\|\frac{w - \frac{h}{2}\mu(x)}{\sqrt{h}}\right\|_x^2\right)$ .  
 Compute  $p(y \xrightarrow{w'} x)$  similarly.  
  
 With probability  $\min\left(1, \frac{p(y \xrightarrow{w'} x)}{p(x \xrightarrow{w} y)}\right)$ , go to  $y$ ; otherwise, stay at  $x$ .

---

### ■ 14.4.7 Smoothness of $p_x$

Here we prove Lemma 14.4.6. Recall that the probability density of going from  $x$  to  $y$  is given by the following formula:

$$p_x(y) = \sum_{v_x: \exp_x(v_x)=y} \det(d\exp_x(v_x))^{-1} \sqrt{\frac{\det(g(y))}{(2\pi h)^n}} \exp\left(-\frac{1}{2}\left\|\frac{v_x - \frac{h}{2}\mu(x)}{\sqrt{h}}\right\|_x^2\right)$$

To simplify the calculation, we apply Lemma 14.4.14 and consider the following estimate of  $p_x(y)$  instead

$$\tilde{p}_x(y) = \sum_{v_x: \exp_x(v_x)=y} \sqrt{\frac{\det(g(y))}{(2\pi h)^n}} \exp\left(-\int_0^\ell \frac{t(\ell-t)}{\ell} \text{Ric}(\gamma'(t)) dt - \frac{1}{2}\left\|\frac{v_x - \frac{h}{2}\mu(x)}{\sqrt{h}}\right\|_x^2\right)$$

where  $\gamma(t)$  be the geodesic from  $\gamma(0) = x$  to  $\gamma(\ell) = y$ . Lemma 14.4.14 shows that  $|\log(\tilde{p}_x(y)/p_x(y))| \leq \frac{(nhR_1)^2}{6}$  and hence it suffices to prove the smoothness of  $\tilde{p}_x(y)$ .

Let  $c(s)$  be an unit speed geodesic going from  $x$  to some point  $z$  infinitesimally close to  $x$ . For (almost) any  $v_x$  such that  $\exp_x(v_x) = y$ , Lemma 14.4.24 shows that there is an unique vector field  $v(s)$

on  $c(s)$  such that  $\exp_{c(s)}(v(s)) = y$  and  $v(0) = v_x$ . Now, we define

$$\zeta(s) = \sqrt{\frac{\det(g(y))}{(2\pi h)^n}} \exp\left(-\int_0^\ell \frac{t(\ell-t)}{\ell} \text{Ric}(\gamma'_s(t)) dt - \frac{1}{2} \left\| \frac{v(s) - \frac{h}{2}\mu(c(s))}{\sqrt{h}} \right\|_{c(s)}^2\right)$$

where  $\exp_{c(s)}(v(s)) = y$  and  $\gamma_s(t) = \exp_{c(s)}(\frac{t}{\ell}v(s))$ . Then, we have that

$$\tilde{p}_{c(s)}(y) = \sum_{v(s): \exp_{c(s)}(v(s))=y} \zeta(s)$$

and hence

$$\frac{d}{ds} \tilde{p}_{c(s)}(y) = \sum_{v(s): \exp_{c(s)}(v(s))=y} \left( -\int_0^\ell \frac{t(\ell-t)}{\ell} \frac{d}{ds} \text{Ric}(\gamma'_s(t)) dt - \frac{1}{2} \frac{d}{ds} \left\| \frac{v(s) - \frac{h}{2}\mu(c(s))}{\sqrt{h}} \right\|_{c(s)}^2 \right) \zeta(s).$$

Hence, it suffices to bound the terms in the bracket.

In Section 14.4.7.1, we analyze  $\frac{d}{ds} \text{Ric}(\gamma'_s(t))$  and prove that

$$\left| \int_0^\ell \frac{t(\ell-t)}{\ell} \frac{d}{ds} \text{Ric}(\gamma'_s(t)) dt \right| \leq 2nhR_2.$$

In Section 14.4.7.2, we analyze  $\frac{d}{ds} \left\| \frac{v(s) - \frac{h}{2}\mu(c(s))}{\sqrt{h}} \right\|_{c(s)}^2$  and prove that

$$\mathbf{E}_\gamma \left| \frac{d}{ds} \left\| \frac{v(s) - \frac{h}{2}\mu(c(s))}{\sqrt{h}} \right\|_{c(s)}^2 \right|_{s=0} = O\left(D_2\sqrt{h} + \frac{1}{\sqrt{h}} + D_0\right).$$

Therefore, we have that

$$\left| \frac{d}{ds} \tilde{p}_c(y) \right| = O\left(nhR_2 + \sqrt{h}D_2 + \frac{1}{\sqrt{h}} + nhR_1D_0\right) \tilde{p}_c(y).$$

Using  $|\log(\tilde{p}_x(y)/p_x(y))| \leq \frac{(nhR_1)^2}{6}$ , we have our result.

**Theorem 14.4.15.** *We have that*

$$d_{TV}(p_x, p_y) = O\left(nhR_2 + \sqrt{h}D_2 + \frac{1}{\sqrt{h}} + nhR_1D_0\right) d(x, y) + O(nhR_1)^2.$$

#### 14.4.7.1 Smoothness of exponential map

**Definition 14.4.16.** Given a manifold  $M$  and a geodesic walk  $\gamma(t)$  on  $M_L$  with step size  $h$ . Let  $R_2$  be a constant depending on the manifold  $M$  and the step size  $h$  such that for any  $t$  such that  $0 \leq t \leq \ell$ , any curve  $\alpha(s)$  starting from  $\gamma(t)$  and any vector field  $u(s)$  on  $\alpha(s)$ , we have that

$$\left| \frac{d}{ds} \text{Ric}(u(s)) \Big|_{s=0} \right| \leq (\|\alpha'(0)\| + \|D_s u(0)\|) R_2.$$

**Lemma 14.4.17.** *For  $\frac{1}{n} < h \leq \frac{1}{2nR_1}$ , we have*

$$\left| \int_0^\ell \frac{t(\ell-t)}{\ell} \frac{d}{ds} \text{Ric}(\gamma'_s(t)) dt \right| \leq 2nhR_2$$



where  $\gamma_s$  is as defined in the beginning of Section 14.4.7.

*Proof.* Fix any  $t$  such that  $0 \leq t \leq \ell$  and let  $\alpha(s) = \gamma_s(t)$  and  $u(s) = \gamma'_s(t)$  be a vector field on  $\alpha$ . By Definition 14.4.16, we have that

$$\left| \frac{d}{ds} Ric(u(s))|_{s=0} \right| \leq (\|\alpha'(0)\| + \|D_s u(0)\|) R_2.$$

Recall that  $\gamma_s = \exp_{c(s)}(\frac{t}{\ell}v(s))$  and hence

$$\begin{aligned} \frac{d\alpha}{ds}(0) &= \frac{\partial}{\partial s} \exp_{c(s)}(\frac{t}{\ell}v(s))|_{s=0}, \\ D_s u(0) &= D_s \frac{\partial}{\partial t} \exp_{c(s)}(\frac{t}{\ell}v(s))|_{s=0} \\ &= D_t \frac{\partial}{\partial s} \exp_{c(s)}(\frac{t}{\ell}v(s))|_{s=0}. \end{aligned}$$

Let  $\psi(t) = \frac{\partial}{\partial s} \exp_{c(s)}(\frac{t}{\ell}v(s))|_{s=0}$ . Since  $\exp_{c(s)}(\frac{t}{\ell}v(s))$  is a geodesic for each  $s$ ,  $\psi(t)$  satisfies the Jacobi field equation. It is easier to calculate using orthogonal frame. So, we pick an arbitrary orthogonal frame along the curve  $\gamma(t)$  and let  $\phi(t)$  be  $\psi(t)$  represented in that orthogonal frame. Then, we know that

$$\begin{aligned} \frac{d^2}{dt^2} \phi(t) + R(\phi(t), \gamma'(t))\gamma'(t) &= 0, \\ \phi(0) &= \alpha, \\ \phi'(0) &= \frac{\beta}{\ell} \end{aligned}$$

for some  $\alpha$  and  $\beta$ . Since  $c$  is a unit speed geodesic, we have that  $\alpha = \phi(0) = \frac{\partial c}{\partial s}(0)$  has  $\ell_2$  norm 1. Next, we note that  $\beta = \ell\phi'(0) = D_s v(s)$ . By lemma 14.4.22, we have that  $\|\beta\|_2 \leq \frac{3}{2}$ .

Therefore, Lemma 14.4.20 shows that

$$\|\phi(t) - \alpha - \frac{\beta}{\ell}t\|_2 \leq 2R_1\ell^2$$

for all  $0 \leq t \leq \frac{1}{\sqrt{R_1}}$ . Hence, we have  $\|\phi(t)\|_2 \leq 4$  for all  $0 \leq t \leq \frac{1}{\sqrt{R_1}}$ . By Lemma 14.4.20 again, we have

$$\|\phi'(t) - \beta\|_2 \leq 3R_1\ell$$

for all  $0 \leq t \leq \frac{1}{\sqrt{R_1}}$ . Hence, we have  $\|\phi'(t)\|_2 \leq 4$  for all  $0 \leq t \leq \frac{1}{\sqrt{R_1}}$ .

Therefore, we have

$$\begin{aligned} \left| \frac{d}{ds} Ric(u(s))|_{s=0} \right| &\leq (\|\alpha'(0)\| + \|D_s u(0)\|) R_2 \\ &= (\|\phi(t)\|_2 + \|\phi'(t)\|_2) R_2 \\ &\leq 8R_2. \end{aligned}$$

This gives the result. □

### 14.4.7.2 One-step overlap

We parametrize the geodesic from  $x$  to  $y$  as  $\gamma(s)$ . For a fixed  $z$ , we have  $z = \exp_x(d_x) = \exp_{\gamma(s)}(d_{\gamma(s)}) = \exp_{\gamma(s)}(d(s))$ .

**Lemma 14.4.18.** *We have that*

$$\mathbf{E}_\gamma \left| \frac{d}{ds} \left\| \frac{v(s) - \frac{h}{2}\mu(c(s))}{\sqrt{h}} \right\|_{c(s)}^2 \right|_{s=0} = O \left( D_2\sqrt{h} + \frac{1}{\sqrt{h}} + D_0 \right).$$

where  $c(s)$  is as defined in the beginning of Section 14.4.7.

*Proof.* Note that

$$\frac{d}{ds} \left\| \frac{v(s) - \frac{h}{2}\mu(c(s))}{\sqrt{h}} \right\|_{c(s)}^2 \Big|_{s=0} = \frac{2}{h} \left\langle D_s v|_{s=0} - \frac{h}{2} D_s \mu|_{s=0}, v - \frac{h}{2}\mu \right\rangle. \quad (14.14)$$

Since  $v - \frac{h}{2}\mu$  is a random Gaussian vector  $N(0, nhI)$  in  $T_cM$  independent of  $D_s\mu$ , we have that

$$\mathbf{E} \left| \left\langle \frac{h}{2} D_s \mu|_{s=0}, v - \frac{h}{2}\mu \right\rangle \right| \leq O \left( \frac{\| \frac{h}{2} D_s \mu|_{s=0} \| \sqrt{nh}}{\sqrt{n}} \right) \leq O \left( \frac{h D_2 \sqrt{nh}}{2\sqrt{n}} \right). \quad (14.15)$$

By Lemma 14.4.22, we know that  $D_s v|_{s=0} = -c' + \zeta$  where  $\zeta \perp v(0)$  and  $\|\zeta\| \leq \frac{3}{2}nhR_1\|c'\|$ . Since  $v - \frac{h}{2}\mu$  is a random Gaussian vector independent of  $c'$ , we have that

$$\begin{aligned} \mathbf{E} \left| \left\langle D_s v, v - \frac{h}{2}\mu \right\rangle \right| &\leq \mathbf{E} \left| \left\langle c', v - \frac{h}{2}\mu \right\rangle \right| + \mathbf{E} \frac{h}{2} |\langle \zeta, \mu \rangle| \\ &\leq O \left( \frac{\|c'\| \sqrt{nh}}{\sqrt{n}} \right) + nh^2 R_1 \|c'\| \|\mu\| \\ &\leq O \left( \sqrt{h} + nh^2 R_1 D_0 \right). \end{aligned} \quad (14.16)$$

Combining the bounds (14.14), (14.15) and (14.16), we have that

$$\mathbf{E} \left| \frac{d}{ds} \left\| \frac{v(s) - \frac{h}{2}\mu(c(s))}{\sqrt{h}} \right\|_{c(s)}^2 \right|_{s=0} \leq O \left( D_2\sqrt{h} + \frac{1}{\sqrt{h}} + nhR_1 D_0 \right).$$

□

### ■ 14.4.8 Approximate Solution of Jacobi Field

Let  $\gamma(t)$  be a geodesic and  $\{X_i(t)\}_{i=1}^n$  be the parallel transport of some orthonormal frame along  $\gamma(t)$ . As we demonstrated in the proof of Lemma 14.4.12, Jacobi fields are solutions the following matrix ODE:

$$\begin{aligned} \frac{d^2}{dt^2} \psi(t) + R(t)\psi(t) &= 0, \\ \frac{d}{dt} \psi(0) &= \beta, \\ \psi(0) &= \alpha \end{aligned} \quad (14.17)$$

where  $\psi(t), \alpha, \beta \in \mathbb{R}^n$  and  $R(t)$  is a matrix given by  $R_{ji}(t) = \langle R(X_i(t), \gamma'(t))\gamma'(t), X_j(t) \rangle$ . Recall that  $R(t)$  is a symmetric matrix (Fact 14.2.3).

In this section, we show give an approximate solution of the the Jacobi equation (14.17). First, we give a basic bound on the solution in terms of hyperbolic sine and cosine functions.

**Lemma 14.4.19.** *Let  $\psi$  be the solution of (14.17). Suppose that  $\|R(t)\|_2 \leq \lambda$ . Then, we have that*

$$\|\psi(t)\|_2 \leq \|\alpha\|_2 \cosh(\sqrt{\lambda}t) + \frac{\|\beta\|_2}{\sqrt{\lambda}} \sinh(\sqrt{\lambda}t)$$

for all  $t \geq 0$ .

*Proof.* Let  $a(t) = \|\psi(t)\|_2$ . Using  $\|R(t)\|_2 \leq \lambda$  and  $\frac{d^2}{dt^2}\psi(t) + R(t)\psi(t) = 0$ , we have that

$$a''(t) \leq \lambda a(t).$$

Since  $a(0) = \|\alpha\|_2$  and  $a'(0) \leq \|\beta\|_2$ , for all  $t \geq 0$ , we have that  $a(t) \leq b(t)$  where  $b(t)$  is the solution of the  $b''(t) = \lambda b(t)$ ,  $b(0) = \|\alpha\|_2$ ,  $b'(0) = \|\beta\|_2$ . Solving these equations, we have

$$a(t) \leq b(t) = \|\alpha\|_2 \cosh(\sqrt{\lambda}t) + \frac{\|\beta\|_2}{\sqrt{\lambda}} \sinh(\sqrt{\lambda}t)$$

for all  $t \geq 0$ . □

Next, we give an approximate solution of (14.17).

**Lemma 14.4.20.** *Let  $\psi$  be the solution of (14.17). Suppose that  $\|R(t)\|_2 \leq \lambda$ . For any  $0 \leq t \leq \frac{1}{\sqrt{\lambda}}$ , we have that*

$$\|\psi(t) - \alpha - \beta t\|_2 \leq \lambda t^2 \|\alpha\|_2 + \frac{\lambda t^3}{5} \|\beta\|_2 \text{ and } \|\psi'(t) - \beta\|_2 \leq 2\lambda t \|\alpha\|_2 + \frac{3\lambda t^2}{5} \|\beta\|_2.$$

*Proof.* Note that

$$\begin{aligned} \psi(t) &= \psi(0) + \psi'(0)t + \int_0^t (t-s)\psi''(s)ds \\ &= \alpha + \beta t - \int_0^t (t-s)R(s)\psi(s)ds. \end{aligned}$$

Using Lemma 14.4.19 and  $\|R(t)\|_2 \leq \lambda$ , we have that

$$\begin{aligned} \|\psi(t) - \alpha - \beta t\|_2 &\leq \lambda \int_0^t (t-s)\|\psi(s)\|_2 ds \\ &\leq \lambda \int_0^t (t-s) \left( \|\alpha\|_2 \cosh(\sqrt{\lambda}s) + \frac{\|\beta\|_2}{\sqrt{\lambda}} \sinh(\sqrt{\lambda}s) \right) ds \\ &= \|\alpha\|_2 \left( \cosh(\sqrt{\lambda}t) - 1 \right) + \frac{\|\beta\|_2}{\sqrt{\lambda}} \left( \sinh(\sqrt{\lambda}t) - \sqrt{\lambda}t \right). \end{aligned}$$

Since  $0 \leq t \leq \frac{1}{\sqrt{\lambda}}$ , we have that  $|\cosh(\sqrt{\lambda}t) - 1| \leq \lambda t^2$  and  $|\sinh(\sqrt{\lambda}t) - \sqrt{\lambda}t| \leq \frac{\lambda^{3/2}t^3}{5}$ . This gives the result.

Similarly, we have that  $\psi'(t) = \psi'(0) + \int_0^t \psi''(s)ds = \psi'(0) + \int_0^t R(s)\psi(s)ds$  and hence

$$\begin{aligned} \|\psi'(t) - \beta\|_2 &\leq \lambda \int_0^t \left( \|\alpha\|_2 \cosh(\sqrt{\lambda}s) + \frac{\|\beta\|_2}{\sqrt{\lambda}} \sinh(\sqrt{\lambda}s) \right) ds \\ &\leq \sqrt{\lambda} \|\alpha\|_2 \sinh(\sqrt{\lambda}t) + \|\beta\|_2 \left( \cosh(\sqrt{\lambda}t) - 1 \right) \\ &\leq 2\lambda t \|\alpha\|_2 + \frac{3}{5} \lambda t^2 \|\beta\|_2 \end{aligned}$$

□

The following is a matrix version of the above result.

**Lemma 14.4.21.** *Let  $\psi$  be the solution of*

$$\begin{aligned}\frac{d^2}{dt^2}\Psi(t) + R(t)\Psi(t) &= 0, \\ \frac{d}{dt}\Psi(0) &= B, \\ \Psi(0) &= A\end{aligned}$$

where  $\Psi, A$  and  $B$  are matrices with a compatible size. Suppose that  $\|R(t)\|_2 \leq \lambda$ . For any  $0 \leq t \leq \frac{1}{\sqrt{\lambda}}$ , we have that

$$\|\Psi(t) - A - Bt\|_F \leq \max_{0 \leq s \leq t} \|R(s)\|_F \left( t^2 \|A\|_2 + \frac{t^3}{5} \|B\|_2 \right).$$

*Proof.* Note that  $\Psi(t)x$  is the solution of (14.17) with  $\beta = Bx$  and  $\alpha = Ax$ . Therefore, Lemma 14.4.19 shows that

$$\|\Psi(t)x\|_2 \leq \left( \|A\|_2 \cosh(\sqrt{\lambda}t) + \frac{\|B\|_2}{\sqrt{\lambda}} \sinh(\sqrt{\lambda}t) \right) \|x\|_2$$

for all  $x$ . Therefore, we have that

$$\Psi(t)^T \Psi(t) \preceq \left( \|A\|_2 \cosh(\sqrt{\lambda}t) + \frac{\|B\|_2}{\sqrt{\lambda}} \sinh(\sqrt{\lambda}t) \right)^2 I.$$

Hence, we have that

$$\Psi(t)\Psi(t)^T \preceq \left( \|A\|_2 \cosh(\sqrt{\lambda}t) + \frac{\|B\|_2}{\sqrt{\lambda}} \sinh(\sqrt{\lambda}t) \right)^2 I.$$

Using this, we have that

$$\begin{aligned}\|\Psi(t) - A - Bt\|_F &= \left\| \int_0^t (t-s)R(s)\Psi(s)ds \right\|_F \\ &\leq \int_0^t (t-s) \sqrt{\text{Tr} \Psi^T(s)R^T(s)R(s)\Psi(s)} ds \\ &= \int_0^t (t-s) \sqrt{\text{Tr} R(s)\Psi(s)\Psi^T(s)R^T(s)} ds \\ &\leq \int_0^t (t-s) \left( \|A\|_2 \cosh(\sqrt{\lambda}s) + \frac{\|B\|_2}{\sqrt{\lambda}} \sinh(\sqrt{\lambda}s) \right) \|R(s)\|_F ds. \\ &\leq \max_{0 \leq s \leq t} \|R(s)\|_F \int_0^t (t-s) \left( \|A\|_2 \cosh(\sqrt{\lambda}s) + \frac{\|B\|_2}{\sqrt{\lambda}} \sinh(\sqrt{\lambda}s) \right) ds \\ &= \max_{0 \leq s \leq t} \|R(s)\|_F \left( \|A\|_2 \left( \cosh(\sqrt{\lambda}t) - 1 \right) + \frac{\|B\|_2}{\sqrt{\lambda}} \left( \sinh(\sqrt{\lambda}t) - \sqrt{\lambda}t \right) \right).\end{aligned}$$

Since  $0 \leq t \leq \frac{1}{\sqrt{\lambda}}$ , we have that  $|\cosh(\sqrt{\lambda}t) - 1| \leq \lambda t^2$  and  $|\sinh(\sqrt{\lambda}t) - \sqrt{\lambda}t| \leq \frac{\lambda^{3/2}t^3}{5}$ . This gives the result.  $\square$

### ■ 14.4.9 Almost One-to-One Correspondence of Geodesics

We do not know if every pair  $x, y \in M$  has a unique geodesic connecting  $x$  and  $y$ . Due to this, the probability density  $p_x$  at  $y$  on  $\mathcal{M}$  can be the sum over all possible geodesics connect  $x$  and  $y$ . The goal

of this section is to show there is a 1-1 map between geodesics paths connecting  $x$  to  $y$  as we move  $x$ .

**Lemma 14.4.22.** *Given a geodesic walk  $\gamma(t) = \exp_x(\frac{t}{\ell}v_x) \in \Pi$  with step size  $h$  satisfying  $\frac{1}{n} \leq h \leq \frac{1}{nR_1}$ . Let the end points  $x = \gamma(0)$  and  $y = \gamma(\ell)$ . There is an unique smooth invertible function  $v : U \subset M \rightarrow V \subset T_x M$  such that*

$$y = \exp_z(v(z))$$

for any  $z \in U$  where  $U$  is a neighborhood of  $x$  and  $V$  is a neighborhood of  $v_x = v(x)$ . Furthermore, for any  $\eta = \eta_1 + \alpha v_x$  with  $\eta_1 \perp v_x$  and scale  $\alpha$ , we have that

$$\nabla_\eta v(x) = -\eta + \zeta$$

where  $\|\zeta\| \leq \frac{3}{2}nhR_1\|\eta_1\| \leq \frac{3}{2}\|\eta_1\|$  and  $\zeta \perp v(x)$ . In particular, we have that  $\|\nabla_\eta v(x)\| \leq \frac{5}{2}\|\eta\|$ .

*Proof.* Consider the function  $f(z, w) = \exp_z(w)$ . It is continuously differentiable. From Lemma 14.4.14, the differential of  $w$  at  $(x, v_x)$  on the  $w$  variables, i.e.  $d\exp_x(v_x)$ , is invertible. Hence, the implicit function theorem shows that there is a open neighborhood  $U$  of  $x$  and a unique function  $v$  on  $U$  such that  $f(z, v(z)) = f(z, v_x)$ , i.e.  $y = \exp_x(v_x) = \exp_z(v(z))$ .

To compute  $\nabla_\eta v$ , let  $c(t, s) = \exp_{\gamma(s)}(\frac{t}{\ell}v(\gamma(s)))$  be a family of geodesics with the end points  $c(\ell, s) = \exp_{\gamma(s)}(v(\gamma(s))) = y$ . Note that  $\psi(t) \stackrel{\text{def}}{=} \frac{\partial c(t, s)}{\partial s} \Big|_{s=0}$  satisfies the Jacobi field equation

$$D_t^2 \psi(t) + R(\psi, \frac{\partial c}{\partial t}) \frac{\partial c}{\partial t} = 0.$$

Moreover, we know that

$$\begin{aligned} \psi(0) &= \frac{\partial c(0, s)}{\partial s} \Big|_{s=0} = \gamma'(0) = \eta, \\ D_t \psi(0) &= D_t \frac{\partial c(t, s)}{\partial s} \Big|_{t,s=0} = D_s \frac{\partial c(t, s)}{\partial t} \Big|_{t,s=0} = \frac{1}{\ell} \nabla_\eta v(x) \stackrel{\text{def}}{=} \frac{\chi}{\ell}, \\ \psi(\ell) &= \frac{\partial c(\ell, s)}{\partial s} \Big|_{s=0} = 0 \end{aligned}$$

where we used  $D_t \frac{\partial}{\partial s} = D_s \frac{\partial}{\partial t}$  (torsion free-ness, Fact 14.2.1).

From Fact 14.2.4, we know that  $\psi$  can be split into the tangential part  $\psi_1$  and the normal part  $\psi_2$ . For the tangential part, we know that

$$\begin{aligned} \psi_1(t) &= \left( \left\langle \psi(0), \frac{dc}{dt}(0, 0) \right\rangle_{c(0)} + \left\langle D_t \psi(0), \frac{dc}{dt}(0, 0) \right\rangle_{c(0)} t \right) \frac{dc}{dt}(t) \\ &= \left( \langle \eta, v_x \rangle_x + \left\langle \frac{\chi}{\ell}, v_x \right\rangle_x t \right) \frac{dc}{dt}(t) \\ &= \left( \alpha \|v_x\|^2 + \left\langle \frac{\chi}{\ell}, v_x \right\rangle_x t \right) \frac{dc}{dt}(t). \end{aligned}$$

Since  $\psi_1(\ell) = 0$ , we have that

$$\langle \chi, v_x \rangle = -\alpha \|v_x\|^2. \quad (14.18)$$

For the normal part  $\psi_2$ , it is easier to calculate using orthogonal frames. So, we pick an arbitrary orthogonal frame parallel transported along the curve  $c(t, 0)$  and let  $\bar{\psi}(t)$  be  $\psi_2(t)$  represented in that

orthogonal frame. Hence, we have that

$$\begin{aligned} \frac{d^2}{dt^2} \bar{\psi}(t) + R(\bar{\psi}(t), \frac{\partial c}{\partial t}) \frac{\partial c}{\partial t} &= 0, \\ \bar{\psi}(0) &= \eta_1, \\ \bar{\psi}'(0) &= \frac{\chi_1}{\ell}, \\ \bar{\psi}(\ell) &= 0 \end{aligned}$$

where  $\chi_1 = \chi - \left\langle \chi, \frac{v_x}{\|v_x\|} \right\rangle \frac{v_x}{\|v_x\|} = \chi + \alpha v_x$  (14.18). Using  $h \leq \frac{1}{nR_1}$ , Lemma 14.4.20 shows that

$$\begin{aligned} \|\eta_1 + \chi_1\|_2 = \|\bar{\psi}(\ell) - \eta_1 - \chi_1\|_2 &\leq R_1 \ell^2 \|\eta_1\|_2 + \frac{R_1 \ell^2}{5} \|\chi_1\|_2 \\ &\leq \frac{6}{5} R_1 \ell^2 \|\eta_1\|_2 + \frac{1}{5} \|\chi_1\|_2. \end{aligned}$$

Hence, we have that  $\|\eta_1 + \chi_1\|_2 \leq \frac{3}{2} R_1 \ell^2 \|\eta_1\|_2$ . Therefore, we have that

$$\begin{aligned} \nabla_\eta v(x) &= \chi \\ &= -\alpha v_x + \chi_1 \\ &= -\eta + (\eta_1 + \chi_1) \end{aligned}$$

where  $\|\chi_1 + \eta_1\|_2 \leq \frac{3}{2} \|\eta_1\|_2$ . Furthermore, we have that both  $\eta_1$  and  $\chi_1$  are orthogonal to  $\eta$ .  $\square$

*Remark.* If the above lemma holds without the assumption  $\gamma \in \Pi$ , this would imply uniqueness of geodesics.

The following lemma shows there is a 1-1 map between geodesics paths connecting  $x$  to  $y$  as we move  $x$ . When we move  $x$ , the geodesic  $\gamma$  from  $x$  to  $y$  no longer insides  $\Xi$ . To avoid this issue, we need to assume  $\langle R(X_i, \gamma') \gamma', X_j \rangle$  is mildly smooth in the following sense.

**Assumption 14.4.23.** *Given a manifold  $M$ . There is some constant  $C$  such that for any a curve  $\alpha(s) \in M$  and a vector field  $u(s)$  on  $\alpha(s)$ , let  $A(s)_{ij} = \langle R(X_i, u(s))u(s), X_j \rangle$  where  $\{X_i\}$  is any orthonormal frame at  $\alpha(s)$ , we have*

$$\left| \frac{d}{ds} \|A(s)\|_F \right| \leq (\|\alpha'(0)\| + \|D_s u(0)\|) n^C.$$

Note that this assumption implies  $R_1 = n^{\Theta(1)}$  by picking the curve is a singleton and the vector field  $u(s) = s \cdot u$ .

**Lemma 14.4.24.** *Given a geodesic  $\gamma(t) = \exp_x(\frac{t}{\ell} v_x)$  on  $M$  with step size  $h$  satisfying  $\frac{1}{n} \leq h \leq \frac{1}{2nR_1}$ , let  $c(s)$  be any geodesic starting at  $\gamma(0)$ . Let  $x = c(0) = \gamma(0)$  and  $z = c(1)$ . Suppose that the length of the geodesic is less than  $n^{-\Theta(1)}$ . Then, there is a unique vector field  $v$  on  $c$  such that*

$$y = \exp_{c(s)}(v(s)).$$

*This vector field is uniquely determined by the geodesic and any  $v(s)$  on this vector field.*

*Proof.* Let  $s_{\max}$  be the supremum of  $s$  such that  $v(s)$  can be defined continuously and  $y = \exp_{c(s)}(v(s))$ . Note that  $s_{\max} > 0$  as Lemma 14.4.22 shows that there is a neighborhood  $N$  at  $x$  and a vector field  $u$  on  $N$  such that for any  $z \in N$ , we have that

$$y = \exp_z(u(z)).$$

Note that the the vector field  $u$  is unique and hence the extension we found for  $v$  is unique. Also, we know that  $\|D_t v(t)\| = \|\nabla_{c'} u(x)\| \leq \frac{5}{2}L$  where  $L$  is the length of  $c$ . Therefore, by continuity,  $v(s_{\max})$  is well-defined.

If  $s_{\max} > 1$ , then we are done. Otherwise, we want to apply Lemma 14.4.22 at  $s_{\max}$ . However, the curve  $\exp_{c(s_{\max})}(\frac{t}{\ell}v(s_{\max}))$  may not lie in  $\Pi$ . Note that the only place Lemma 14.4.22 uses the assumption  $\gamma = \exp_x(\frac{t}{\ell}v_x) \in \Pi$  is to apply the bound  $R_1 = \sup_{\gamma \in \Pi, 0 \leq t \leq \ell} \|R(t)\|_F$ .

To bound  $\|R(t)\|_F$ , we note that  $\|D_t v(t)\| = \|\nabla_{c'} u(x)\| \leq \Theta(L)$  where  $L$  is the length of  $c$ . Hence, both the starting point and the direction of the geodesic move by  $\Theta(L)$ . Now, we can apply Lemma 14.4.19 to prove that along the geodesic, all the points move by  $\Theta(L)$  for  $0 \leq t \leq \ell$ . Then, we can apply Assumption 14.4.23 to conclude that  $\|R(t)\|_F$  only changes by at most  $\Theta(n^C L)$ . Therefore, if  $L \leq \frac{R_1}{n^C} = n^{-C}$ , then  $\|R(t)\|_F$  is still bounded by  $\Theta(R_1)$  and Lemma 14.4.22 can be applied at  $s_{\max} < 1$ . This is a contradiction. Hence,  $s_{\max} \geq 1$ . Its uniqueness follows from Lemma 14.4.22.  $\square$

## ■ 14.5 Logarithmic Barrier

For any polytope  $M = \{Ax > b\}$ , the logarithmic barrier function  $\phi(x)$  is defined as

$$\phi(x) = - \sum_{i=1}^m \log(a_i^T x - b_i).$$

We denote the Hessian manifold induced by the logarithmic barrier on  $M$  by  $M_L$ . The goal of this section is to analyze the geodesic walk on  $M_L$ .

In section 14.5.1, we give explicit formulas of various Riemannian geometry concepts on  $M_L$ . In Section 14.5.2, we explain the geodesic walk on  $M_L$ . In Sections 14.5.3 to 14.5.7, we bound the necessary constant required by Theorem 14.4.2. This shows that

**Theorem 14.5.1.** *The geodesic walk on  $M_L$  with step size  $h \leq \frac{c}{n^{3/4}}$  has mixing time is  $O(m/h)$  for some universal constant  $c$ .*

In the sections afterward, we show how to implement geodesic walk and calculate the rejection probability. To implement these, we apply the techniques developed in Section 14.6 to solve the corresponding ODEs. In this section, we first show that geodesic, parallel transport and Jacobi Field are complex analytic (Section 14.5.9), then we bound their radius of convergence (Section 14.5.9.1). This will help us bound the approximation error of the collocation method developed in Section 14.6.

**Theorem 14.5.2.** *We can implement one step of geodesic walk with step size  $h \leq \frac{c}{\sqrt{n}}$  in time  $O(mn^{\omega-1} \log^2(n))$  for a universal constant  $c$ .*

### ■ 14.5.1 Riemannian geometry on $M_L$

We use the following definitions throughout this section:

- $s_x = Ax - b$ ,  $S_x = \text{Diag}(s_x)$ .
- $A_x = S_x^{-1}A$ .
- $s_{x,v} = A_x v$ ,  $S_{x,v} = \text{Diag}(A_x v)$ .
- $P_x = A_x(A_x^T A_x)^{-1}A_x^T$ ,  $\sigma_x = \text{Diag}(P_x)$ .
- $\left(P_x^{(2)}\right)_{ij} = (P_x)_{ij}^2$ .

- For brevity (overloading notation), we define  $s_{\gamma'} = s_{\gamma,\gamma'}$ ,  $s_{\gamma''} = s_{\gamma,\gamma''}$ ,  $S_{\gamma'} = S_{\gamma,\gamma'}$  and  $S_{\gamma''} = S_{\gamma,\gamma''}$  for a curve  $\gamma(t)$ .

Since the Hessian manifold  $M_L$  is naturally embedded in  $\mathbb{R}^n$ , we identify  $T_x M_L$  by Euclidean coordinates unless otherwise stated. Therefore, we have that

$$\begin{aligned}\langle u, v \rangle_x &= u^T \nabla^2 \phi(x) v \\ &= u^T A_x^T A_x v.\end{aligned}$$

**Lemma 14.5.3.** *Let  $u(t)$  be a vector field defined on a curve  $\gamma(t)$  in  $M_L$ . Then, we have that*

$$\nabla_{\gamma'} u = \frac{du}{dt} - (A_{\gamma}^T A_{\gamma})^{-1} A_{\gamma}^T S_{\gamma'} s_{\gamma,u}.$$

In particular, the geodesic equation on  $M_L$  is given by

$$\gamma'' = (A_{\gamma}^T A_{\gamma})^{-1} A_{\gamma}^T s_{\gamma'}^2 \quad (14.19)$$

and the equation for parallel transport on a curve  $\gamma(t)$  is given by

$$\frac{d}{dt} v(t) = (A_{\gamma}^T A_{\gamma})^{-1} A_{\gamma}^T S_{\gamma'} A_{\gamma} v. \quad (14.20)$$

*Proof.* By Lemma 14.2.5, the Christoffel symbols respect to the Euclidean coordinates is given by

$$\begin{aligned}\Gamma_{ij}^k &= \frac{1}{2} \sum_l g^{kl} \phi_{ijl} \\ &= - \sum_z \sum_l e_k^T (A_x^T A_x)^{-1} e_l (A_x e_i)_z (A_x e_j)_z (A_x e_l)_z \\ &= -e_k^T (A_x^T A_x)^{-1} A_x^T ((A_x e_i) (A_x e_j)).\end{aligned}$$

Recall that the Levi-Civita connection is given by

$$\nabla_v u = \sum_{ik} v_i \frac{\partial u_k}{\partial x^i} e_k + \sum_{ijk} v_i u_j \Gamma_{ij}^k e_k.$$

Therefore, we have

$$\nabla_v u = \sum_{ik} v^i \frac{\partial u^k}{\partial x^i} e_k - (A_x^T A_x)^{-1} A_x^T S_{x,v} s_{x,u}.$$

Since  $U$  is a vector field defined on a curve  $\gamma(t)$ , we have that  $\sum_{ik} \gamma_i' \frac{\partial u_k}{\partial x^i} e_k = \frac{du}{dt}$ .

The geodesic equation follows from  $\nabla_{\gamma'} \gamma' = 0$  and the parallel transport equation follows from  $\nabla_{\gamma'} v = 0$ .  $\square$

**Lemma 14.5.4.** *Given  $u, v, w, z \in T_x M_L$ , the Riemann Curvature Tensor is given by*

$$\langle R(u, v)w, z \rangle = (s_{x,u} s_{x,w})^T A_x (A_x^T A_x)^{-1} A_x^T (s_{x,v} s_{x,z}) - (s_{x,u} s_{x,z})^T A_x (A_x^T A_x)^{-1} A_x^T (s_{x,v} s_{x,w})$$

and the Ricci curvature is given by

$$\text{Ric}(u) = s_{x,u}^T P_x^{(2)} s_{x,u} - \sigma_x^T P_x s_{x,u}^2.$$



*Proof.* By Lemma 14.2.5, we have that

$$\begin{aligned} \langle R(u, v)w, z \rangle &= \frac{1}{4} \sum_{pqijkl} g^{pq} (\phi_{jkp}\phi_{ilq} - \phi_{ikp}\phi_{jlq}) u^i v^j w^l z^k \\ &= \sum_{pq} g^{pq} (e_p^T A_x^T (s_{x,v} s_{x,z}) e_q^T A_x^T (s_{x,u} s_{x,w}) - e_p^T A_x^T (s_{x,u} s_{x,z}) e_q^T A_x^T (s_{x,v} s_{x,w})) \\ &= (s_{x,u} s_{x,w})^T A_x (A_x^T A_x)^{-1} A_x^T (s_{x,v} s_{x,z}) - (s_{x,u} s_{x,z})^T A_x (A_x^T A_x)^{-1} A_x^T (s_{x,v} s_{x,w}) \end{aligned}$$

and

$$\begin{aligned} Ric(u) &= \sum_{jl} g^{jl} \left\langle R(u, \frac{\partial}{\partial x^j}), \frac{\partial}{\partial x^l}, u \right\rangle \\ &= \sum_{jl} g^{jl} (s_{x,u} s_{x,e_l})^T A_x (A_x^T A_x)^{-1} A_x^T (s_{x,e_j} s_{x,u}) - (s_{x,u}^2)^T A_x (A_x^T A_x)^{-1} A_x^T (s_{x,e_j} s_{x,e_l}) \\ &= \text{Tr} S_{x,u} A_x (A_x^T A_x)^{-1} A_x^T S_{x,u} A_x (A_x^T A_x)^{-1} A_x - \text{Tr} \left( A_x^T \text{Diag}(P_x s_{x,u}^2) A_x (A_x^T A_x)^{-1} \right) \\ &= s_{x,u}^T P_x^{(2)} s_{x,u} - \sigma_x^T P_x s_{x,u}^2. \end{aligned}$$

□

**Lemma 14.5.5.** *Given a geodesic  $\gamma(t)$  on  $M_L$  and an orthogonal frame  $\{x_i\}_{i=1}^n$  on  $\gamma(t)$ . The Jacobi field equation (in the orthogonal frame coordinate) is given by*

$$\frac{d^2 u}{dt^2} + X^{-1} (A_\gamma^T A_\gamma)^{-1} (A_\gamma^T S_{\gamma'} P_\gamma S_{\gamma'} A_\gamma - A_\gamma^T \text{Diag}(P_\gamma s_{\gamma'}^2) A_\gamma) X u = 0$$

where  $X(t) = [x_1(t), x_2(t), \dots, x_n(t)]$ .

*Proof.* The equation for Jacobi field along is

$$D_t^2 v + R(v, \gamma') \gamma' = 0.$$

By Lemma 14.5.4, under Euclidean coordinates, we have

$$D_t^2 v + (A_\gamma^T A_\gamma)^{-1} (A_\gamma^T S_{\gamma'} P_\gamma S_{\gamma'} A_\gamma - A_\gamma^T \text{diag}(P_\gamma s_{\gamma'}^2) A_\gamma) v = 0$$

We write  $v$  in terms of the orthogonal frame, namely,  $v(t) = X(t)u(t)$  where  $u(t) \in \mathbb{R}^n$ . Then, we have that  $D_t^2 v = X(t) \frac{d^2 u(t)}{dt^2}$ . Hence, under the orthogonal frame coordinate, we have that

$$\frac{d^2 u}{dt^2} + X^{-1} (A_\gamma^T A_\gamma)^{-1} (A_\gamma^T S_{\gamma'} P_\gamma S_{\gamma'} A_\gamma - A_\gamma^T \text{Diag}(P_\gamma s_{\gamma'}^2) A_\gamma) X u = 0.$$

□

Now, we show the relation between the distance induced by the metric and the Hilbert metric.

**Lemma 14.5.6.** *For any  $x, y \in M_L$ , we have that*

$$\frac{d(x, y)}{d_H(x, y)} \leq \sqrt{m}.$$

Hence, we have that  $G_2 \leq \sqrt{m}$ .

*Proof.* First, we note that it suffices to prove that  $\frac{d(x, y)}{d_H(x, y)} \leq (1 + O(\varepsilon))\sqrt{m}$  for any  $x, y \in M_L$  with  $d(x, y) \leq \varepsilon$ . Then, one can run a limiting argument as follows. Let  $x_t = t \cdot x + (1 - t) \cdot y$ , then we have

that

$$d_H(x, y) = \lim_{n \rightarrow \infty} \sum_{i=0}^{n-1} d_H(x_{k/n}, x_{(k+1)/n}).$$

Since  $d_H(x_{k/n}, x_{(k+1)/n}) = O_{x,y}(\frac{1}{n})$ , we have that

$$\begin{aligned} d_H(x, y) &= \lim_{n \rightarrow \infty} \sum_{i=0}^{n-1} d_H(x_{k/n}, x_{(k+1)/n}) \\ &\geq \lim_{n \rightarrow \infty} \frac{(1 - O_{x,y}(\frac{1}{n}))}{\sqrt{m}} \sum_{i=0}^{n-1} d(x_{k/n}, x_{(k+1)/n}) \\ &= \frac{1}{\sqrt{m}} \lim_{n \rightarrow \infty} \sum_{i=0}^{n-1} d(x_{k/n}, x_{(k+1)/n}) \\ &\geq \frac{1}{\sqrt{m}} \lim_{n \rightarrow \infty} d(x, y) = \frac{d(x, y)}{\sqrt{m}}. \end{aligned}$$

Now, we can assume  $d(x, y) \leq \varepsilon$ . Let  $p$  and  $q$  are on the boundary of  $M_L$  such that  $p, x, y, q$  are on the straight line  $\overline{xy}$  and are in order. Without loss of generality, we assume  $p$  is closer to  $x$ . Then, we have that  $p \in M_L \cap (x - M_L)$ , equivalently, we have that  $|a_i^T p - a_i^T x| \leq a_i^T x - b_i$  for all  $i$  and hence  $\|A_x(p - x)\|_\infty \leq 1$ . Therefore, we have that

$$\|A_x(p - x)\|_2 \leq \sqrt{m} \|A_x(p - x)\|_\infty \leq \sqrt{m}.$$

Since  $p, x, y$  are on the same line, we have that

$$\frac{\|x - y\|_2 \|p - q\|_2}{\|p - x\|_2 \|y - q\|_2} \geq \frac{\|x - y\|_2}{\|p - x\|_2} = \frac{\|A_x(x - y)\|_2}{\|A_x(p - x)\|_2} \geq \frac{\|A_x(x - y)\|_2}{\sqrt{m}}.$$

Since  $d(x, y) < \varepsilon$ , Lemma 3.1 in [214] shows that

$$d(x, y) \leq -\log(1 - \|A_x(x - y)\|_2).$$

Hence, we have that  $\|A_x(x - y)\|_2 \geq (1 - O(\varepsilon))d(x, y)$  and hence

$$d_H(x, y) = \frac{\|x - y\|_2 \|p - q\|_2}{\|p - x\|_2 \|y - q\|_2} \geq (1 - O(\varepsilon)) \frac{d(x, y)}{\sqrt{m}}.$$

□

### ■ 14.5.2 Geodesic Walk on $M_L$

Recall that the geodesic walk is given by

$$x^{(j+1)} = \exp_{x^{(j)}}(\sqrt{h}w + \frac{h}{2}\mu(x^{(j)}))$$

where  $w \sim N(0, I)$  (in the local coordinates at  $x^{(j)}$ ). In many proofs in this section, we consider the geodesic  $\gamma$  from  $x^{(j)}$  to  $x^{(j+1)}$  with the initial velocity

$$\gamma'(0) = \frac{w}{\sqrt{n}} + \frac{1}{2}\sqrt{\frac{h}{n}}\mu(x).$$

The scaling is to make the speed of geodesic  $\gamma$  close to one. Since  $w$  is a Gaussian vector, we have that  $0.9 \leq \|\gamma'(0)\|_{\gamma(0)} \leq 1.1$  with high probability. Due to this rescaling, the geodesic is defined from 0 to  $\ell = \sqrt{nh}$ .

We often work on the Euclidean coordinates. In this case, the geodesic walk is given by the following formula:

**Lemma 14.5.7.** *Given  $x \in M_L$  and step size  $h > 0$ , one step of the geodesic walk starting at  $x$  is given by the solution  $\gamma(\ell)$  of the following geodesic equation*

$$\begin{aligned}\gamma''(t) &= (A_\gamma^T A_\gamma)^{-1} A_\gamma^T s_{\gamma'}^2 \text{ for } 0 \leq t \leq \ell \\ \gamma'(0) &= \frac{w}{\sqrt{n}} + \frac{1}{2} \sqrt{\frac{h}{n}} \mu(x) \\ \gamma(0) &= x\end{aligned}$$

where  $\ell = \sqrt{nh}$ ,  $w \sim N(0, (A_x^T A_x)^{-1})$  and  $\mu(x) = (A_x^T A_x)^{-1} A_x^T \sigma_x$ .

*Proof.* The geodesic equation is given by Lemma 14.5.3. From (14.2), the drift term is given by

$$\begin{aligned}\mu_i(x) &= \frac{1}{2} \sum_{j=1}^n \frac{\partial}{\partial x_j} \left( (\nabla^2 \phi(X_t))^{-1} \right)_{ij} = \frac{1}{2} \sum_{j=1}^n \frac{\partial}{\partial x_j} \left( (A_x^T A_x)^{-1} \right)_{ij} \\ &= \sum_j e_i^T (A_x^T A_x)^{-1} A_x^T S_{x, e_j} A_x (A_x^T A_x)^{-1} e_j \\ &= \sum_j \sum_k V_{ik} e_k^T A_x e_j V_{jk} \quad \text{where } V = (A_x^T A_x)^{-1} A_x^T \\ &= \sum_k V_{ik} e_k^T A_x \sum_j e_j (e_k^T A_x (A_x^T A_x)^{-1} e_j) \\ &= \sum_k V_{ik} e_k^T A_x \sum_j e_j e_j^T (A_x^T A_x)^{-1} A_x^T e_k \\ &= \sum_k V_{ik} e_k^T A_x (A_x^T A_x)^{-1} A_x^T e_k \\ &= V_i \sigma_x = e_i^T (A_x^T A_x)^{-1} A_x^T \sigma_x.\end{aligned}$$

□

### ■ 14.5.3 Randomness and Stability of the Geodesic

In this and subsequent sections, we will frequently use the following elementary calculus facts (using only the chain/product rules and formula for derivative of inverse of a matrix):

$$\begin{aligned}\frac{dA_x}{dt} &= -S_{\gamma'} A_x \\ \frac{dP_x}{dt} &= -S_{\gamma'} P_x - P_x S_{\gamma'} + 2P_x S_{\gamma'} P_x \\ \frac{dS_{\gamma'}}{dt} &= \text{Diag}(-S_{\gamma'} A_x \gamma' + A_x \gamma'') = -S_{\gamma'}^2 + \text{Diag}(P_x) S_{\gamma'}^2 \\ \frac{d\sigma_x}{dt} &= \text{Diag}\left(\frac{dP_x}{dt}\right)\end{aligned}$$

We also use these matrix inequalities:  $\text{Tr}(AB) = \text{Tr}(BA)$ ,  $\text{Tr}(PAP) \leq \text{Tr}(A)$  for any psd matrix  $A$ ;  $\text{Tr}(ABA^T) \leq \text{Tr}(AZA^T)$  for any  $B \preceq Z$ ; the Cauchy-Schwartz, namely,  $\text{Tr}(AB) \leq \text{Tr}(AA^T)^{\frac{1}{2}} \text{Tr}(B^T B)^{\frac{1}{2}}$ . We also use  $P^2 = P$  since  $P$  is a projection matrix.

We note that  $\|\gamma'(0)\|_4$  is small because it is the sum of a random vector plus a small drift term. Here, we would like to prove that the geodesic is stable in  $\ell^4$  norm. This fact is very important in the proof and hence we first establish this here.

**Lemma 14.5.8.** *Let  $\gamma$  be a geodesic in  $M_L$  starting at  $x$ . Let  $v_4 = \|s_{\gamma'(0)}\|_4$ . Then, for  $0 \leq t \leq \frac{1}{10v_4}$ , we have that*

1.  $\|s_{\gamma'(t)}\|_4 \leq 1.25v_4$ .
2.  $\|\gamma''\|_\gamma^2 \leq 3v_4^4$ .

*Proof.* Let  $u(t) = \|s_{\gamma'(t)}\|_4$ . Then, we have

$$\begin{aligned} \frac{du}{dt} &\leq \left\| \frac{d}{dt} (A_\gamma \gamma') \right\|_4 \\ &= \|A_\gamma \gamma'' - (A_\gamma \gamma')^2\|_4 \\ &\leq \|A_\gamma \gamma''\|_4 + u^2(t). \end{aligned} \tag{14.21}$$

Under the Euclidean coordinates, the geodesic equation is given by

$$\gamma'' = (A_\gamma^T A_\gamma)^{-1} A_\gamma^T s_{\gamma'}^2.$$

Hence, we have that

$$\begin{aligned} \|\gamma''\|_\gamma^2 &= (s_{\gamma'}^2)^T A_\gamma (A_\gamma^T A_\gamma)^{-1} (A_\gamma^T A_\gamma) (A_\gamma^T A_\gamma)^{-1} A_\gamma^T s_{\gamma'}^2 \\ &\leq \sum_i (s_{\gamma'}^4)_i = u^4(t). \end{aligned} \tag{14.22}$$

Therefore, we have

$$\|A_\gamma \gamma''\|_4 \leq \|A_\gamma \gamma''\|_2 \leq u^2(t).$$

Plugging it into (14.21), we have that

$$\frac{du}{dt} \leq u^2(t) + u^2(t) \leq 2u^2(t).$$

Since  $u(0) = v_4$ , for all  $t \geq 0$ , we have

$$u(t) \leq \frac{v_4}{1 - 2v_4 t}. \tag{14.23}$$

For  $0 \leq t \leq \frac{1}{10v_4}$ , we have  $u(t) \leq 1.25v_4$  and this gives the first inequality. Using (14.22), we get the second inequality.  $\square$

**Lemma 14.5.9.** *For  $p \geq 1$ , we have*

$$P_{x \sim N(0, I)} \left( \|Ax\|_p^p \leq \left( \left( \frac{2^{p/2} \Gamma(\frac{p+1}{2})}{\sqrt{\pi}} \sum_i \|a_i\|_2^p \right)^{1/p} + \|A\|_{2 \rightarrow p} \right)^p \right) \leq 1 - \exp\left(-\frac{t^2}{2}\right).$$

*In particular, we have*

$$P_{x \sim N(0, I)} \left( \|Ax\|_2^2 \leq (\|A\|_F + \|A\|_{2t})^2 \right) \leq 1 - \exp\left(-\frac{t^2}{2}\right)$$

and

$$P_{x \sim N(0, I)} \left( \|Ax\|_4^4 \leq \left( \left( 3 \sum_i \|a_i\|_2^4 \right)^{1/4} + \|A\|_{2 \rightarrow 4} t \right)^4 \right) \leq 1 - \exp\left(-\frac{t^2}{2}\right).$$

*Proof.* Let  $F(x) = \|Ax\|_p$ . Since  $|F(x) - F(y)| \leq \|A\|_{2 \rightarrow p} \|x - y\|_2$ , Gaussian concentration shows that

$$P(F(x) \leq \mathbf{E}F(x) + \|A\|_{2 \rightarrow p} t) \leq 1 - \exp\left(-\frac{t^2}{2}\right).$$

Since  $x^p$  is convex, we have that

$$\begin{aligned} \mathbf{E}\|Ax\|_p &\leq (\mathbf{E}\|Ax\|_p^p)^{1/p} = \left( \sum_i \mathbf{E}|a_i^T x|^p \right)^{1/p} \\ &= \left( \mathbf{E}_{t \sim N(0,1)} |t|^p \sum_i \|a_i\|_2^p \right)^{1/p} \\ &= \left( \frac{2^{p/2} \Gamma(\frac{p+1}{2})}{\sqrt{\pi}} \sum_i \|a_i\|_2^p \right)^{1/p}. \end{aligned}$$

Hence, we have that

$$P \left( \|Ax\|_p^p \leq \left( \left( \frac{2^{p/2} \Gamma(\frac{p+1}{2})}{\sqrt{\pi}} \sum_i \|a_i\|_2^p \right)^{1/p} + \|A\|_{2 \rightarrow p} t \right)^p \right) \leq 1 - \exp\left(-\frac{t^2}{2}\right).$$

□

Using the lemma above, we prove that within small time, the direction of geodesic is a random vector plus a small vector.

**Lemma 14.5.10.** *Let  $\gamma$  given by a geodesic walk on  $M_L$  with step size  $h$  satisfying  $h \leq \frac{1}{256\sqrt{n}}$ . Let  $u$  be the Gaussian part of the initial velocity  $\gamma'(0)$ . If  $\|A_{\gamma(0)}u\|_2 \leq \frac{3}{2}$  and  $\|A_\gamma u\|_4 \leq \frac{1.55}{n^{1/4}}$ , then for any  $0 \leq t \leq \ell$ ,*

1.  $\|A_\gamma \gamma'\|_2 \leq 2$ .
2.  $\|A_\gamma \gamma'\|_4 \leq 2n^{-1/4}$ .
3.  $\|\gamma''\|_\gamma^2 \leq 15n^{-1}$ .
4.  $\|A_\gamma \gamma'\|_\infty \leq 6\sqrt{h}$ .

Furthermore,  $\|A_{\gamma(0)}u\|_2 \leq \frac{3}{2}$  and  $\|A_\gamma u\|_4 \leq \frac{1.55}{n^{1/4}}$  holds with probability at least  $1 - \exp(-\sqrt{n}/100)$ .

*Proof.* From the definition of the geodesic walk (Lemma 14.5.7), we have that

$$\gamma'(0) = u + v$$

where  $u \sim N\left(0, \frac{1}{n} (A_\gamma^T A_\gamma)^{-1}\right)$  and  $v = \frac{1}{2} \sqrt{\frac{h}{n}} (A_\gamma^T A_\gamma)^{-1} A_\gamma^T \sigma_\gamma$ . Note that

$$\begin{aligned} \|A_\gamma \gamma'(0)\|_2 &\leq \|A_\gamma u\|_2 + \left\| \frac{1}{2} \sqrt{\frac{h}{n}} A_\gamma (A_\gamma^T A_\gamma)^{-1} A_\gamma^T \sigma_\gamma \right\|_2 \\ &= \|A_\gamma u\|_2 + \frac{1}{2} \sqrt{\frac{h}{n}} \sqrt{\sigma_\gamma^T A_\gamma (A_\gamma^T A_\gamma)^{-1} A_\gamma^T \sigma_\gamma}. \end{aligned}$$

Note that  $A_\gamma u = Bx$  where  $B = \frac{1}{\sqrt{n}} A_\gamma (A_\gamma^T A_\gamma)^{-1/2}$  and  $x \sim N(0, I)$ . Note that  $\|Bx\|_2^2 = \frac{1}{n} \|x\|_2^2$  for all  $x$  and hence  $\|B\|_F^2 = 1$  and  $\|B\|_{2 \rightarrow 2} = \frac{1}{\sqrt{n}}$ . Hence, Lemma 14.5.9 shows that

$$P\left(\|A_\gamma u\|_2^2 \leq \left(1 + \frac{t}{\sqrt{n}}\right)^2\right) \leq 1 - \exp\left(-\frac{t^2}{2}\right)$$

Hence, we have

$$\|A_\gamma \gamma'\|_2 \leq \frac{3}{2} + \frac{1}{2} \sqrt{h} \leq 2$$

with probability  $1 - \exp(-n/8)$ . Since geodesic preserves the speed, this gives the first inequality.

Next, we note that  $\sum_i \|e_i^T B\|_2^4 = \frac{1}{n^2} \sum_i (\sigma_\gamma)_i^2 \leq \frac{1}{n}$  and  $\|B\|_{2 \rightarrow 4} \leq \|B\|_{2 \rightarrow 2} = \frac{1}{\sqrt{n}}$ . Hence, Lemma 14.5.9 shows that

$$P\left(\|A_\gamma u\|_4^4 \leq \left(\left(\frac{3}{n}\right)^{1/4} + \frac{t}{\sqrt{n}}\right)^4\right) \leq 1 - \exp\left(-\frac{t^2}{2}\right).$$

Hence, we have that

$$v_4 \stackrel{\text{def}}{=} \|A_\gamma \gamma'(0)\|_4 \leq \frac{1.55}{n^{1/4}} + \frac{1}{2} \sqrt{h} \leq \frac{1.6}{n^{1/4}}$$

with probability at least  $1 - \exp(-\sqrt{n}/100)$ . Since  $\ell = \sqrt{nh} \leq \frac{n^{1/4}}{16} = \frac{1}{10v_4}$ , Lemma 14.5.8 shows the middle two inequalities.

Lemma 14.5.8 shows that  $\|A_\gamma \gamma''\|_\infty \leq \|\gamma''\|_{\gamma(t)} \leq \sqrt{3} v_4^2 \leq 5n^{-1/2}$  and hence

$$\begin{aligned} \|A_{\gamma(t)} \gamma'(t)\|_\infty &\leq \|A_\gamma \gamma'(0)\|_\infty + \int_0^\ell \|A_{\gamma(t)} \gamma''(t)\|_\infty dt \\ &\leq \sqrt{h} + 5\ell n^{-1/2} \\ &= 6\sqrt{h}. \end{aligned}$$

□

### ■ 14.5.4 Stability of Drift

We begin with  $D_0$  and  $D_1$ .

**Lemma 14.5.11.** *For any  $x \in M_L$ , we have that  $\|\mu(x)\|_x^2 \leq n$ . Hence,  $D_0 \leq \sqrt{n}$ .*

*Proof.* We have

$$\|\mu(x)\|_x^2 = \sigma_x^T A_x (A_x^T A_x)^{-1} A_x \sigma_x \leq \sum_i (\sigma_x)_i^2 \leq n.$$

□

**Lemma 14.5.12.** *Let  $\gamma(t)$  be a geodesic walk on  $M_L$  with step size  $h$  satisfying  $\frac{1}{n} \leq h \leq \frac{1}{256\sqrt{n}}$ . Then, we have that*

$$D_1 = \sup_{0 \leq t \leq \ell} \left| \frac{d}{dt} \|\mu(\gamma(t))\|^2 \right| \leq 32n\sqrt{h}.$$

with probability at least  $1 - \exp(-\sqrt{n}/100)$  over the choice of the next step  $\gamma \in \tilde{\Xi}_x$  for any  $x \in M_L$ .

*Proof.* We want to bound  $|\|\mu(x)\|_x^2 - \|\mu(y)\|_y^2|$ . Note that

$$\begin{aligned} \|\mu(\gamma(t))\|_\gamma^2 &= \sigma_\gamma^T A_\gamma (A_\gamma^T A_\gamma)^{-1} A_\gamma \sigma_\gamma \\ &= 1^T \text{Diag}(P_\gamma) P_\gamma \text{Diag}(P_\gamma) 1. \end{aligned}$$

Using  $\frac{d}{dt} P_\gamma = -S_{\gamma'} P_\gamma - P_\gamma S_{\gamma'} + 2P_\gamma S_{\gamma'} P_\gamma$ , we have

$$\begin{aligned} \frac{d}{dt} \|\mu(\gamma(t))\|_\gamma^2 &= -2 \cdot 1^T \text{Diag}(S_{\gamma'} P_\gamma) P_\gamma \text{Diag}(P_\gamma) 1 \\ &\quad -2 \cdot 1^T \text{Diag}(P_\gamma S_{\gamma'}) P_\gamma \text{Diag}(P_\gamma) 1 \\ &\quad +4 \cdot 1^T \text{Diag}(P_\gamma S_{\gamma'} P_\gamma) P_\gamma \text{Diag}(P_\gamma) 1 \\ &\quad -1^T \text{Diag}(P_\gamma) S_{\gamma'} P_\gamma \text{Diag}(P_\gamma) 1 \\ &\quad -1^T \text{Diag}(P_\gamma) P_\gamma S_{\gamma'} \text{Diag}(P_\gamma) 1 \\ &\quad +2 \cdot 1^T \text{Diag}(P_\gamma) P_\gamma S_{\gamma'} P_\gamma \text{Diag}(P_\gamma) 1 \\ &= -2 \cdot 1^T \text{Diag}(S_{\gamma'} P_\gamma) P_\gamma \text{Diag}(P_\gamma) 1 \\ &\quad -4 \cdot 1^T \text{Diag}(P_\gamma) S_{\gamma'} P_\gamma \text{Diag}(P_\gamma) 1 \\ &\quad +4 \cdot 1^T \text{Diag}(P_\gamma S_{\gamma'} P_\gamma) P_\gamma \text{Diag}(P_\gamma) 1 \\ &\quad +2 \cdot 1^T \text{Diag}(P_\gamma) P_\gamma S_{\gamma'} P_\gamma \text{Diag}(P_\gamma) 1. \end{aligned}$$

Now, we note that

$$\begin{aligned} |1^T \text{Diag}(S_{\gamma'} P_\gamma) P_\gamma \text{Diag}(P_\gamma) 1| &\leq \|\gamma'\|_\gamma \sqrt{n}, \\ |1^T \text{Diag}(P_\gamma) S_{\gamma'} P_\gamma \text{Diag}(P_\gamma) 1| &\leq \|\gamma'\|_\gamma \sqrt{n}, \\ |1^T \text{Diag}(P_\gamma S_{\gamma'} P_\gamma) P_\gamma \text{Diag}(P_\gamma) 1| &\leq \|\gamma'\|_\gamma \sqrt{n}, \\ |1^T \text{Diag}(P_\gamma) P_\gamma S_{\gamma'} P_\gamma \text{Diag}(P_\gamma) 1| &\leq \|S_{\gamma'}\|_\infty n. \end{aligned}$$

Since  $\|\gamma'\|_\gamma \leq 2$  and  $\|A_\gamma \gamma'\|_\infty \leq 6\sqrt{h}$  (Lemma 14.5.10), we have

$$D_1 = \sup \left| \frac{d}{dt} \|\mu(\gamma(t))\|^2 \right| \leq 20\sqrt{n} + 12n\sqrt{h} \leq 32n\sqrt{h}.$$

□

### ■ 14.5.5 Smoothness of the metric

**Lemma 14.5.13.** *Let  $\gamma$  given by a geodesic walk on  $M_L$  with step size  $h$  satisfying  $\frac{1}{n} \leq h \leq \frac{1}{256\sqrt{n}}$ . Let  $f(t) = \log \det(A_{\gamma(t)}^T A_{\gamma(t)})$ . Then, we have*

$$G_1 = \sup_{0 \leq t \leq \ell} |f'''(t)| \leq 1000\sqrt{h}$$

with probability at least  $1 - \exp(-\sqrt{n}/100)$ .

*Proof.* Note that

$$f'(t) = -2\text{Tr}(A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'} A_\gamma.$$

Hence, we have

$$\begin{aligned} f''(t) &= -4\text{Tr}(A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'} A_\gamma (A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'} A_\gamma \\ &\quad + 6\text{Tr}(A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'}^2 A_\gamma \\ &\quad - 2\text{Tr}(A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma''} A_\gamma. \end{aligned}$$

So, we have

$$\begin{aligned} f'''(t) &= -16\text{Tr}(A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'} A_\gamma (A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'} A_\gamma (A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'} A_\gamma \\ &\quad + 24\text{Tr}(A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'}^2 A_\gamma (A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'} A_\gamma \\ &\quad - 8\text{Tr}(A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma''} A_\gamma (A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'} A_\gamma \\ &\quad + 12\text{Tr}(A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'} A_\gamma (A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'}^2 A_\gamma \\ &\quad - 24\text{Tr}(A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'}^3 A_\gamma \\ &\quad + 12\text{Tr}(A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'} S_{\gamma''} A_\gamma \\ &\quad - 4\text{Tr}(A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'} A_\gamma (A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma''} A_\gamma \\ &\quad + 4\text{Tr}(A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma''} S_{\gamma'} A_\gamma \\ &\quad - 2\text{Tr}(A_\gamma^T A_\gamma)^{-1} A_\gamma^T \frac{d}{dt} S_{\gamma''} A_\gamma \\ &= -16\text{Tr} P_\gamma S_{\gamma'} P_\gamma S_{\gamma'} P_\gamma S_{\gamma'} + 36\text{Tr} P_\gamma S_{\gamma'}^2 P_\gamma S_{\gamma'} \\ &\quad - 12\text{Tr} P_\gamma S_{\gamma''} P_\gamma S_{\gamma'} - 24\text{Tr} P_\gamma S_{\gamma'}^3 \\ &\quad + 16\text{Tr} P_\gamma S_{\gamma'} S_{\gamma''} - 2\text{Tr} P_\gamma \frac{d}{dt} S_{\gamma''}. \end{aligned}$$

Hence, we have

$$|f'''(t)| \leq (16 + 36 + 24)\sqrt{\text{Tr} S_{\gamma'}^4 \text{Tr} S_{\gamma'}^2} + (12 + 16)\sqrt{\text{Tr} S_{\gamma''}^2 \text{Tr} S_{\gamma'}^2} + 2 \left| \text{Tr} P_\gamma \frac{d}{dt} S_{\gamma''} \right|.$$

Since  $\text{Tr} S_{\gamma'}^2 \leq 2$ ,  $\text{Tr} S_{\gamma'}^4 \leq 16n^{-1}$ ,  $\text{Tr} S_{\gamma''}^2 \leq 15n^{-1}$  (Lemma 14.5.10), we have

$$\begin{aligned} |f'''(t)| &\leq 76\sqrt{2 \cdot 16n^{-1}} + 28\sqrt{2 \cdot 15n^{-1}} + 2 \left| \text{Tr} P_\gamma \frac{d}{dt} S_{\gamma''} \right| \\ &\leq 600n^{-1/2} + 2 \left| \text{Tr} P_\gamma \frac{d}{dt} S_{\gamma''} \right|. \end{aligned}$$

To bound the last term, we start with the geodesic equation:  $s_{\gamma''} = P_\gamma s_{\gamma'}^2$ . Since  $\frac{d}{dt} P_\gamma = -S_{\gamma'} P_\gamma - P_\gamma S_{\gamma'} + 2P_\gamma S_{\gamma'} P_\gamma$ , we have that

$$\begin{aligned} \frac{d}{dt} s_{\gamma''} &= -S_{\gamma'} P_\gamma s_{\gamma'}^2 - P_\gamma S_{\gamma'} s_{\gamma'}^2 + 2P_\gamma S_{\gamma'} P_\gamma s_{\gamma'}^2 \\ &\quad + 2P_\gamma s_{\gamma'} s_{\gamma''} - 2P_\gamma s_{\gamma'}^3 \\ &= -S_{\gamma'} P_\gamma s_{\gamma'}^2 - 3P_\gamma s_{\gamma'}^3 + 2P_\gamma S_{\gamma'} P_\gamma s_{\gamma'}^2 + 2P_\gamma s_{\gamma'} s_{\gamma''} \\ &= -S_{\gamma'} P_\gamma s_{\gamma'}^2 - 3P_\gamma s_{\gamma'}^3 + 4P_\gamma S_{\gamma'} P_\gamma s_{\gamma'}^2 \end{aligned}$$



Hence, we have that

$$\begin{aligned}
\left| \text{Tr} P_\gamma \frac{d}{dt} S_{\gamma''} \right| &\leq |\sigma_\gamma^T S_{\gamma'} P_\gamma s_{\gamma'}^2| + 3 |\sigma_\gamma^T P_\gamma s_{\gamma'}^3| + 4 |\sigma_\gamma^T P_\gamma S_{\gamma'} P_\gamma s_{\gamma'}^2| \\
&\leq |(S_{\gamma'} \sigma_\gamma)^T P_\gamma s_{\gamma'}^2| + 3 |(S_{\gamma'} P_\gamma \sigma_\gamma)^T s_{\gamma'}^2| + 4 (S_{\gamma'} P_\gamma \sigma_\gamma)^T (P_\gamma s_{\gamma'}^2) \\
&\leq \sqrt{\sum s_{\gamma'}^2 \sum s_{\gamma'}^4} + 28n^{-1/2} \sqrt{\sigma_\gamma^T P_\gamma S_{\gamma'}^2 P_\gamma \sigma_\gamma} \\
&\leq 8n^{-1/2} + 168\sqrt{h} \leq 180\sqrt{h}.
\end{aligned}$$

Hence, we have that  $|f'''(t)| \leq 600n^{-1/2} + 360\sqrt{h} \leq 1000\sqrt{h}$ .  $\square$

### ■ 14.5.6 Stability of Jacobi Field

Now, we give a bound on  $R_1$ .

**Lemma 14.5.14.** *Let  $\gamma$  given by a geodesic walk on  $M_L$  with step size  $h$  satisfying  $\frac{1}{n} \leq h \leq \frac{1}{256\sqrt{n}}$ . With probability at least  $1 - \exp(-\sqrt{n}/100)$  over the choice of the next step  $\gamma \in \tilde{\Xi}_x$ , we have that*

$$R_1 = \sup_{0 \leq t \leq \ell} \|R(t)\|_F \leq \frac{8}{\sqrt{n}}$$

where  $R(t)_{ij} = \langle R(X_i(t), \gamma'(t)) \gamma'(t), X_j(t) \rangle$  and  $\{X_i\}_{i=1}^n$  is any orthogonal basis at  $\gamma(t)$ .

*Proof.* Let  $\overline{R(t)}$  such that  $\langle R(u, \gamma'(t)) \gamma'(t), v \rangle = u^T \overline{R(t)} v$ . Lemma 14.5.4 shows that

$$u^T \overline{R(t)} v = s_u^T S_{\gamma'} A_\gamma (A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'} s_v - (s_v s_u)^T A_\gamma (A_\gamma^T A_\gamma)^{-1} A_\gamma^T s_{\gamma'}^2.$$

Hence, we have

$$\overline{R(t)} = A_\gamma^T S_{\gamma'} P_\gamma S_{\gamma'} A_\gamma - A_\gamma^T \text{Diag}(P_\gamma s_{\gamma'}^2) A_\gamma.$$

Pick  $X_i = (A_\gamma^T A_\gamma)^{-1/2} e_i$ . Then, we have

$$R(t) = (A_\gamma^T A_\gamma)^{-1/2} (A_\gamma^T S_{\gamma'} P_\gamma S_{\gamma'} A_\gamma - A_\gamma^T \text{Diag}(P_\gamma s_{\gamma'}^2) A_\gamma) (A_\gamma^T A_\gamma)^{-1/2}.$$

The claim follows from Lemma 14.5.10 and the calculation

$$\begin{aligned}
&\|R(t)\|_F^2 \\
&\leq 2\|(A_\gamma^T A_\gamma)^{-1/2} A_\gamma^T S_{\gamma'} P_\gamma S_{\gamma'} A_\gamma (A_\gamma^T A_\gamma)^{-1/2}\|_F^2 + 2\|(A_\gamma^T A_\gamma)^{-1/2} A_\gamma^T \text{Diag}(P_\gamma s_{\gamma'}^2) A_\gamma (A_\gamma^T A_\gamma)^{-1/2}\|_F^2 \\
&= 2\text{Tr} P_\gamma S_{\gamma'} P_\gamma S_{\gamma'} P_\gamma S_{\gamma'} P_\gamma S_{\gamma'} + 2\text{Tr} P_\gamma \text{Diag}(P_\gamma s_{\gamma'}^2) P_\gamma \text{Diag}(P_\gamma s_{\gamma'}^2) \leq 4\|s_{\gamma'}\|_4^4.
\end{aligned}$$

$\square$

### ■ 14.5.7 Smoothness of One-Step Distribution

We begin with  $R_2$ .

**Lemma 14.5.15.** *Given a geodesic  $\gamma(t)$  on  $M_L$  with step size  $h$  satisfying  $\frac{1}{n} \leq h \leq \frac{1}{256\sqrt{n}}$ , with probability at least  $1 - \exp(-\sqrt{n}/100)$  over the choice of the next step  $\gamma \in \tilde{\Xi}_x$ , for any  $t$  such that  $0 \leq t \leq \ell$ , any curve  $c(s)$  starting from  $\gamma(t)$  and any vector field  $v(s)$  on  $c(s)$ , we have that*

$$\left| \frac{d}{ds} \text{Ric}(v(s)) \Big|_{s=0} \right| \leq 64\sqrt{nh} \|c'(0)\| + 6\sqrt{nh} \|D_s v(0)\|.$$

Therefore,  $R_2 = O(\sqrt{nh})$ .

*Proof.* By Lemma 14.5.4, we know that

$$\begin{aligned} Ric(v(s)) &= s_{c(s),v(s)}^T P_{c(s)}^{(2)} s_{c(s),v(s)} - \sigma_{c(s)}^T P_{c(s)} s_{c(s),v(s)}^2 \\ &= \text{Tr}(S_{c(s),v(s)} P_{c(s)} S_{c(s),v(s)} P_{c(s)}) - \text{Tr}(\text{Diag}(P_{c(s)} s_{c(s),v(s)}^2) P_{c(s)}). \end{aligned}$$

For simplicity, we suppress the parameter  $s$  and hence, we have

$$Ric(v) = \text{Tr}(S_{c,v} P_c S_{c,v} P_c) - \text{Tr}(\text{Diag}(P_c s_{c,v}^2) P_c)$$

Now, we use that

$$\begin{aligned} \frac{d}{ds} P_c &= -S_{c'} P_c - P_c S_{c'} + 2P_c S_{c'} P_c, \\ \frac{d}{ds} S_{c,v} &= -S_{c'} S_{c,v} + S_{c,v}' \end{aligned}$$

and get

$$\begin{aligned} &\frac{d}{ds} Ric(v) \\ &= -2\text{Tr}(S_{c,v} S_{c'} P_c S_{c,v} P_c) - 2\text{Tr}(S_{c,v} P_c S_{c'} S_{c,v} P_c) + 2\text{Tr}(S_{c,v} P_c S_{c'} P_c S_{c,v} P_c) \\ &\quad - 2\text{Tr}(S_{c'} S_{c,v} P_c S_{c,v} P_c) + 2\text{Tr}(S_{c,v}' P_c S_{c,v} P_c) \\ &\quad + \text{Tr}(\text{Diag}(P_c s_{c,v}^2) S_{c'} P_c) + \text{Tr}(\text{Diag}(P_c s_{c,v}^2) P_c S_{c'}) - 2\text{Tr}(\text{Diag}(P_c s_{c,v}^2) P_c S_{c'} P_c) \\ &\quad + \text{Tr}(\text{Diag}(P_c S_{c'} s_{c,v}^2) P_c) + \text{Tr}(\text{Diag}(S_{c'} P_c s_{c,v}^2) P_c) - 2\text{Tr}(\text{Diag}(P_c S_{c'} P_c s_{c,v}^2) P_c) \\ &\quad + 2\text{Tr}(\text{Diag}(P_c S_{c,v} S_{c'} s_{c,v}) P_c) - 2\text{Tr}(\text{Diag}(P_c S_{c,v} s_{c,v}') P_c) \\ &= -6\text{Tr}(S_{c,v} S_{c'} P_c S_{c,v} P_c) + 2\text{Tr}(S_{c,v} P_c S_{c'} P_c S_{c,v} P_c) + 2\text{Tr}(S_{c,v}' P_c S_{c,v} P_c) \\ &\quad + 2\text{Tr}(\text{Diag}(P_c s_{c,v}^2) S_{c'} P_c) - 2\text{Tr}(\text{Diag}(P_c s_{c,v}^2) P_c S_{c'} P_c) \\ &\quad + 3\text{Tr}(\text{Diag}(P_c S_{c'} s_{c,v}^2) P_c) + \text{Tr}(\text{Diag}(S_{c'} P_c s_{c,v}^2) P_c) - 2\text{Tr}(\text{Diag}(P_c S_{c'} P_c s_{c,v}^2) P_c) \\ &\quad - 2\text{Tr}(\text{Diag}(P_c S_{c,v} s_{c,v}') P_c). \end{aligned}$$

Let  $\frac{d}{ds} Ric(v) = (1) + (2)$  where (1) is the sum of all terms not involving  $v'$  and (2) is the sum of other terms.

For the first term (1), let  $L = \|s_{c'}(s)\|_2$ . Then, we have that

$$\begin{aligned} |(1)| &\leq 6 |\text{Tr}(S_{c,v} S_{c'} P_c S_{c,v} P_c)| + 2 |\text{Tr}(S_{c,v} P_c S_{c'} P_c S_{c,v} P_c)| \\ &\quad + 2 |\text{Tr}(\text{Diag}(P_c s_{c,v}^2) S_{c'} P_c)| + 2 |\text{Tr}(\text{Diag}(P_c s_{c,v}^2) P_c S_{c'} P_c)| \\ &\quad + 3 |\text{Tr}(\text{Diag}(P_c S_{c'} s_{c,v}^2) P_c)| + |\text{Tr}(\text{Diag}(S_{c'} P_c s_{c,v}^2) P_c)| + 2 |\text{Tr}(\text{Diag}(P_c S_{c'} P_c s_{c,v}^2) P_c)| \\ &\leq 6L \sqrt{\sum_i (s_{c,v})_i^2} \sqrt{\sum_i (s_{c,v})_i^2} + 2L |\text{Tr}(S_{c,v} P_c S_{c,v} P_c)| \\ &\quad + 2L \sqrt{\sum_i (s_{c,v})_i^4} \sqrt{\sum_i (P_c)_{ii}^2} + 2L \sqrt{\sum_i (s_{c,v})_i^4} \sqrt{\sum_i (P_c S_{c'} P_c)_{ii}^2} \\ &\quad + 3L \sqrt{\sum_i (s_{c,v})_i^4} \sqrt{\sum_i (P_c)_{ii}^2} + L \sqrt{\sum_i (s_{c,v})_i^4} \sqrt{\sum_i (P_c)_{ii}^2} + 2L \sqrt{\sum_i (s_{c,v})_i^4} \sqrt{\sum_i (P_c)_{ii}^2} \\ &\leq 8L \sqrt{\sum_i (s_{c,v})_i^2} \sqrt{\sum_i (s_{c,v})_i^2} + 10\sqrt{n}L \sqrt{\sum_i (s_{c,v})_i^4}. \end{aligned}$$

Since  $v(0) = \gamma'(t)$ , Lemma 14.5.10 shows that  $\|s_{c,v}\|_2 \leq 2$  and  $\|s_{c,v}\|_4 \leq 2n^{-1/4}$ . Hence, we have

$$|(1)| \leq 52L.$$

For the second term (2), we note that

$$D_s v = \frac{dv}{ds} - (A_c^T A_c)^{-1} A_c^T S_{c'} s_{c,v}.$$

Therefore, we have

$$s_{c,v'} = A_c (D_s v) - A_c (A_c^T A_c)^{-1} A_c^T S_{c'} s_{c,v}$$

and hence

$$\begin{aligned} \|s_{c,v'}\|_2 &\leq \|D_s v\| + \|A_c (A_c^T A_c)^{-1} A_c^T S_{c'} s_{c,v}\|_2 \\ &\leq \|D_s v\| + L \|s_{c,v}\|_2 \\ &\leq 2L + \|D_s v\|. \end{aligned}$$

Also, note that Lemma 14.5.10 shows that  $\|s_{c,v}\|_\infty \leq 6\sqrt{h}$ . Hence, we have that

$$\begin{aligned} |(2)| &\leq 2 |\text{Tr}(S_{c,v'} P_c S_{c,v} P_c)| + 2 |\text{Tr}(\text{Diag}(P_c S_{c,v} s_{c,v'}) P_c)| \\ &\leq 2 \|s_{c,v'}\|_2 \|s_{c,v}\|_2 + 2\sqrt{n} \sqrt{\sum_i (s_{c,v'} s_{c,v})_i^2} \\ &\leq 4(2 + \|D_s v\|) + 2\sqrt{nh} (2 + \|D_s v\|) \\ &\leq 8 + 4\|D_s v\| + 4\sqrt{nh} + 2\sqrt{nh} \|D_s v\| \\ &\leq 12\sqrt{nh} + 6\sqrt{nh} \|D_s v\|. \end{aligned}$$

Therefore, we have

$$\begin{aligned} \left| \frac{d}{ds} \text{Ric}(v) \right| &\leq 52L + 12\sqrt{nh}L + 6\sqrt{nh} \|D_s v\| \\ &\leq 64\sqrt{nh}L + 6\sqrt{nh} \|D_s v\|. \end{aligned}$$

□

Finally, we bound  $D_2$ .

**Lemma 14.5.16.** *Given  $x \in M_L$  and any curve  $c(t)$  starting at  $x$  with unit speed. We have that*

$$\left\| \frac{D}{dt} \mu(c(t)) \right\| \leq 9\sqrt{n}.$$

Hence,  $D_2 \leq 9\sqrt{n}$ .

*Proof.* Recall that  $\frac{D}{dt} \mu(c(t)) = \frac{d\mu(c)}{dt} - (A_c^T A_c)^{-1} A_c^T S_{c,\mu} s_{c'}$ . We bound the terms separately. For the second term, since  $c$  is an unit speed curve, we have that  $\|s_{c'}\|_2 = 1$  and

$$\begin{aligned} \|(A_c^T A_c)^{-1} A_c^T S_{c'} A_c \mu(c)\|_c &\leq \|S_{c,\mu} s_{c'}\|_2 \leq \|S_{c,\mu}\|_\infty \\ &= \|A_c (A_c^T A_c)^{-1} A_c^T \sigma_c\|_\infty \leq \sqrt{n}. \end{aligned}$$

For the first term,

$$\begin{aligned}
\frac{d}{dt}(A_c\mu(c)) &= \frac{d}{dt}(P_c \mathbf{diag}(P_c)) \\
&= -S_{c'} P_c \mathbf{diag}(P_c) - P_c S_{c'} \mathbf{diag}(P_c) + 2P_c S_{c'} P_c \mathbf{diag}(P_c) \\
&\quad - P_c \mathbf{diag}(S_{c'} P_c) - P_c \mathbf{diag}(P_c S_{c'}) + 2P_c \mathbf{diag}(P_c S_{c'} P_c) \\
&= -S_{c'} P_c \mathbf{diag}(P_c) - 2P_c S_{c'} \mathbf{diag}(P_c) + 2P_c S_{c'} P_c \mathbf{diag}(P_c) \\
&\quad - P_c \mathbf{diag}(P_c S_{c'}) + 2P_c \mathbf{diag}(P_c S_{c'} P_c).
\end{aligned}$$

Using  $\|s_{c'}\|_2 = 1$ , we have

$$\begin{aligned}
\left\| \frac{d}{dt} \mu(c) \right\|_c &\leq \|S_{c'} A_c \mu(c)\|_2 \\
&\quad + \|S_{c'} P_c \mathbf{diag}(P_c)\|_2 + 2\|P_c S_{c'} \mathbf{diag}(P_c)\|_2 + 2\|P_c S_{c'} P_c \mathbf{diag}(P_c)\|_2 \\
&\quad + \|P_c \mathbf{diag}(P_c S_{c'})\|_2 + \|P_c \mathbf{diag}(P_c S_{c'} P_c)\|_2 \\
&\leq 8\sqrt{n}.
\end{aligned}$$

Combining both terms, we have the result.  $\square$

### ■ 14.5.8 Mixing time

*Proof of Theorem 14.5.1.* In the last previous sections, we proved that if  $\frac{1}{n} \leq h \leq \frac{1}{256\sqrt{n}}$

1.  $D_0 \leq O(\sqrt{n})$  (Lemma 14.5.11)
2.  $D_1 = O(n\sqrt{h})$  (Lemma 14.5.12)
3.  $D_2 = O(\sqrt{n})$  (Lemma 14.5.16)
4.  $G_1 = O(\sqrt{h})$  (Lemma 14.5.13)
5.  $G_2 = O(\sqrt{m})$  (Lemma 14.5.6)
6.  $R_1 = O(1/\sqrt{n})$  (Lemma 14.5.14)
7.  $R_2 = O(\sqrt{nh})$  (Lemma 14.5.15)

It is easy to see that all these parameters decrease if  $h$  decreases. So, the condition  $\frac{1}{n} \leq h$  can be ignored. We only use randomness to prove that  $\|A_{\gamma(0)}u\|_2 \leq \frac{3}{2}$  and  $\|A_{\gamma}u\|_4 \leq \frac{1.55}{n^{1/4}}$  (Lemma 14.5.10) and they can be checked. These bounds hold with exponentially small failure probability. Hence, the measure of  $\Xi$  is large as claimed and the result follows. Theorem 14.4.2 implies that the walk has mixing time  $O(G_2^2/h)$  as long as

$$h \leq \min \left\{ \frac{1}{(nD_0R_1)^{2/3}}, \frac{1}{D_2}, \frac{1}{nR_1}, \frac{1}{n^{1/3}D_1^{2/3}}, \frac{1}{nG_1^{2/3}}, \frac{1}{(nR_2)^{2/3}} \right\} \leq \frac{1}{\Theta(n^{3/4})}.$$

$\square$

### ■ 14.5.9 Complex Analyticity of geodesics, parallel transport and Jacobi fields

To show geodesics, parallel transport and Jacobi fields on  $M_L$  are complex analytic, we note that the Cauchy–Kowalevski theorem shows that if a differential equation is complex analytic, then the equation has a unique complex analytic solution.

**Theorem 14.5.17** (Simplified Version of Cauchy–Kowalevski theorem). *If  $f$  is a complex analytic function defined on a neighborhood of  $(z_0, \alpha) \in \mathbb{C}^{n+1}$ , then the problem*

$$\frac{dw}{dz} = f(z, w), \quad w(z_0) = \alpha,$$

*has a unique complex analytic solution  $w$  defined on a neighborhood around  $z_0$ .*

*Similarly, for a complex analytic function  $f$  defined in a neighborhood of  $(z_0, \alpha, \beta) \in \mathbb{C}^{2n+1}$ , the ODE*

$$\frac{d^2w}{dz^2} = f\left(z, w, \frac{dw}{dz}\right), \quad w(z_0) = \alpha, \quad \frac{dw}{dz}(z_0) = \beta$$

*has a unique complex analytic solution  $w$  defined in a neighborhood around  $z_0$ .*

As we see in section 14.5.1, the equations for geodesic, parallel transport and Jacobi field involve only rational functions. Since rational functions are complex analytic, the Cauchy-Kowalevski theorem shows that geodesic, orthogonal frame and Jacobi field are complex analytic.

**Lemma 14.5.18.** *Geodesics, parallel transport and Jacobi fields are complex analytic for the Hessian manifold induced by the logarithmic barrier.*

### 14.5.9.1 Radius of Convergence

Next we bound the higher-order derivatives of geodesic, parallel transport and Jacobi field, using techniques developed in Section 14.7 to. Since these are complex analytic, we get a bound for their radius of convergence.

The purpose of these derivative bounds is to show that the solutions of the corresponding ODEs are well-approximated by low-degree polynomials, where the degree of the polynomial grows as  $\log(\frac{1}{\varepsilon})$  for desired accuracy  $\varepsilon$ . The bound on the degree also implies that the Collocation method for solving ODEs is efficient (roughly matrix multiplication time).

### 14.5.9.2 Geodesic

Motivated from the geodesic equation under Euclidean coordinate (14.19), we define the following auxiliary function

$$\bar{F}_\eta(y, x, t) = (A_{x+t\eta}^T A_{x+t\eta})^{-1} A_{x+t\eta}^T s_{x+t\eta, y+\eta}^2. \quad (14.24)$$

The derivative bounds on geodesic rely on the smoothness of this auxiliary function.

**Lemma 14.5.19.** *Under the normalization  $A^T A = I$ ,  $S_x = I$ , we have that*

$$\bar{F}_\eta(y, x, t) \leq_{(0, x, 0)} \frac{4(\|A\eta\|_4 + 1)^2}{1 - \max(8 + 8\|A\eta\|_\infty, 1)t}$$

where  $\bar{F}_\eta$  is defined in (14.24).

*Proof.* By the assumption that  $S_x = I$ , we have that  $\|S_x\|_2 = 1$ . Using  $A^T A = I$ , we have that

$$\begin{aligned} \|DS_{x+t\eta}(x, t)[(d_x, d_t)]\|_2 &= \|\text{Diag}(Ad_x + d_t A\eta)\|_2 \\ &\leq \|Ad_x\|_\infty + \|d_t A\eta\|_\infty \\ &\leq \|d_x\|_2 + |d_t| \|A\eta\|_\infty \\ &\leq (1 + \|A\eta\|_\infty)\|(d_x, d_t)\|_2. \end{aligned}$$

Let  $\beta = 1 + \|A\eta\|_\infty$ . Then, we have that  $S_{x+t\eta} \leq_{(x, t=0)} 1 + \beta t$ .

By using the inverse formula (Lemma 14.7.3),

$$S_{x+t\eta}^{-1} \leq_{(x,0)} \frac{1}{1-\beta t}.$$

Using  $A^T A = I$ , we have that  $A \leq_{(x,0)} 1$  and hence product formula (Lemma 14.7.3) shows that

$$A_{x+t\eta} = S_{x+t\eta}^{-1} A \leq_{(x,0)} \frac{1}{1-\beta t}.$$

Since  $A^T A = I$ , we have that  $AA^T \leq I$  and hence  $A_{x+t\eta}^T \leq_{(x,0)} \frac{1}{1-\beta t}$ . Therefore, we have that

$$A_{x+t\eta}^T A_{x+t\eta} \leq_{(x,0)} \frac{1}{(1-\beta t)^2}.$$

By using the inverse formula again,

$$(A_{x+t\eta}^T A_{x+t\eta})^{-1} \leq_{(x,0)} \frac{1}{2 - \frac{1}{(1-\beta t)^2}} = \frac{(1-\beta t)^2}{2(1-\beta t)^2 - 1}.$$

Hence, we have that

$$(A_{x+t\eta}^T A_{x+t\eta})^{-1} A_{x+t\eta}^T \leq_{(x,0)} \frac{(1-\beta t)^2}{2(1-\beta t)^2 - 1} \frac{1}{1-\beta t} = \frac{1-\beta t}{2(1-\beta t)^2 - 1}.$$

Now, we consider the function  $H(y) = (A(y + \eta))^2$ . Note that

$$\begin{aligned} \|H(y)\|_2 &\leq \|A(y + \eta)\|_4^2, \\ \|DH(y)[d]\|_2 &= 2\|(A(y + \eta))Ad\|_2 \leq 2\|A(y + \eta)\|_4 \|d\|_2, \\ \|D^2H(y)[d, d]\|_2 &\leq 2\|(Ad)^2\|_2 \leq 2\|d\|_2^2. \end{aligned}$$

Therefore, we have that

$$H \leq_{y=0} \|A\eta\|_4^2 + 2\|A\eta\|_4 t + t^2 = (\|A\eta\|_4 + t)^2.$$

Hence, we have that

$$\bar{F}_\eta(y, x, t) = (A_{x+t\eta}^T A_{x+t\eta})^{-1} A_{x+t\eta}^T s_{x+t\eta, y+\eta}^2 \leq_{(0, x, 0)} \frac{(1-\beta t)(\|A\eta\|_4 + t)^2}{2(1-\beta t)^2 - 1} \frac{1}{(1-\beta t)^2}.$$

Let  $\phi(t) = \frac{(\|A\eta\|_4 + t)^2}{(2(1-\beta t)^2 - 1)(1-\beta t)}$  and we write  $\phi(t) = \sum_{k=0}^{\infty} a_k t^k$ . For any  $|z| = \frac{1}{8} \min\left(\frac{1}{\beta}, \|A\eta\|_4 + 8\right)$ , we have that

$$\begin{aligned} |\phi(z)| &\leq 3 \left( \|A\eta\|_4 + \frac{\|A\eta\|_4}{8} + 1 \right)^2 \\ &\leq 4 (\|A\eta\|_4 + 1)^2. \end{aligned}$$

Theorem 14.2.11 shows that

$$\begin{aligned} |a_k| &\leq 4 (\|A\eta\|_4 + 1)^2 \left( 8 \max\left(\beta, (\|A\eta\|_4 + 8)^{-1}\right) \right)^k \\ &\leq 4 (\|A\eta\|_4 + 1)^2 (\max(8\beta, 1))^k \end{aligned}$$

Hence, we can instead bound  $\bar{F}_\eta$  by

$$\bar{F}_\eta \leq_{(0,x,0)} \frac{4(\|A\eta\|_4 + 1)^2}{1 - \max(8 + 8\|A\eta\|_\infty, 1)t}.$$

□

Now, we prove the geodesic has large radius of convergence.

**Lemma 14.5.20.** *Under the normalization  $A^T A = I$ ,  $S_x = I$  and Euclidean coordinate, any geodesic starting at  $x$  satisfies the bound*

$$\|\gamma^{(k)}(0)\|_2 \leq k!c^k$$

for all  $k \geq 2$  where  $c = 512\|A\gamma'(0)\|_4$ .

*Proof.* Recall that under Euclidean coordinate, the geodesic equation is given by

$$\gamma'' = (A_\gamma^T A_\gamma)^{-1} A_\gamma^T s_{\gamma'}^2 \stackrel{\text{def}}{=} F(\gamma', \gamma).$$

So, we have that  $\bar{F}_\eta(y, x, t) = \alpha^2 F(\alpha^{-1}y(t) + \gamma'(0), x(t) + \alpha t\gamma'(0))$  with  $\eta = \alpha\gamma'(0)$ .

Now, we estimate  $\bar{F}_\eta$ . Lemma 14.5.19 shows that

$$\bar{F}_\eta \leq_{(0,x,0)} \frac{4(\|A\eta\|_4 + 1)^2}{1 - \max(8 + 8\|A\eta\|_\infty, 1)t} = \frac{4(\alpha\|A\gamma'(0)\|_4 + 1)^2}{1 - \max(8 + 8\alpha\|A\gamma'(0)\|_\infty, 1)t}.$$

Setting  $\alpha = \|A\gamma'(0)\|_4^{-1}$  and using  $\|A\gamma'(0)\|_\infty \leq \|A\gamma'(0)\|_4$ , we have that

$$\bar{F}_\eta \leq_{(0,x,0)} \frac{16}{1 - 16t}.$$

Lemma 14.7.5 shows that

$$\|\gamma^{(k)}(0)\|_2 \leq \frac{\psi^{(k)}(0)}{\alpha^k}$$

for all  $k \geq 2$  where  $\psi(t)$  is the solution of

$$\psi'(t) = \frac{16}{1 - 16\psi(t)} \text{ with } \psi(0) = 0.$$

Solving it, we get that

$$\psi(t) = \frac{1}{16}(1 - \sqrt{1 - 512t}).$$

By Theorem 14.2.11, we have that

$$\left| \psi^{(k)}(0) \right| \leq k!(512)^k.$$

Hence, we have that

$$\|\gamma^{(k)}(0)\|_2 \leq \frac{k!(512)^k}{\alpha^k} = k!(512)^k \|A\gamma'(0)\|_4^k$$

for all  $k \geq 2$ . □

### 14.5.9.3 Parallel Transport

Motivated from the equation for parallel transport under Euclidean coordinate 14.20, we define the following auxiliary function

$$F(t) = (A_{\gamma(t)}^T A_{\gamma(t)})^{-1} A_{\gamma(t)}^T S_{\gamma'(t)} A_{\gamma(t)}. \quad (14.25)$$

The derivative bounds on parallel transport rely on the smoothness of this auxiliary function.

**Lemma 14.5.21.** *Given a geodesic  $\gamma(t)$ . Under the normalization that  $A^T A = I$ ,  $S_{\gamma(0)} = I$ , we have that*

$$F(t) \leq_0 \frac{c}{2(1-2ct)^2 - (1-ct)^2} \frac{1-ct}{1-2ct}$$

where  $F$  is defined in (14.25) and  $c = 512\|A\gamma'(0)\|_4$ .

*Proof.* Lemma 14.5.20 shows that  $\|\gamma^{(k)}(0)\|_2 \leq k!c^k$  for all  $k \geq 2$ . Therefore, we have that

$$S_{\gamma(t)} \leq_0 1 + t\|A\gamma'(0)\|_\infty + \sum_{k \geq 2} (ct)^k \leq_0 \frac{1}{1-ct}.$$

Using Lemma 14.7.3, we have that

$$A_{\gamma(t)} \leq_0 \frac{1}{2 - \frac{1}{1-ct}} = \frac{1-ct}{1-2ct}$$

and hence

$$(A_{\gamma(t)}^T A_{\gamma(t)})^{-1} A_{\gamma(t)}^T \leq_0 \frac{(1-2ct)(1-ct)}{2(1-2ct)^2 - (1-ct)^2}.$$

Now, we note that

$$\text{Diag}(A\gamma'(t)) \leq_0 \|A\gamma'(0)\|_\infty + \sum_{k \geq 1} \frac{(k+1)!c^{k+1}t^k}{k!} \leq_0 \frac{c}{(1-ct)^2} \quad (14.26)$$

and hence

$$S_{\gamma'} = S_\gamma^{-1} \text{Diag}(A\gamma'(t)) \leq_0 \frac{1-ct}{1-2ct} \frac{c}{(1-ct)^2} = \frac{c}{(1-2ct)(1-ct)}.$$

This gives the result.  $\square$

Now, we prove parallel transport has large radius of convergence.

**Lemma 14.5.22.** *Given a geodesic with  $\gamma(0) = x$ . Let  $v(t)$  be the parallel transport of a unit vector along  $\gamma(t)$ . Under the normalization that  $A^T A = I$ ,  $S_{\gamma(0)} = I$ , we have that*

$$\|v^{(k)}(0)\|_2 \leq k!(16c)^k$$

for all  $k \geq 1$  where  $c = 512\|A\gamma'(0)\|_4$ .

*Proof.* From 14.20, we have that

$$\frac{d}{dt}v(t) = (A_{\gamma(t)}^T A_{\gamma(t)})^{-1} A_{\gamma(t)}^T S_{\gamma'(t)} A_{\gamma(t)} v(t).$$

Let  $u(t) = v(\alpha t)$ , then we have that

$$\begin{aligned} u'(t) &= \alpha v'(\alpha t) \\ &= \alpha (A_{\gamma(\alpha t)}^T A_{\gamma(\alpha t)})^{-1} A_{\gamma(\alpha t)}^T S_{\gamma'(\alpha t)} A_{\gamma(\alpha t)} u(t). \end{aligned}$$

Let  $F(x, t) = \alpha (A_{\gamma(\alpha t)}^T A_{\gamma(\alpha t)})^{-1} A_{\gamma(\alpha t)}^T S_{\gamma'(\alpha t)} A_{\gamma(\alpha t)} x$ . Then, Lemma 14.5.21 shows that

$$F(x, t) \leq_{(X(0),0)} \frac{\alpha c}{2(1-2\alpha ct)^2 - (1-\alpha ct)^2} \frac{1-\alpha ct}{1-2\alpha ct} (1+t).$$



Setting  $\alpha = \frac{1}{8c}$ , we have that

$$F(x, t) \leq_{(X(0), 0)} \frac{\frac{1}{8}}{2(1 - \frac{t}{4})^2 - (1 - \frac{t}{8})^2} \frac{1 - \frac{t}{8}}{1 - \frac{t}{4}} (1 + t) \leq_0 \frac{1}{1 - t}.$$

Lemma 14.7.4 shows that

$$\|u^{(k)}(0)\|_2 \leq \psi^{(k)}(0)$$

for all  $k \geq 1$  where  $\psi(t)$  is the solution of

$$\psi'(t) = \frac{1}{1 - \psi(t)} \text{ with } \psi(0) = 0.$$

Solving it, we get that

$$\psi(t) = 1 - \sqrt{1 - 2t}.$$

By Theorem 14.2.11, we have that for any  $0 \leq t \leq \frac{1}{2}$ , we have that

$$\|u^{(k)}(0)\|_2 \leq \left| \psi^{(k)}(0) \right| \leq k! 2^k$$

for all  $k \geq 1$ . For  $k \geq 1$ , we have that

$$\|v^{(k)}(0)\|_2 \leq k! (16c)^k.$$

□

#### 14.5.9.4 Jacobi field

Motivated from the equation for Jacobi field under orthogonal frame basis (Lemma 14.5.5), we define the following auxiliary function

$$F(t) = X^{-1} (A_\gamma^T A_\gamma)^{-1} (A_\gamma^T S_{\gamma'} P_\gamma S_{\gamma'} A_\gamma - A_\gamma^T \text{diag}(P_\gamma s_{\gamma'}^2) A_\gamma) X. \quad (14.27)$$

The derivative bounds on Jacobi field rely on the smoothness of this auxiliary function.

**Lemma 14.5.23.** *Given a geodesic  $\gamma(t)$  and an orthogonal frame  $\{X_i\}_{i=1}^n$ . Under the normalization that  $A^T A = I$ ,  $S_{\gamma(0)} = I$ , we have that*

$$F(t) \leq_0 \frac{12c^2}{1 - 64ct}$$

where  $F$  is defined in (14.27) and  $c = 512 \|A\gamma'(0)\|_4$ .

*Proof.* We first bound the derivatives of  $s_{\gamma'}^2$ . Using  $\|\gamma^{(k)}(0)\|_2 \leq k! c^k$  for all  $k \geq 2$  (Lemma 14.5.20), we have that

$$\gamma'(t) - \gamma'(0) \leq_0 \frac{d}{dt} \frac{1}{1 - ct} = \frac{c}{(1 - ct)^2}.$$

Using  $A^T A = I$  and  $\text{Diag}(A\gamma'(t)) \leq_0 \frac{c}{(1 - ct)^2}$  (14.26), we have that

$$\text{Diag}(A\gamma'(t)) A(\gamma'(t) - \gamma'(0)) \leq_0 \left( \frac{c}{(1 - ct)^2} \right)^2. \quad (14.28)$$

Next, we note that

$$\text{Diag}(A\gamma'(t)) A\gamma'(0) = \text{Diag}(A\gamma'(0)) A(\gamma'(t) - \gamma'(0)) + \text{Diag}(A\gamma'(0))^2.$$

and hence

$$\text{Diag}(A\gamma'(t))A\gamma'(0) \leq_0 \frac{c^2}{(1-ct)^2} + c^2. \quad (14.29)$$

Combining (14.28) and (14.28), we get

$$\text{Diag}(A\gamma'(t))A(\gamma'(t)) \leq_0 \left( \frac{c}{(1-ct)^2} \right)^2 + \frac{c^2}{(1-ct)^2} + c^2 \leq_0 \frac{3c^2}{(1-ct)^4}$$

In the proof of Lemma 14.5.21, we showed that  $S_\gamma \leq_0 \frac{1}{1-ct}$  and hence

$$\begin{aligned} s_{\gamma'}^2 &= S_\gamma^{-2} \text{Diag}(A\gamma'(t))A(\gamma'(t)). \\ &\leq_0 \left( \frac{1}{2 - \frac{1}{1-ct}} \right)^2 \frac{3c^2}{(1-ct)^4} \\ &= \frac{3c^2}{(1-ct)^2(2(1-ct) - 1)^2}. \end{aligned}$$

In the proof of Lemma 14.5.21, we showed that  $A_\gamma \leq_0 \frac{1-ct}{1-2ct}$  and  $(A_\gamma^T A_\gamma)^{-1} \leq_0 \frac{(1-2ct)^2}{2(1-2ct)^2 - (1-ct)^2}$ . Therefore, we have that

$$P_\gamma = A_\gamma(A_\gamma^T A_\gamma)^{-1} A_\gamma^T \leq_0 \frac{(1-ct)^2}{2(1-2ct)^2 - (1-ct)^2}.$$

Hence, we have

$$\begin{aligned} P_\gamma s_{\gamma'}^2 &\leq_0 \frac{(1-ct)^2}{2(1-2ct)^2 - (1-ct)^2} \frac{3c^2}{(1-ct)^2(1-2ct)^2} \\ &= \frac{3c^2}{(1-2ct)^2(2(1-2ct)^2 - (1-ct)^2)}. \end{aligned}$$

Let  $Y = (A_\gamma^T A_\gamma)^{-1} \left( A_\gamma^T S_{\gamma'} P_\gamma S_{\gamma'} A_\gamma - A_\gamma^T \text{diag}(P_\gamma s_{\gamma'}^2) A_\gamma \right)$ . By a similar proof, we have that

$$\begin{aligned} Y_0 &\leq_0 \frac{(1-2ct)^2}{2(1-2ct)^2 - (1-ct)^2} \left( \frac{1-ct}{1-2ct} \right)^2 \\ &\quad \left( \left( \frac{c}{(1-2ct)(1-ct)} \right)^2 \frac{(1-ct)^2}{2(1-2ct)^2 - (1-ct)^2} + \frac{3c^2}{(1-2ct)^2(2(1-2ct)^2 - (1-ct)^2)} \right) \\ &= \frac{4c^2}{(2(1-2ct)^2 - (1-ct)^2)^2} \left( \frac{1-ct}{1-2ct} \right)^2 \end{aligned}$$

Next, we let  $Z(t) = X(t)v$  for some unit vector  $v$ . Since  $X(t)$  is a parallel transport of  $X(0)$ ,  $Z(t)$  is a parallel transport of  $X(0)v$  and hence Lemma 14.5.22 shows that  $\|Z^{(k)}(t)\|_2 \leq (16c)^k$  for all  $i$ . For any  $k$ , there is unit vector  $v_k$  such that  $\|\frac{d}{dt^k} X(0)\|_2 = \|\frac{d}{dt^k} X(0)v_k\|_2$ . Hence, we have that

$$\|\frac{d}{dt^k} X(0)\|_2 \leq (16c)^k.$$

Therefore, we have that  $X \leq_0 \frac{1}{1-16ct}$ . Since  $X(0) = I$ , we have that  $X^{-1} \leq_0 \frac{1}{2 - \frac{1}{1-16ct}} = \frac{1-16ct}{1-32ct}$ . Thus,

we have that

$$\begin{aligned} F(t) &\leq_0 \frac{4c^2}{(2(1-2ct)^2 - (1-ct)^2)^2} \left( \frac{1-ct}{1-2ct} \right)^2 \frac{1}{1-32ct} \\ &\leq_0 \frac{12c^2}{1-64ct}. \end{aligned}$$

□

Now, we prove Jacobi field has large radius of convergence.

**Lemma 14.5.24.** *Given a geodesic  $\gamma(t)$  and an orthogonal frame  $\{X_i\}_{i=1}^n$  along  $\gamma$ . Let  $V(t)$  be a Jacobi field along  $\gamma(t)$  with  $V(0) = 0$  and let  $U(t)$  be the Jacobi field under the  $X_i$  coordinates, namely,  $V(t) = X(t)U(t)$ . Under the normalization that  $A^T A = I$ ,  $S_{\gamma(0)} = I$ , we have that*

$$\|U^{(k)}(0)\|_2 \leq k! (256c)^k$$

for all  $k \geq 2$  where  $c = 512\|A\gamma'(0)\|_4$ .

*Proof.* Note that  $\frac{d^2U(t)}{dt^2} + F(t)U(t) = 0$  where  $F(t)$  defined in (14.27). Since the equation is linear, we can rescale  $U$  and assume that  $\|U'(0)\|_2 = \frac{1}{\alpha}$ .

Let  $\bar{F}(U(t), t) = \alpha^2 F(\alpha t) (U(t) + \alpha t U'(0))$  with  $\alpha = \frac{1}{64c}$ . Using  $\alpha^2 F(\alpha t) \leq_0 \frac{12c^2\alpha^2}{1-64c\alpha t}$  (Lemma 14.5.23),  $U \leq_{U=0} t$  and  $\alpha t U'(0) \leq_{t=0} \|\alpha U'(0)\|_2 t = t$ , we have that

$$\begin{aligned} \bar{F} &\leq_{(0,0)} \frac{24c^2\alpha^2 t}{1-64c\alpha t} \\ &\leq_0 \frac{t}{1-t}. \end{aligned}$$

Lemma 14.7.5 shows that that

$$\|U^{(k)}(t)\|_2 \leq \frac{\psi^{(k)}(\alpha^{-1}t)}{\alpha^k}$$

for all  $k \geq 2$  where  $\psi(t)$  is the solution of

$$\psi'(t) = \frac{2}{1-\psi(t)} \text{ with } \psi(0) = 0.$$

Solving it, we get that

$$\psi(t) = 1 - \sqrt{1-4t}.$$

By Theorem 14.2.11, we have that for any  $0 \leq t \leq \frac{1}{4}$ , we have that

$$\left| \psi^{(k)}(0) \right| \leq k!(4)^k.$$

Hence, we have that

$$\|U^{(k)}(0)\|_2 \leq \frac{k!4^k}{\alpha^k} \leq k! (256c)^k$$

for all  $k \geq 2$ .

□

## ■ 14.5.10 Implementation

### 14.5.10.1 Computing Geodesic Equation

To apply Theorem 14.6.4, we define

$$F(u, s) = A(A^T S^{-2} A)^{-1} A^T S^{-1} u^2 \quad (14.30)$$

The following lemma bounds the Lipschitz constant of  $F$ .

**Lemma 14.5.25.** *Assuming  $\frac{1}{2} \leq s_i \leq 2$  for all  $i$ . Then, we have*

$$\|DF(u, s)[d_u, d_s]\|_2^2 \leq 48\|u\|_4^2 \|d_u\|_4^2 + 240\|u\|_4^4 \|d_s\|_\infty^2.$$

*Proof.* Let  $A_s = S^{-1}A$  and  $S_d = S^{-1}\text{Diag}(d_s)$ . Then, we have  $F(u, s) = A(A_s^T A_s)^{-1} A_s^T u^2$ . Hence, we have that

$$\begin{aligned} DF(u, s)[d_u, d_s] &= 2A(A_s^T A_s)^{-1} A_s^T S_d A_s (A_s^T A_s)^{-1} A_s^T u^2 \\ &\quad - A(A_s^T A_s)^{-1} A_s^T S_d u^2 \\ &\quad + 2A(A_s^T A_s)^{-1} A_s^T U d_u \end{aligned}$$

Let  $P = A_s (A_s^T A_s)^{-1} A_s^T$ , then, we have that

$$\begin{aligned} \|DF(u, s)[d_u, d_s]\|_2^2 &\leq 12(u^2)^T P S_d P S^2 P S_d P u^2 \\ &\quad + 3(u^2)^T S_d P S^2 P S_d u^2 \\ &\quad + 12d_u^T U P S^2 P U d_u. \end{aligned}$$

Using that  $P \preceq I$ , we have that

$$\begin{aligned} &\|DF(u, s)[d_u, d_s]\|_2^2 \\ &\leq 12\|S\|_\infty^2 (u^2)^T P S_d^2 P u^2 + 3\|S\|_\infty^2 (u^2)^T S_d^2 u^2 + 12\|S\|_\infty^2 \sum_i (d_u)_i^2 u_i^2 \\ &\leq 15\|S\|_\infty^2 \|S_d\|_\infty^2 \|u\|_4^4 + 12\|S\|_\infty^2 \|u\|_4^2 \|d_u\|_4^2 \\ &\leq 15\|S\|_\infty^2 \|S^{-1}\|_\infty^2 \|d_s\|_\infty^2 \|u\|_4^4 + 12\|S\|_\infty^2 \|u\|_4^2 \|d_u\|_4^2. \end{aligned}$$

Now, we use that  $\frac{1}{2} \leq s_i \leq 2$  and get

$$\|DF(u, s)[d_u, d_s]\|_2^2 \leq 48\|u\|_4^2 \|d_u\|_4^2 + 240\|d_s\|_\infty^2 \|u\|_4^4.$$

□

Now, we can apply the collocation method to obtain a good approximation of the solution  $u^*$ .

**Lemma 14.5.26.** *Let  $\gamma$  be a random geodesic generated by the geodesic walk with step size  $h \leq \frac{1}{10^7 \sqrt{n}}$ . In time  $O(mn^{\omega-1} \log^2(1/\varepsilon))$ , we can find  $\bar{\gamma}$  such that  $\max_{0 \leq t \leq \ell} \|\gamma(t) - \bar{\gamma}(t)\|_\infty \leq \varepsilon$  and  $\max_{0 \leq t \leq \ell} \|\gamma'(t) - \bar{\gamma}'(t)\|_\infty \leq \varepsilon$  with probability at least  $1 - O(\exp(-\sqrt{n}/100))$ . Furthermore,  $\bar{\gamma}$  is a  $O(\log(1/\varepsilon))$  degree polynomial.*

*Proof.* Let  $s(t) = A\gamma(t) - b$ . By rotating the space and rescaling the rows of  $A$ , we assume that

$s(0)_i = 1$  for all  $i$  and  $A^T A = I$ . We define  $F$  as (14.30). Then, we have that

$$\begin{aligned} s''(t) &= F(s', s), \\ s'(0) &= A\gamma'(0), \\ s(0) &= 1. \end{aligned}$$

We let  $\alpha = 40\ell$  and

$$\begin{aligned} K &\stackrel{\text{def}}{=} \alpha \|F(A\gamma'(0), 1)\|_4 + \|A\gamma'(0)\|_4 \\ &\leq \alpha \|A(A^T A)^{-1} A^T (A\gamma'(0))^2\|_2 + \|A\gamma'(0)\|_4 \\ &\leq \alpha \|A\gamma'(0)\|_4^2 + \|A\gamma'(0)\|_4. \end{aligned}$$

Using  $\|A_x \gamma'\|_4 \leq 2n^{-1/4}$  (Lemma 14.5.10), we have that

$$\begin{aligned} K &\leq 2\|A\gamma'(0)\|_4 \leq 4n^{-1/4}, \\ \alpha K &\leq 160\ell n^{-1/4} \leq \frac{1}{2}. \end{aligned}$$

For any  $\|u - s'(0)\|_4 \leq K \leq 4n^{-1/4}$ ,  $\|s - 1\|_4 \leq \alpha K \leq \frac{1}{2}$ , Lemma 14.5.25 shows that

$$\begin{aligned} \|DF(u, s)[d_u, d_s]\|_4 &\leq \|DF(u, s)[d_u, d_s]\|_2 \\ &\leq 8\|u\|_4 \|d_u\|_4 + 16\|u\|_4^2 \|d_s\|_\infty \\ &\leq 48n^{-1/4} \|d_u\|_4 + 576n^{-1/2} \|d_s\|_4. \end{aligned}$$

Therefore, for any  $\|u_1 - s'(0)\|_4 \leq K$ ,  $\|s_1 - 1\|_4 \leq \alpha K$ ,  $\|u_2 - s'(0)\|_4 \leq K$ ,  $\|s_2 - 1\|_4 \leq \alpha K$ , we have that

$$\begin{aligned} \|F(u_1, s_1) - F(u_2, s_2)\|_4 &\leq 48n^{-1/4} \|u_1 - u_2\|_4 + 576n^{-1/2} \|s_1 - s_2\|_4 \\ &\leq \frac{1}{\alpha} \|u_1 - u_2\|_4 + \frac{1}{\alpha^2} \|s_1 - s_2\|_4. \end{aligned}$$

Since  $\gamma$  is analytic and  $\|\gamma^{(k)}(0)\|_2 = O(k!n^{-k/4})$  (Lemma 14.5.22),  $\gamma(t)$  is  $\varepsilon$  close to a polynomial with degree  $O(\log(1/\varepsilon))$  for  $0 \leq t \leq o(n^{1/4})$ . Hence, we can apply Theorem 14.6.6 and find  $\bar{\gamma}$  such that  $\|\gamma - \bar{\gamma}\|_4 \leq \varepsilon$  and  $\|\gamma' - \bar{\gamma}'\|_4 \leq \varepsilon$  in  $O(n \log^3(K/\varepsilon))$  time plus  $O(\log^2(K/\varepsilon))$  evaluations of  $F$ . Note that each evaluation of  $F$  involves solving a linear system and hence it takes  $O(mn^{\omega-1})$ . Therefore, the total running time is  $O(mn^{\omega-1} \log^2(1/\varepsilon))$ .  $\square$

### 14.5.10.2 Computing Parallel Transport

**Lemma 14.5.27.** *Given  $\gamma$  be a random geodesic generated by the geodesic walk with step size  $h \leq \frac{1}{10^7 \sqrt{n}}$  and an unit vector  $v$ . Let  $v(t)$  be the parallel transport of a unit vector along  $\gamma(t)$ . In time  $O(mn^{\omega-1} \log^2(1/\varepsilon))$ , we can find  $\bar{v}$  such that  $\max_{0 \leq t \leq \ell} \|v(t) - \bar{v}(t)\|_\infty \leq \varepsilon$  with probability at least  $1 - O(\exp(-\sqrt{n}/100))$ . Furthermore,  $\bar{v}$  is a  $O(\log(1/\varepsilon))$  degree polynomial.*

*Similarly, given a basis  $\{v_i\}_{i=1}^n$ , in time  $O(mn^{\omega-1} \log^2(1/\varepsilon))$ , we can find an approximate parallel transport  $\bar{v}_i(t)$  of  $\{v_i\}_{i=1}^n$  along  $\gamma(t)$  with probability at least  $1 - O(\exp(-\sqrt{n}/100))$ .*

*Proof.* Recall that the equation for parallel transport (14.20) is given by

$$\frac{d}{dt} v(t) = \left( A_{\gamma(t)}^T A_{\gamma(t)} \right)^{-1} A_{\gamma(t)}^T S_{\gamma'(t)} A_{\gamma(t)} v.$$

By rotating the space and rescaling the rows of  $A$ , we assume that  $s(\gamma(0))_i = 1$  for all  $i$  and  $A^T A = I$ .

In the proof of Lemma 14.5.26, we know that  $\frac{1}{2} \leq s(\gamma(t))_i \leq 2$  for all  $0 \leq t \leq \ell$ . For any unit vector  $u$ , we have

$$\begin{aligned} \|(A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'} A_\gamma u\|_2 &\leq 2 \|(A_\gamma^T A_\gamma)^{-1/2} A_\gamma^T S_{\gamma'} A_\gamma u\|_2 \\ &\leq 2 \|S_{\gamma'}\|_\infty \|A_\gamma u\|_2 \\ &\leq 4 \|S_{\gamma'}\|_\infty \leq 24\sqrt{h} \end{aligned}$$

where we used Lemma 14.5.10 in the last line. Using  $h \leq \frac{1}{10^7 \sqrt{n}}$ , we have that

$$\|(A_\gamma^T A_\gamma)^{-1} A_\gamma^T S_{\gamma'} A_\gamma\|_2 \leq \frac{1}{20\ell}.$$

Since  $v$  is analytic and  $\|v^{(k)}(0)\|_2 = O(k!n^{-k/4})$  (Lemma 14.5.22),  $v(t)$  is  $\varepsilon$  close to a polynomial with degree  $O(\log(1/\varepsilon))$  for  $0 \leq t \leq o(n^{1/4})$ . Hence, we can apply Theorem 14.6.4 and find  $\bar{v}$  such that  $\|v - \bar{v}\|_2 \leq \varepsilon$  in  $O(n \log^3(1/\varepsilon))$  time plus  $O(\log^2(1/\varepsilon))$  evaluations of  $F$ . Note that each evaluation of  $F$  involves solving a linear system and hence it takes  $O(mn^{\omega-1})$ . Therefore, the total running time is  $O(mn^{\omega-1} \log^2(1/\varepsilon))$ .

For the last result, we note that each evaluation of  $F$  becomes computing matrix inverse and performing matrix multiplication and they can be done in again  $O(mn^{\omega-1})$  time.  $\square$

### 14.5.10.3 Computing Jacobi field

**Lemma 14.5.28.** *Given  $\gamma$  be a random geodesic generated by the geodesic walk with step size  $h \leq \frac{1}{10^7 \sqrt{n}}$ . Let  $\{X_i\}_{i=1}^n$  be an orthogonal frame along  $\gamma$ . Let  $v(t)$  be a Jacobi field along  $\gamma(t)$  with  $v(0) = 0$  and let  $u(t)$  be the Jacobi field under the  $X_i$  coordinates, namely,  $v(t) = X(t)u(t)$ . In time  $O(mn^{\omega-1} \log^2(1/\varepsilon))$ , we can find  $\bar{u}$  such that  $\max_{0 \leq t \leq \ell} \|u(t) - \bar{u}(t)\|_\infty \leq \varepsilon$  with probability at least  $1 - O(\exp(-\sqrt{n}/100))$ . Furthermore,  $\bar{u}$  is a  $O(\log(1/\varepsilon))$  degree polynomial.*

*Proof.* Recall that the equation for Jacobi field 14.5.5 is given by

$$\frac{d^2 u}{dt^2} + X^{-1} (A_\gamma^T A_\gamma)^{-1} (A_\gamma^T S_{\gamma'} P_\gamma S_{\gamma'} A_\gamma - A_\gamma^T \text{Diag}(P_\gamma s_{\gamma'}^2) A_\gamma) X u = 0.$$

By rotating the space and rescaling the rows of  $A$ , we assume that  $s(\gamma(0))_i = 1$  for all  $i$  and  $A^T A = I$ . In the proof of Lemma 14.5.26, we know that  $\frac{1}{2} \leq s(\gamma(t))_i \leq 2$  for all  $0 \leq t \leq \ell$ . Since  $X$  is an orthogonal frame, we have  $XX^T = XX^T = I$ . Hence, for any unit vector  $v$ , we have

$$\begin{aligned} &\|X^{-1} (A_\gamma^T A_\gamma)^{-1} (A_\gamma^T S_{\gamma'} P_\gamma S_{\gamma'} A_\gamma - A_\gamma^T \text{Diag}(P_\gamma s_{\gamma'}^2) A_\gamma) X v\|_2 \\ &\leq 2 \|(A_\gamma^T A_\gamma)^{-1/2} (A_\gamma^T S_{\gamma'} P_\gamma S_{\gamma'} A_\gamma - A_\gamma^T \text{Diag}(P_\gamma s_{\gamma'}^2) A_\gamma) X v\|_2 \\ &\leq 2 \|S_{\gamma'} P_\gamma S_{\gamma'} A_\gamma v\|_2 + 2 \|\text{Diag}(P_\gamma s_{\gamma'}^2) A_\gamma X v\|_2. \end{aligned}$$

Using  $\|S_{\gamma'}\|_2 \leq 6\sqrt{h}$  and  $\|\text{Diag}(P_\gamma s_{\gamma'}^2)\|_2 \leq \|P_\gamma s_{\gamma'}^2\|_2 \leq \|s_{\gamma'}\|_4^2 \leq 4n^{-1/2}$  (Lemma 14.5.10), we have that

$$\begin{aligned} &\|X^{-1} (A_\gamma^T A_\gamma)^{-1} (A_\gamma^T S_{\gamma'} P_\gamma S_{\gamma'} A_\gamma - A_\gamma^T \text{Diag}(P_\gamma s_{\gamma'}^2) A_\gamma) X v\|_2 \\ &\leq 12\sqrt{h} \|P_\gamma S_{\gamma'} A_\gamma v\|_2 + 8n^{-1/2} \|A_\gamma X v\|_2 \\ &\leq 72h \|A_\gamma v\|_2 + 16n^{-1/2} \|X v\|_2 \\ &\leq 144h + 16n^{-1/2} \leq \frac{1}{(40\ell)^2} \end{aligned}$$

where we used  $h \leq \frac{1}{10^7 \sqrt{n}}$  in the last line.

Since  $u$  is analytic and  $\|u^{(k)}(0)\|_2 = O(k!n^{-k/4})$  (Lemma 14.5.24),  $u(t)$  is  $\varepsilon$  close to a polynomial with degree  $O(\log(1/\varepsilon))$  for  $0 \leq t \leq o(n^{1/4})$ . Hence, we can apply Theorem 14.6.4 and find  $\bar{u}$  such that  $\|u - \bar{u}\|_2 \leq \varepsilon$  in  $O(n \log^3(1/\varepsilon))$  time plus  $O(\log^2(1/\varepsilon))$  evaluations of  $F$ . Note that each evaluation of  $F$  involves computing matrix inversions and matrix multiplications and hence it takes  $O(mn^{\omega-1})$ . Therefore, the total running time is  $O(mn^{\omega-1} \log^2(1/\varepsilon))$ .  $\square$

#### 14.5.10.4 Computing Geodesic Walk

*Proof of Theorem 14.5.2.* To implement the geodesic walk, we use Lemma 14.5.26 to compute the geodesic, Lemma 14.5.27 to compute an orthogonal frame along the geodesic and Lemma 14.5.28 to compute the Jacobi field along. Using the Jacobi field, we can use (14.9) and Lemma 14.4.12 to compute the probability from  $x$  to  $y$  and the probability from  $y$  to  $x$ . Using these probabilities, we can implement the rejection sampling. It suffices to compute the geodesic and the probability up to  $1/n^{O(1)}$  accuracy and hence these operations can be done in time  $O(mn^{\omega-1} \log^2(n))$ .

Note that we only use randomness to prove that  $\|A_{\gamma(0)}u\|_2 \leq \frac{3}{2}$  and  $\|A_{\gamma}u\|_4 \leq \frac{1.55}{n^{1/4}}$  (Lemma 14.5.10) and they can be checked. When we condition our walk to that, we only change the distribution by exponential small amount. Hence, this result is stated without mentioning the success probability.  $\square$

## ■ 14.6 Collocation Method for ODE

In this section, we study a collocation method for solving ordinary differential equation (ODE) and show how to solve a nice enough ODE in nearly constant number of iterations without explicitly computing higher derivatives.

### ■ 14.6.1 First Order ODE

We first consider the following first order ODE

$$\begin{aligned} \frac{d}{dt}u(t) &= F(u(t), t), \text{ for } 0 \leq t \leq \ell \\ u(0) &= v \end{aligned} \tag{14.31}$$

where  $F : \mathbb{R}^{n+1} \rightarrow \mathbb{R}^n$  and  $u(t) \in \mathbb{R}^n$ . The idea of collocation methods is to find a degree  $d$  polynomial  $p$  such that

$$\begin{aligned} \frac{d}{dt}p(t) &= F(p(t), t), \text{ for } t = c_1, c_2, \dots, c_d \\ p(0) &= v \end{aligned} \tag{14.32}$$

where  $c_1, c_2, \dots, c_d$  are carefully chosen distinct points on  $[0, \ell]$ . Here, we call  $p : \mathbb{R} \rightarrow \mathbb{R}^n$  is a degree  $d$  polynomial if  $p(t) = [p_1(t); p_2(t); \dots; p_n(t)]$  and each  $p_i(t)$  is an univariate polynomial with degree at most  $d$ . The first part of the proof shows the existence of a solution for the systems (14.32). To describe the algorithm, it is easier to consider an equivalent integral equation.

**Lemma 14.6.1.** *Given distinct points  $c_1, c_2, \dots, c_d \in \mathbb{R}$  and  $F : \mathbb{R}^{n+1} \rightarrow \mathbb{R}^n$ , consider the nonlinear*

map  $T : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{d \times n}$  defined by

$$T(\zeta)_{(i,k)} = \int_0^{c_i} \sum_{j=1}^d F(\zeta_j, c_j)_k \phi_j(s) ds \text{ for } i \in [d], k \in [n]$$

where  $\phi_i(s) = \prod_{j \neq i} \frac{s-c_j}{c_i-c_j}$  are the Lagrange basis polynomials. Given any  $\zeta \in \mathbb{R}^{d \times n}$  such that

$$\zeta_i = v + T(\zeta)_i, \text{ for } i \in [d] \quad (14.33)$$

the polynomial

$$p(t) = v + \int_0^t \sum_{j=1}^d F(\zeta_j, c_j) \phi_j(s) ds$$

is a solution of the system (14.32).

*Proof.* Define the polynomials  $\phi_i(s) = \prod_{j \neq i} \frac{s-c_j}{c_i-c_j}$ . Note that  $\phi_i(c_j) = \delta_{ij}$ . Therefore, we have

$$\sum_{j=1}^d \alpha_j \phi_j(c_i) = \alpha_i.$$

Therefore,  $p(0) = v$  and

$$\frac{d}{dt} p(c_i) = \sum_{j=1}^d F(\zeta_j, c_j) \phi_j(c_i) = F(\zeta_j, c_j).$$

Since  $\zeta_i = v + T(\zeta)_i$ , we have that

$$\zeta_i = v + \int_0^{c_i} \sum_{j=1}^d F(\zeta_j, c_j) \phi_j(s) ds = p(c_i).$$

Hence, we have  $\frac{d}{dt} p(c_i) = F(p(c_i), c_i)$ . Therefore,  $p$  is a solution to the system (14.32).  $\square$

From Lemma 14.6.1, we see that it suffices to solve the system (14.33). We solve it by a simple fix point iteration shown in Algorithm 50.

---

**Algorithm 50: CollocationMethod**

---

**Input:** An ordinary differential equation  $\frac{d}{dt} u(t) = F(u(t), t)$  for  $0 \leq t \leq \ell$  with initial condition  $u(0) = v$ .

Define  $T(\zeta)_{(i,k)} = \int_0^{c_i} \sum_{j=1}^d F(\zeta_j, c_j)_k \phi_j(s) ds$  for  $i \in [d], k \in [n]$  as defined in Lemma 14.6.1.

$c_i = \frac{\ell}{2} + \frac{\ell}{2} \cos(\frac{2i-1}{2d}\pi)$  for all  $i \in [d]$ ,  $K = 40\ell \max_{t \in [0, \ell]} \|F(v, t)\|_p$ ,  $Z = \log_2(K/\varepsilon)$

$\zeta_i^{(0)} = v + T(\bar{v})_i$  for all  $i \in [d]$  where  $\bar{v} = (v, v, \dots, v) \in \mathbb{R}^{d \times n}$ .

**for**  $z = 0, \dots, Z$  **do**

    |  $\zeta_i^{(z+1)} = v + T(\zeta^{(z)})_i$  for  $i \in [d]$ .

**end**

**Output:**  $p(t) = v + \int_0^t \sum_{j=1}^d F(\zeta_j^{(Z+1)}, c_j) \phi_j(s) ds$ .

---

In the following lemma, we show that  $T$  is a contraction mapping if  $F$  is smooth enough and hence this algorithm converges linearly. We will use the following norm:  $\|x\|_{\infty; p} = \max_{i \in [d]} \left( \sum_{k \in [n]} |x_{(i,k)}|^p \right)^{1/p}$ .



**Lemma 14.6.2.** *Given  $x, y \in \mathbb{R}^{d \times n}$  with  $d \geq 2$ , suppose that  $\|F(x_i, t) - F(y_i, t)\|_p \leq L\|x_i - y_i\|_p$  for all  $0 \leq t \leq \ell$  and all  $i \in [d]$ . For  $c_k = \frac{\ell}{2} + \frac{\ell}{2} \cos(\frac{2k-1}{2d}\pi)$  and  $T$  defined Lemma 14.6.1, we have that*

$$\|T(x) - T(y)\|_{\infty; p} \leq 10\ell L\|x - y\|_{\infty; p}.$$

Also, we have that

$$\|T(x)\|_{\infty; p} \leq 10\ell \max_i \|F(x_i, c_i)\|_{\infty; p}$$

*Proof.* Using the definition of  $T$  and  $\|\cdot\|_{\infty; p}$ , we have

$$\begin{aligned} \|T(x) - T(y)\|_{\infty; p} &= \max_i \left\| \int_0^{c_i} \sum_{j=1}^d F(x_j, c_j) \phi_j(s) ds - \int_0^{c_i} \sum_{j=1}^d F(y_j, c_j) \phi_j(s) ds \right\|_p \\ &\leq \left( \max_i \sum_{j=1}^d \left| \int_0^{c_i} \phi_j(s) ds \right| \right) \left( \max_i \|F(x_i, c_i) - F(y_i, c_i)\|_p \right) \\ &\leq L\|x - y\|_{\infty; p} \left( \max_i \sum_{j=1}^d \left| \int_0^{c_i} \phi_j(s) ds \right| \right). \end{aligned}$$

To bound  $\max_i \sum_{j=1}^d \left| \int_0^{c_i} \phi_j(s) ds \right|$ , it is easier to work with the function  $\psi_j(x) \stackrel{\text{def}}{=} \phi_j(\frac{2}{\ell}x - 1)$ . Note that  $\psi_k$  is the Lagrange basis polynomials on the nodes  $\{\cos(\frac{2k-1}{2d}\pi)\}_{k=1}^d$  and hence we have

$$\psi_k(x) = \frac{(-1)^{k-1} \sqrt{1 - x_k^2} \cos(d \cos^{-1} x)}{d(x - x_k)}$$

where  $x_k = \cos(\frac{2k-1}{2d}\pi)$ . Lemma 14.6.5 at the end of this section shows that  $\left| \int_{-1}^t \psi_k(x) dx \right| \leq \frac{20}{d}$  for all  $t$ . Therefore, we have that  $\max_i \sum_{j=1}^d \left| \int_0^{c_i} \phi_j(s) ds \right| \leq 10\ell$ . This gives the first inequality. The second inequality is similar.  $\square$

In each iteration of the collation method, we need to compute  $\int_0^{c_i} \sum \alpha_j \phi_j(s) ds$  for some  $\alpha_j$ . The following theorem shows that this can be done in  $O(d \log(d/\varepsilon))$  time using the multipole method.

**Theorem 14.6.3** ([77, Sec 5]). *Let  $\phi_i(s)$  be the Lagrange basis polynomials on the Chebyshev nodes on  $[0, \ell]$ , namely,  $\phi_i(s) = \prod_{j \neq i} \frac{s - c_j}{c_i - c_j}$  with  $c_i = \frac{\ell}{2} + \frac{\ell}{2} \cos(\frac{2i-1}{2d}\pi)$ . Given a polynomial  $p(s) = \sum_j \alpha_j \phi_j(s)$  and a point set  $\{x_1, x_2, \dots, x_d\}$ , one can compute  $t_i$  such that*

$$\left| t_i - \int_0^{x_i} p(s) ds \right| \leq \varepsilon \ell \sqrt{\sum_{i \neq j} (\alpha_i - \alpha_j)^2} \text{ for } i \in [d]$$

in time  $O(d \log(d/\varepsilon))$ .

Now we have everything to state our main result in this subsection.

**Theorem 14.6.4.** *Let  $u(t) \in \mathbb{R}^n$  be the solution of the ODE (14.31). Suppose we are given  $\varepsilon, \ell > 0$ ,  $d \geq 2$  and  $1 \leq p \leq \infty$  such that*

1. *There is a degree  $d$  polynomial  $q$  from  $\mathbb{R}$  to  $\mathbb{R}^n$  such that  $q(0) = v$  and  $\|\frac{d}{dt}u(t) - \frac{d}{dt}q(t)\|_p \leq \frac{\varepsilon}{\ell}$  for all  $0 \leq t \leq \ell$ .*
2. *We have that  $\|F(x, t) - F(y, t)\|_p \leq \frac{1}{20\ell}\|x - y\|_p$  for all  $\|x - v\|_p \leq K$  and  $\|y - v\|_p \leq K$  with  $K = 40\ell \max_{t \in [0, \ell]} \|F(v, t)\|_p$ .*

Then, for any  $t$  such that  $0 \leq t \leq \ell$ , Algorithm *CollocationMethod* outputs a degree  $d$  polynomial  $p(t)$  such that  $\|u(t) - p(t)\|_p \leq 43\varepsilon$  in  $O(nd \log^2(dK/\varepsilon))$  time and  $O(d \log(K/\varepsilon))$  evaluations of  $F$ .

*Proof.* First, we estimate the initial error. Let  $\zeta^{(\infty)}$  be the solution to (14.33),  $\zeta^{(0)} = v + T(\bar{v})_i$  be the initial vector and  $\bar{v} = (v, v, \dots, v) \in \mathbb{R}^{d \times n}$ . Then, we have that

$$\begin{aligned} \|\zeta^{(\infty)} - \zeta^{(0)}\|_{\infty;p} &= \max_i \|(v + T(\zeta^{(\infty)})_i) - (v + T(\bar{v})_i)\|_p \\ &= \|T(\zeta^{(\infty)}) - T(\bar{v})\|_{\infty;p} \\ &\leq \frac{1}{2} \|\zeta^{(\infty)} - \bar{v}\|_{\infty;p} \\ &\leq \frac{1}{2} \left( \|\zeta^{(\infty)} - \zeta^{(0)}\|_{\infty;p} + \|\zeta^{(0)} - \bar{v}\|_{\infty;p} \right). \end{aligned}$$

Therefore, we have that

$$\begin{aligned} \|\zeta^{(\infty)} - \zeta^{(0)}\|_{\infty;p} &\leq \|\zeta^{(0)} - \bar{v}\|_{\infty;p} \\ &= \|T((v, v, \dots, v))\|_{\infty;p} \\ &\leq 10\ell \max_i \|F(v, c_i)\|_p \\ &\leq \frac{K}{4}. \end{aligned}$$

Hence, we have that  $\|\zeta^{(\infty)} - \bar{v}\|_{\infty;p} \leq \frac{K}{2}$ .

Using the assumption on  $F$ , Lemma 14.6.2 shows that  $\|T(x) - T(y)\|_{\infty;p} \leq \frac{1}{2} \|x - y\|_{\infty;p}$  and hence

$$\|\zeta^{(\infty)} - \zeta^{(k)}\|_{\infty;p} \leq \frac{K}{2^{k+1}}.$$

Thus, it takes  $\log_2(K/\varepsilon)$  iteration to get a point  $\zeta^{(k)}$  with

$$\|\zeta^{(\infty)} - \zeta^{(k)}\|_{\infty;p} \leq \varepsilon. \quad (14.34)$$

Also, this shows that  $\|\zeta^{(\infty)} - \bar{v}\|_{\infty;p} \leq K$ . Hence, we only requires the assumption (2) to be satisfied in this region.

Now, we show that  $\zeta^{(k)}$  is close to the solution by using the existence of  $q$ . By the assumption on  $q$ , we have that  $\|\frac{d}{dt}u(t) - \frac{d}{dt}q(t)\|_p \leq \frac{\varepsilon}{\ell}$  and hence  $\|u(t) - q(t)\|_p \leq \varepsilon$ . Using the smoothness of  $F$ , we have that  $\|F(u(t), t) - F(q(t), t)\|_p \leq \frac{\varepsilon}{\ell}$  for all  $0 \leq t \leq \ell$ . Therefore, we have that

$$\begin{aligned} \left\| \frac{d}{dt}q(t) - F(q(t), t) \right\|_p &\leq \left\| \frac{d}{dt}q(t) - \frac{d}{dt}u(t) \right\|_p + \|F(q(t), t) - F(u(t), t)\|_p \\ &\leq 2\frac{\varepsilon}{\ell} \end{aligned}$$

for all  $0 \leq t \leq 1$ . Therefore, we have

$$\begin{aligned} q(t) &= v + \int_0^t \sum_{j=1}^d \frac{d}{dt}q(c_j) \phi_j(s) ds \\ &= v + \int_0^t \sum_{j=1}^d (F(q(c_j), c_j) + \delta_j) \phi_j(s) ds \end{aligned}$$

where  $\|\delta_i\|_p \leq 2\frac{\varepsilon}{\ell}$  for all  $i \in [d]$ . By Lemma 14.6.2, we have that

$$\|q(t) - v - \int_0^t \sum_{j=1}^d F(q(c_j), c_j) \phi_j(s) ds\|_p \leq 20\varepsilon.$$

Now, we compare  $\bar{q}_j \stackrel{\text{def}}{=} q(c_j)$  with the approximate solution constructed by the fixpoint algorithm

$$q^{(k)}(t) = v + \int_0^t \sum_{j=1}^d F(\zeta_j^{(k)}, c_j) \phi_j(s) ds.$$

For  $0 \leq t \leq \ell$ , we have

$$\begin{aligned} \|q(t) - q^{(k)}(t)\|_p &\leq \|T(\bar{q}) - T(\zeta^{(k)})\|_{\infty;p} + 20\varepsilon \\ &\leq \frac{1}{2} \|\bar{q} - \zeta^{(k)}\|_{\infty;p} + 20\varepsilon. \end{aligned} \quad (14.35)$$

Setting  $t = c_i$ , we have

$$\begin{aligned} \|\bar{q} - \zeta^{(k+1)}\|_{\infty;p} &= \max_i \|\bar{q}(c_i) - q^{(k)}(c_i)\|_p \\ &\leq \frac{1}{2} \|q - \zeta^{(k)}\|_{\infty;p} + 20\varepsilon. \end{aligned}$$

Since  $\|\zeta^{(k+1)} - \zeta^{(k)}\|_{\infty;p} \leq \|\zeta^{(k+1)} - \zeta^{(\infty)}\|_{\infty;p} + \|\zeta^{(k)} - \zeta^{(\infty)}\|_{\infty;p} \leq 2\varepsilon$ , we have

$$\|\bar{q} - \zeta^{(k)}\|_{\infty;p} \leq \frac{1}{2} \|\bar{q} - \zeta^{(k)}\|_{\infty;p} + 2\varepsilon + 20\varepsilon$$

Therefore, we have

$$\|\bar{q} - \zeta^{(k)}\|_{\infty;p} \leq 44\varepsilon.$$

Putting it into (14.35), we have

$$\|q(t) - q^{(k)}(t)\|_p \leq 42\varepsilon$$

for all  $0 \leq t \leq \ell$ . Using that  $\|u(t) - q(t)\|_p \leq \varepsilon$ , we have

$$\|u(t) - q^{(k)}(t)\|_p \leq 43\varepsilon.$$

Hence, it proves the guarantee.

Each iteration involves computing  $v_k + \int_0^{c_i} \sum_{j=1}^d F(\zeta_j^{(z)}, c_j)_k \phi_j(s) ds$  for all  $i \in [d], k \in [n]$ . Note that  $\sum_{j=1}^d F(\zeta_j^{(z)}, c_j)_k \phi_j(s)$  is a polynomial expressed by Lagrange polynomials. Theorem 14.6.3 shows they can be computed in  $O(d \log(dK/\varepsilon))$  with  $\frac{\varepsilon}{Kd^{O(1)}}$  accuracy. Since there are  $n$  coordinates, it takes  $O(nd \log(dK/\varepsilon))$  time plus  $d$  evaluation per iteration.  $\square$

The theorem above essentially says that if the solution is well approximated by a polynomial and if the  $F$  has small enough Lipschitz constant, then we can reconstruct the solution efficiently. Note that this method is not useful for stochastic differential equation because Taylor expansion of the solution involves the high moments of probability distributions which is very expensive to store.

We conclude this subsection with a lemma used in the proof of Lemma 14.6.2.

**Lemma 14.6.5.** *Given  $d \geq 2$ . Let  $\psi(x) = \frac{\sqrt{1-y^2} \cos(d \cos^{-1} x)}{d(x-y)}$  where  $y = \cos(\frac{2k-1}{2d}\pi)$  for some integer  $k \in [d]$ . Then, we have that*

$$\left| \int_{-1}^t \psi(x) dx \right| \leq \frac{20}{d}.$$

for any  $-1 \leq t \leq 1$ .

*Proof.* Let  $z_j = \cos(j\pi/d)$ . For  $j < k$ , we have that

$$\begin{aligned} \left| \int_{z_{j-1}}^{z_j} \psi(x) dx \right| &\leq \frac{\sqrt{1-y^2}}{d} \left| \int_{z_{j-1}}^{z_j} \frac{\cos(d \cos^{-1} x)}{z_j - y} dx \right| + \frac{\sqrt{1-y^2}}{d} \int_{z_{j-1}}^{z_j} \left| \frac{\cos(d \cos^{-1} x) (z_j - x)}{(z_j - y)(x - y)} \right| dx \\ &\leq \frac{\sqrt{1-y^2}}{d} \left| \int_{z_{j-1}}^{z_j} \frac{\cos(d \cos^{-1} x)}{z_j - y} dx \right| + \frac{\sqrt{1-y^2}}{d} \frac{(z_j - z_{j-1})^2}{(z_j - y)^2} \\ &= \frac{\sqrt{1-y^2}}{d} \frac{z_j - z_{j-1}}{d^2 - 1} + \frac{\sqrt{1-y^2}}{d} \frac{(z_j - z_{j-1})^2}{(z_j - y)^2} \\ &\leq \frac{\pi}{d^2(d^2 - 1)} + \frac{4}{d} \frac{1}{(k - \frac{1}{2} - j)^2}. \end{aligned}$$

Hence, we have that

$$\begin{aligned} \left| \int_{-1}^t \psi(x) dx \right| &\leq \sum_{k \neq j} \left| \int_{z_{j-1}}^{z_j} \psi(x) dx \right| + \int_{z_{k-1}}^{z_k} |\psi(x)| dx \\ &\leq \frac{6}{d} + \int_{z_{k-1}}^{z_k} |\psi(x)| dx. \end{aligned}$$

By simple calculation, one can check that  $|\psi(x)| \leq 2$  and hence

$$\left| \int_{-1}^t \psi(x) dx \right| \leq \frac{20}{d}.$$

□

## ■ 14.6.2 Second Order ODE

Now, we consider the following second order ODE

$$\begin{aligned} \frac{d^2}{dt^2} u(t) &= F\left(\frac{d}{dt} u(t), u(t), t\right), \text{ for } 0 \leq t \leq \ell \\ \frac{d}{dt} u(0) &= w, \\ u(0) &= v \end{aligned} \tag{14.36}$$

where  $F : \mathbb{R}^{2n+1} \rightarrow \mathbb{R}^n$  and  $u(t) \in \mathbb{R}^n$ . Using a standard reduction from second order ODE to first order ODE, we show how to apply our first order ODE method to second order ODE.

**Theorem 14.6.6.** *Let  $x(t) \in \mathbb{R}^n$  be the solution of the ODE (14.36). Given some  $\varepsilon, \ell > 0$ ,  $d \geq 2$  and  $1 \leq p \leq \infty$ , let  $\alpha = 40\ell$  and suppose that*

1. *There is a degree  $d$  polynomial  $q$  from  $\mathbb{R}$  to  $\mathbb{R}^n$  such that  $q(0) = v$ ,  $q'(0) = w$ ,  $\|\frac{d}{dt^2} x(t) - \frac{d}{dt^2} q(t)\|_p \leq \frac{\varepsilon}{\ell^2}$  for all  $0 \leq t \leq \ell$ .*
2. *We have that  $\|F(x, \gamma, t) - F(y, \eta, t)\|_p \leq \frac{1}{\alpha} \|x - y\|_p + \frac{1}{\alpha^2} \|\gamma - \eta\|_p$  for all  $\|x - w\|_p \leq K$ ,  $\|y - w\|_p \leq K$ ,  $\|\gamma - v\|_p \leq \alpha K$ ,  $\|\eta - v\|_p \leq \alpha K$  where  $K = \alpha \max_{t \in [0, \ell]} \|F(w, v, t)\|_p + \|w\|_p$ .*

*Then, for any  $t$  such that  $0 \leq t \leq \ell$ , in  $O(nd \log^2(dK/\varepsilon))$  time plus  $O(d \log(K/\varepsilon))$  evaluations of  $F$ , we can find  $p(t), p'(t)$  such that  $\|u(t) - p(t)\|_p = O(\varepsilon)$  and  $\|u'(t) - p'(t)\|_p \leq O(\varepsilon/\ell)$ .*

*Proof.* Let  $\alpha = 40\ell$ . Let  $x(t) = (\alpha u'(\alpha t), u(\alpha t)) \in \mathbb{R}^{2n}$ . Note that  $x(t)$  satisfies the following ODE

$$\begin{aligned} \frac{d}{dt}x(t) &= \bar{F}(x(t), t) \text{ for } 0 \leq t \leq \ell \\ x(0) &= (\alpha w, v). \end{aligned} \quad (14.37)$$

where  $\bar{F}(x(t), t) = (\alpha^2 F(\alpha^{-1}x_{(1)}(t), x_{(2)}(t), \alpha t), x_{(1)}(t))$ ,  $x_{(1)}(t)$  is the first  $n$  variables of  $x(t)$  and  $x_{(2)}(t)$  is the last  $n$  variables of  $x(t)$ . Next, we verify the conditions of Theorem 14.6.4 for this ODE. Let  $\bar{\ell} = \frac{1}{60}$  and  $\bar{K} = 40\bar{\ell} \max_{t \in [0, \ell]} \|\bar{F}(x(0), t)\|_p$ . Note that

$$\bar{K} \leq \alpha^2 \max_{t \in [0, \ell]} \|F(w, v, t)\|_p + \alpha \|w\|_p.$$

For any  $y, z$  such that  $\|y - x(0)\|_p \leq \bar{K}$  and  $\|z - x(0)\|_p \leq \bar{K}$ , we apply the assumption on  $F$  and get that

$$\begin{aligned} &\|\bar{F}(y, t) - \bar{F}(z, t)\|_p \\ &\leq \alpha^2 \|F(\alpha^{-1}y_{(1)}, y_{(2)}, t) - F(\alpha^{-1}z_{(1)}, z_{(2)}, t)\|_p + \|y_{(1)} - z_{(1)}\|_p \\ &\leq \alpha^2 \left( \frac{1}{\alpha} \|\alpha^{-1}y_{(1)} - \alpha^{-1}z_{(1)}\|_p + \frac{1}{\alpha^2} \|y_{(2)} - z_{(2)}\|_p \right) + \|y_{(1)} - z_{(1)}\|_p \\ &\leq 3\|y - z\|_p. \end{aligned}$$

Also, by our assumption on  $x$ , we have a polynomial  $\bar{q} = (\alpha q'(\alpha t), q(\alpha t))$  such that

$$\left\| \frac{d}{dt}x(t) - \frac{d}{dt}\bar{q}(t) \right\|_p \leq \frac{28\varepsilon}{\bar{\ell}}$$

where  $x$  is the solution of the ODE (14.37). Therefore, Theorem 14.6.4 shows that we can compute  $x(\bar{\ell})$  with  $O(\varepsilon)$  error in  $O(nd \log^2(dK/\varepsilon))$  time plus  $O(d \log(K/\varepsilon))$  evaluations of  $F$ . Now, we can read off  $u(\ell)$  from  $x_{(2)}$  because

$$x_{(2)}(\bar{\ell}) = u(\alpha\bar{\ell}) = u(\ell).$$

The result for  $u'(\ell)$  is similar. □

## ■ 14.7 Derivative Estimations

For any smooth one dimension function  $f$ , we know by Taylor's theorem that

$$f(x) = \sum_{k=0}^N \frac{f^{(k)}(a)}{k!} (x-a)^k + \frac{1}{N!} \int_a^x (x-t)^N f^{(N+1)}(t) dt.$$

This formula provides a polynomial estimate of  $f$  around  $a$ . To analyze the accuracy of this estimate, we need to bound  $|f^{(N+1)}(t)|$ . In one dimension, we could simply give explicit formulas for the derivatives of  $f$  and use it to estimate the remainder term. However, for functions in high dimension, it is usually too tedious. In this section, we describe some techniques for bounding the derivatives of higher dimensional functions.

The derivatives of one-variable functions can be bounded via Cauchy's estimate (Theorem (14.2.11)). In Subsection 14.7.1, we give calculus rules that reduces the problem of estimating derivatives of high-dimensional functions to derivatives of one-dimensional functions. In Subsection 14.7.2, we show how to reduce bounding the derivative for an arbitrary ODE to an ODE in one dimension.

### ■ 14.7.1 Explicit Function

In this subsection, we show how to bound the derivatives of a complicated explicit function using the following object, generating upper bound. We reduce estimates of the derivatives of functions in high dimension to one variable rational functions. Since rational functions are holomorphic, one can apply Cauchy's estimates (Theorem (14.2.11)) to bound their derivatives.

**Definition 14.7.1.** Given a function  $F$ . We call that  $F \leq_x f$  for some one variable function  $f : \mathbb{R} \rightarrow \mathbb{R}$  if

$$\|D^{(k)}F(x)[\Delta_1, \Delta_2, \dots, \Delta_k]\| \leq f^{(k)}(0) \prod_{i=1}^k \|\Delta_i\|^k \quad (14.38)$$

for any  $k \geq 0$  and any  $\Delta_i$ .

*Remark.* In general,  $F \leq_x f$  and  $f(t) \leq g(t)$  point-wise does NOT imply  $F \leq_x g$ . However,  $F \leq_x f$  and  $f \leq_0 g$  does imply  $F \leq_x g$ .

This concept is useful for us to reduce bounding derivatives of a high dimension function to bounding derivatives of 1 dimension function. First of all, we note that upper bounds are composable.

**Lemma 14.7.2.** Given  $F \leq_x f$  and  $G \leq_{F(x)} g$ , we have that

$$G \circ F \leq_x g \circ \bar{f}$$

where  $\bar{f}(s) = f(s) - f(0)$ .

*Proof.* Fix any  $\Delta_1, \Delta_2, \dots$  be unit vectors in the domain of  $F$ . Let  $H(x) = G \circ F(x)$ . By chain rule, we have that

$$\begin{aligned} DH(x)[\Delta_1] &= DG(F(x))[DF(x)[\Delta_1]], \\ DH(x)[\Delta_1, \Delta_2] &= DG(F(x))[D^2F(x)[\Delta_1, \Delta_2] \\ &\quad + D^2G(F(x))[DF(x)[\Delta_1], DF(x)[\Delta_2]], \\ DH(x)[\Delta_1, \Delta_2, \Delta_3] &= DG(F(x))[D^2F(x)[\Delta_1, \Delta_2, \Delta_3] \\ &\quad + D^2G(F(x))[D^2F(x)[\Delta_1, \Delta_2], DF(x)[\Delta_3]] \\ &\quad + D^2G(F(x))[D^2F(x)[\Delta_1, \Delta_3], DF(x)[\Delta_2]] \\ &\quad + D^2G(F(x))[D^2F(x)[\Delta_2, \Delta_3], DF(x)[\Delta_1]] \\ &\quad + D^3G(F(x))[DF(x)[\Delta_1], DF(x)[\Delta_2], DF(x)[\Delta_3]], \\ &\quad \vdots \end{aligned}$$

Since  $G \leq_{F(x)} g$ , equation (14.38) shows that

$$\begin{aligned} \|DH(x)[\Delta_1]\| &\leq g^{(1)}(0)\|DF(x)[\Delta_1]\|, \\ \|D^2H(x)[\Delta_1, \Delta_2]\| &\leq g^{(1)}(0)\|D^2F(x)[\Delta_1, \Delta_2]\| \\ &\quad + g^{(2)}(0)\|DF(x)[\Delta_1]\|_2\|DF(x)[\Delta_2]\|, \\ \|D^3H(x)[\Delta_1, \Delta_2, \Delta_3]\| &\leq g^{(1)}(0)\|D^2F(x)[\Delta_1, \Delta_2, \Delta_3]\| \\ &\quad + g^{(2)}(0)\|D^2F(x)[\Delta_1, \Delta_2]\|\|DF(x)[\Delta_3]\| \\ &\quad + g^{(2)}(0)\|D^2F(x)[\Delta_1, \Delta_3]\|\|DF(x)[\Delta_2]\| \\ &\quad + g^{(2)}(0)\|D^2F(x)[\Delta_2, \Delta_3]\|\|DF(x)[\Delta_1]\| \\ &\quad \vdots \\ &\quad + g^{(3)}(0)\|DF(x)[\Delta_1]\|\|DF(x)[\Delta_2]\|\|DF(x)[\Delta_3]\|. \end{aligned}$$

Now, we use  $F \leq_x f$  to get

$$\begin{aligned} \|DH(x)[\Delta_1]\| &\leq g^{(1)}(0)f^{(1)}(0) = (g \circ \bar{f})^{(1)}(0) \\ \|D^2H(x)[\Delta_1, \Delta_2]\| &\leq g^{(1)}(0)f^{(2)}(0) + g^{(2)}(0)f^{(1)}(0)^2 = (g \circ \bar{f})^{(2)}(0), \\ \|D^3H(x)[\Delta_1, \Delta_2, \Delta_3]\| &\leq g^{(1)}(0)f^{(3)}(0) + 2g^{(2)}(0)f^{(2)}(0)f^{(1)}(0) \\ &\quad \vdots + g^{(3)}(0)f^{(1)}(0)^3 = (g \circ \bar{f})^{(3)}(0). \end{aligned}$$

Therefore, we have that  $\|D^kH(x)[\Delta_i]\| \leq (g \circ \bar{f})^{(k)}(0)$  for all  $k \geq 1$ . For  $k = 0$ , we have that  $\|H(x)\| = \|G(F(x))\| \leq g(0) = g(\bar{f}(0))$ .  $\square$

Next, we give some extra calculus rule for generating upper bounds.

**Lemma 14.7.3.** *Given that  $H_i \leq_x h_i$  for all  $i = 1, \dots, k$ . Then, we have that*

$$\sum_{i=1}^k H_i \leq_x \sum_{i=1}^k h_i \text{ and } \prod_{i=1}^k H_i \leq_x \prod_{i=1}^k h_i.$$

Given that  $H \leq_x h$  and  $\|H^{-1}(x)\| \leq C$ , we have that

$$H^{-1} \leq_x \frac{1}{C^{-1} - (h(s) - h(0))}.$$

*Proof.* Fix  $\Delta_1, \Delta_2, \dots$  be unit vectors. For the first claim, let  $H = \sum H_i$ , we note that

$$D^j H[\Delta_1, \dots, \Delta_j] = \sum_{i=1}^k D^j H_i[\Delta_1, \dots, \Delta_j].$$

Therefore, we have that  $\|D^j H[\Delta_1, \dots, \Delta_j]\| \leq \sum_{i=1}^k \|D^j H_i[\Delta_1, \dots, \Delta_j]\|$ . Since  $H \leq_x h_i$ , we have that  $\|D^j H[\Delta_1, \dots, \Delta_j]\| \leq \sum_{i=1}^k h^{(j)}(0)$ . Hence, we have  $H \leq_x \sum_{i=1}^k h_i$ .

For the second claim, we note that

$$D^j G = \sum_{i_1+i_2+\dots+i_k=j} \prod_{l=1}^k D^{i_l} H_l.$$

Let  $g = \prod_{i=1}^k h_i$ . Then, we have that

$$\|D^j G\| \leq \sum_{i_1+i_2+\dots+i_k=j} \prod_{l=1}^k \|D^{i_l} H_l\| \leq \sum_{i_1+i_2+\dots+i_k=j} \prod_{l=1}^k h_i^{(i_l)}(0) = D^j g(0).$$

Hence,  $G \leq_x g$ .

For the last claim, we first consider the function  $\Phi(M) = M^{-1}$ . Note that  $D\Phi[\Delta_1] = -M^{-1}\Delta_1M^{-1}$ ,  $D^2\Phi[\Delta] = M^{-1}\Delta_1M^{-1}\Delta_2M^{-1} + M^{-1}\Delta_2M^{-1}\Delta_1M^{-1}$  and hence

$$\|D^j \Phi\| \leq j! \|M^{-1}\|^{j+1} = j! C^{j+1}$$

Hence, we have  $\Phi \leq_M \sum_{j=0}^{\infty} \frac{j! C^{j+1}}{j!} s^j = \frac{1}{C^{-1}-s}$ . By Lemma 14.7.2, we see that  $H^{-1} \leq_x \frac{1}{C^{-1}-(h(s)-h(0))}$ .  $\square$

## ■ 14.7.2 Explicit ODE

In this section, we study the Taylor expansion of the solution of ODE (14.31).

**Lemma 14.7.4.** *Let  $u(t)$  be the solution of the ODE  $u'(t) = F(u(t))$ . Suppose that  $F \leq_{u(0)} f$  and let  $\psi(t)$  be the solution of the ODE  $\psi'(t) = f(\psi(t))$  and  $\psi(0) = 0$ . Then, we have*

$$\|u^{(k)}(0)\|_2 \leq \psi^{(k)}(0)$$

for all  $k \geq 1$ .

*Proof.* Since  $u'(t) = F(u(t))$ , we have that

$$\begin{aligned} u^{(2)}(t) &= DF(u(t))[u^{(1)}(t)], \\ u^{(3)}(t) &= DF(u(t))[u^{(2)}(t)] + D^2F(u(t))[u^{(1)}(t), u^{(1)}(t)], \\ u^{(4)}(t) &= DF(u(t))[u^{(3)}(t)] + 2D^2F(u(t))[u^{(2)}(t), u^{(1)}(t)] \\ &\quad + D^3F(u(t))[u^{(1)}(t), u^{(1)}(t), u^{(1)}(t)] \\ &\quad \vdots \end{aligned}$$

Therefore, we have

$$\begin{aligned} \|u^{(2)}(0)\|_2 &\leq f^{(1)}(0)\|u^{(1)}(0)\|_2, \\ \|u^{(3)}(0)\|_2 &\leq f^{(1)}(0)\|u^{(2)}(0)\|_2 + f^{(2)}(0)\|u^{(1)}(0)\|_2^2, \\ \|u^{(4)}(0)\|_2 &\leq f^{(1)}(0)\|u^{(3)}(0)\|_2 + 2f^{(2)}(0)\|u^{(2)}(0)\|_2\|u^{(1)}(0)\|_2 + f^{(3)}(0)\|u^{(1)}(0)\|_2^3, \\ &\quad \vdots \end{aligned}$$

By expanding  $\psi'(t) = f(\psi(t))$  at  $t = 0$ , we see that

$$\begin{aligned} \psi^{(2)}(0) &= f^{(1)}(0)\psi^{(1)}(0), \\ \psi^{(3)}(0) &= f^{(1)}(0)\psi^{(2)}(0) + f^{(2)}(0)\left(\psi^{(1)}(0)\right)^2, \\ \psi^{(4)}(0) &= f^{(1)}(0)\psi^{(3)}(0) + 2f^{(2)}(0)\psi^{(2)}(0)\psi^{(1)}(0) + f^{(3)}(0)\left(\psi^{(1)}(0)\right)^3, \\ &\quad \vdots \end{aligned}$$

Since  $\|u^{(1)}(0)\|_2 = \|F(u(0))\|_2 \leq f(0) = \psi^{(1)}(0)$ , we have that  $\|u^{(k)}(0)\|_2 \leq \psi^{(k)}(0)$  for all  $k \geq 1$ .  $\square$

Now, we apply Lemma 14.7.4 to second order ODE.

**Lemma 14.7.5.** *Let  $u(t)$  be the solution of the ODE  $u''(t) = F(u'(t), u(t), t)$ . Given some  $\alpha > 0$  and define  $\bar{F}(y, x, t) = \alpha^2 F(\alpha^{-1}y(t) + u'(0), x(t) + \alpha u'(0), \alpha t)$ . Suppose that  $\bar{F} \leq_{(0, u(0), 0)} f$  and let  $\psi(t)$  be the solution of the ODE  $\psi'(t) = 1 + \psi(t) + f(\psi(t))$  and  $\psi(0) = 0$ . Then, we have*

$$\|u^{(k)}(0)\|_2 \leq \frac{\psi^{(k)}(0)}{\alpha^k}$$

for all  $k \geq 2$ .

*Proof.* Let  $x(t) = u(\alpha t) - \alpha u'(0)$  and  $y(t) = \alpha u'(\alpha t) - \alpha u'(0)$ . Then, we can write the problem into



first order ODE

$$\begin{aligned}y'(t) &= \alpha^2 F(\alpha^{-1}y(t) + u'(0), x(t) + \alpha t u'(0), \alpha t) = \bar{F}(y(t), x(t), t), \\x'(t) &= y(t), \\t' &= 1.\end{aligned}$$

Let  $\mathbf{F}(y, x, t) = (\bar{F}(y, x, t), y, 1)$ . Then, we have that  $\|D^k \mathbf{F}\| \leq \|D^k 1\| + \|D^k y\| + \|D^k \bar{F}\|$ . Using  $\bar{F} \leq_{(0, u(0), 0)} f$ , we have that

$$\mathbf{F} \leq_{(0, x(0), 0)} 1 + t + f.$$

By Lemma 14.7.4, we know that

$$\|u^{(k)}(0)\|_2 = \frac{\|x^{(k)}(0)\|_2}{\alpha^k} \leq \frac{\psi^{(k)}(0)}{\alpha^k}.$$

□



---

---

## Bibliography

- [1] Jacob Abernethy and Elad Hazan. “Faster Convex Optimization: Simulated Annealing with an Efficient Universal Barrier”. In: *arXiv preprint arXiv:1507.02528* (2015).
- [2] Dimitris Achlioptas. “Database-friendly random projections: Johnson-Lindenstrauss with binary coins”. In: *Journal of computer and System Sciences* 66.4 (2003), pp. 671–687.
- [3] Rudolf Ahlswede and Andreas Winter. “Strong converse for identification via quantum channels”. In: *Information Theory, IEEE Transactions on* 48.3 (2002), pp. 569–579.
- [4] R. K. Ahuja, T. L. Magnanti, and J. B. Orlin. *Network flows*. New York, NY, USA: Elsevier North-Holland, Inc., 1989.
- [5] Ravindra K. Ahuja et al. “Applications of Network Optimization”. In: *Network Models*. Vol. 7. Handbooks in Operations Research and Management Science. North-Holland, 1995, pp. 1–75.
- [6] Martin Aigner and Thomas A Dowling. “Matching theory for combinatorial geometries”. In: *Transactions of the American Mathematical Society* 158.1 (1971), pp. 231–245.
- [7] Boris Alexeev et al. “Phase retrieval with polarization”. In: *SIAM Journal on Imaging Sciences* 7.1 (2014), pp. 35–66.
- [8] Zeyuan Allen-Zhu, Yin Tat Lee, and Lorenzo Orecchia. “Using Optimization to Obtain a Width-Independent, Parallel, Simpler, and Faster Positive SDP Solver”. In: *arXiv preprint arXiv:1507.02259* (2015).
- [9] Zeyuan Allen-Zhu, Zhenyu Liao, and Lorenzo Orecchia. “Spectral Sparsification and Regret Minimization Beyond Multiplicative Updates”. In: *Proceedings of the 47th Annual ACM Symposium on Theory of Computing (STOC)*. 2015, pp. 237–245.
- [10] Kurt M. Anstreicher. “Large step volumetric potential reduction algorithms for linear programming”. In: *Annals of Operations Research* 62.1 (1996), pp. 521–538.
- [11] Kurt M Anstreicher. “On Vaidya’s volumetric cutting plane method for convex programming”. In: *Mathematics of Operations Research* 22.1 (1997), pp. 63–89.
- [12] Kurt M. Anstreicher. “The volumetric barrier for convex quadratic constraints”. In: *Mathematical Programming* 100.3 (2004), pp. 613–662.
- [13] Kurt M. Anstreicher. “The volumetric barrier for semidefinite programming”. In: *Mathematics of Operations Research* 25.3 (2000), pp. 365–380.
- [14] Kurt M. Anstreicher. “Towards a practical volumetric cutting plane method for convex programming”. In: *SIAM Journal on Optimization* 9.1 (1998), pp. 190–206.
- [15] Kurt M. Anstreicher. “Volumetric path following algorithms for linear programming”. In: *Math. Program.* 76 (1996), pp. 245–263.
- [16] Mica Arie-Nachimson et al. “Global motion estimation from point matches”. In: *3D Imaging, Modeling, Processing, Visualization and Transmission (3DIMPVT), 2012 Second International Conference on*. IEEE. 2012, pp. 81–88.

- [17] Sanjeev Arora, Elad Hazan, and Satyen Kale. “Fast algorithms for approximate semidefinite programming using the multiplicative weights update method”. In: *Foundations of Computer Science, 2005. FOCS 2005. 46th Annual IEEE Symposium on*. IEEE. 2005, pp. 339–348.
- [18] Sanjeev Arora, Elad Hazan, and Satyen Kale. “The Multiplicative Weights Update Method: A Meta-Algorithm and Applications”.
- [19] Sanjeev Arora and Satyen Kale. “A combinatorial, primal-dual approach to semidefinite programs”. In: *Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*. ACM. 2007, pp. 227–236.
- [20] David S Atkinson and Pravin M Vaidya. “A cutting plane algorithm for convex programming that uses analytic centers”. In: *Mathematical Programming* 69.1-3 (1995), pp. 1–43.
- [21] Y. Aumann and Y. Rabani. “An  $O(\log k)$  Approximate Min-Cut Max-Flow Theorem and Approximation Algorithm”. In: *SIAM Journal on Computing* 27.1 (1998), pp. 291–301.
- [22] Haim Avron, Petar Maymounkov, and Sivan Toledo. “Blendenpik: Supercharging LAPACK’s least-squares solver”. In: *SIAM J. Scientific Computing* 32.3 (2010), pp. 1217–1236.
- [23] Haim Avron and Sivan Toledo. “Effective stiffness: Generalizing effective resistance sampling to finite element matrices”. In: *arXiv preprint arXiv:1110.4437* (2011).
- [24] Francis Bach. “Learning with submodular functions: A convex optimization perspective”. In: *Foundations and Trends in Machine Learning* (2013).
- [25] Mihai Badoiu, Sariel Har-Peled, and Piotr Indyk. “Approximate clustering via core-sets”. In: *Proceedings on 34th Annual ACM Symposium on Theory of Computing, May 19-21, 2002, Montréal, Québec, Canada*. 2002, pp. 250–257.
- [26] Olivier Bahn et al. “A cutting plane method from analytic centers for stochastic programming”. In: *Mathematical Programming* 69.1-3 (1995), pp. 45–73.
- [27] Chanderjit Bajaj. “The algebraic degree of geometric optimization problems”. English. In: *Discrete & Computational Geometry* 3.2 (1988), pp. 177–191.
- [28] Egon Balas and Chang-Sung Yu. “A Note on the Weiszfeld-Kuhn Algorithm for the General Fermat Problem”. In: *Managme Sci Res Report* 484 (1982), pp. 1–6.
- [29] Afonso S Bandeira, Amit Singer, and Daniel A Spielman. “A Cheeger inequality for the graph connection Laplacian”. In: *SIAM Journal on Matrix Analysis and Applications* 34.4 (2013), pp. 1611–1630.
- [30] Francisco Barahona and William H Cunningham. “A submodular network simplex method”. In: *Mathematical Programming at Oberwolfach II*. Springer, 1984, pp. 9–31.
- [31] Joshua Batson, Daniel Spielman, and Nikhil Srivastava. “Twice-Ramanujan sparsifiers”. In: *SIAM Journal on Computing* 41.6 (2012), pp. 1704–1721.
- [32] Stephen R Becker, Emmanuel J Candès, and Michael C Grant. “Templates for convex cone problems with applications to sparse signal recovery”. In: *Mathematical programming computation* 3.3 (2011), pp. 165–218.
- [33] András A Benczúr and David R Karger. “Approximating st minimum cuts in  $\tilde{O}(n^2)$  time”. In: *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*. ACM. 1996, pp. 47–55.
- [34] M. Bern et al. “Support-graph preconditioners”. In: *SIAM J. Matrix Anal. & Appl* 27.4 (2006), pp. 930–951.

- [35] Dimitri Bertsekas. “Nonlinear programming”. In: (1999).
- [36] Dimitris Bertsimas and Santosh Vempala. “Solving convex programs by random walks”. In: *Journal of the ACM (JACM)* 51.4 (2004), pp. 540–556.
- [37] Prosenjit Bose, Anil Maheshwari, and Pat Morin. “Fast approximations for sums of distances, clustering and the Fermat-Weber problem”. In: *Computational Geometry* 24.3 (2003), pp. 135–146.
- [38] Jean Bourgain, Joram Lindenstrauss, and V Milman. “Approximation of zonoids by zonotopes”. In: *Acta mathematica* 162.1 (1989), pp. 73–141.
- [39] Carl Brezovec, Gerard Cornuéjols, and Fred Glover. “Two algorithms for weighted matroid intersection”. In: *Mathematical Programming* 36.1 (1986), pp. 39–53.
- [40] Sébastien Bubeck. “Convex Optimization: Algorithms and Complexity”. In: *Foundations and Trends in Machine Learning* 8.3-4 (2015), pp. 231–357.
- [41] Sébastien Bubeck. “Theory of Convex Optimization for Machine Learning”. In: *Arxiv preprint arXiv:1405.4980* (2014).
- [42] Sébastien Bubeck and Ronen Eldan. “The entropic barrier: a simple and optimal universal self-concordant barrier”. In: *arXiv preprint arXiv:1412.1587* (2014).
- [43] Sébastien Bubeck, Ronen Eldan, and Joseph Lehec. “Sampling from a log-concave distribution with Projected Langevin Monte Carlo”. In: *arXiv preprint arXiv:1507.02564* (2015).
- [44] Yang Cai, Constantinos Daskalakis, and S Matthew Weinberg. “Optimal multi-dimensional mechanism design: Reducing revenue to welfare maximization”. In: *Foundations of Computer Science (FOCS), 2012 IEEE 53rd Annual Symposium on*. IEEE. 2012, pp. 130–139.
- [45] Yang Cai, Constantinos Daskalakis, and S Matthew Weinberg. “Reducing revenue to welfare maximization: Approximation algorithms and other generalizations”. In: *Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM. 2013, pp. 578–595.
- [46] Marcel K. de Carli Silva, Nicholas J. A. Harvey, and Cristiane M. Sato. “Sparse Sums of Positive Semidefinite Matrices”. In: *CoRR* abs/1107.0088 (2011).
- [47] R. Chandrasekaran and A. Tamir. “Open questions concerning Weiszfeld’s algorithm for the Fermat-Weber location problem”. English. In: *Mathematical Programming* 44.1-3 (1989), pp. 293–295.
- [48] Chih-Chung Chang and Chih-Jen Lin. “LIBSVM: A library for support vector machines”. In: *ACM Transactions on Intelligent Systems and Technology (TIST)*. 2.3 (2011). Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>, p. 27.
- [49] L. Paul Chew. “There Are Planar Graphs Almost as Good as the Complete Graph”. In: *Journal of Computer and System Sciences* 39.2 (1989), pp. 205–219.
- [50] Hui Han Chin et al. “Runtime guarantees for regression problems”. In: *ITCS*. 2013, pp. 269–282.
- [51] Paul Christiano et al. “Electrical flows, laplacian systems, and faster approximation of maximum flow in undirected graphs”. In: *Proceedings of the 43rd annual ACM symposium on Theory of computing*. ACM. 2011, pp. 273–282.
- [52] Nam-Kee Chung and Dong-Wan Tcha. “A dual algorithm for submodular flow problems”. In: *Operations research letters* 10.8 (1991), pp. 489–495.

- [53] Kenneth Clarkson and David Woodruff. “Low rank approximation and regression in input sparsity time”. In: *Proceedings of the 45th annual ACM symposium on Symposium on theory of computing*. ACM. 2013, pp. 81–90.
- [54] Kenneth Clarkson and David Woodruff. “Numerical linear algebra in the streaming model”. In: *Proceedings of the 41st Annual ACM Symposium on Theory of Computing (STOC)*. 2009, pp. 205–214.
- [55] E. S. Coakley, Vladimir Rokhlin, and Mark Tygert. “A Fast Randomized Algorithm for Orthogonal Projection”. In: *SIAM J. Scientific Computing* 33.2 (2011), pp. 849–868.
- [56] Michael B Cohen and Richard Peng. “Lp Row Sampling by Lewis Weights”. In: *arXiv preprint arXiv:1412.0588* (2014).
- [57] Michael B. Cohen et al. “Solving SDD Linear Systems in Nearly  $M\log^{1/2}N$  Time”. In: *Proceedings of the 46th Annual ACM Symposium on Theory of Computing*. STOC ’14. New York, New York: ACM, 2014, pp. 343–352.
- [58] Leon Cooper and I.Norman Katz. “The Weber problem revisited”. In: *Computers and Mathematics with Applications* 7.3 (1981), pp. 225–234.
- [59] B. Cousins and S. Vempala. “A Cubic Algorithm for Computing Gaussian Volume”. In: *SODA*. 2014, pp. 1215–1228.
- [60] William H Cunningham. “Improved bounds for matroid partition and intersection algorithms”. In: *SIAM Journal on Computing* 15.4 (1986), pp. 948–957.
- [61] William H Cunningham. “On submodular function minimization”. In: *Combinatorica* 5.3 (1985), pp. 185–192.
- [62] William H Cunningham and András Frank. “A primal-dual algorithm for submodular flows”. In: *Mathematics of Operations Research* 10.2 (1985), pp. 251–262.
- [63] Samuel I Daitch and Daniel A Spielman. “Faster approximate lossy generalized flow via interior point algorithms”. In: *Proceedings of the 40th annual ACM symposium on Theory of computing*. ACM. 2008, pp. 451–460.
- [64] Arnak S Dalalyan. “Theoretical guarantees for approximate sampling from smooth and log-concave densities”. In: *arXiv preprint arXiv:1412.7392* (2014).
- [65] Constantinos Daskalakis and S Matthew Weinberg. “Bayesian truthful mechanisms for job scheduling from bi-criterion approximation algorithms”. In: *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM. 2015, pp. 1934–1952.
- [66] James Demmel et al. “Fast matrix multiplication is stable”. In: *Numerische Mathematik* 106.2 (2007), pp. 199–224.
- [67] A. B. Dieker. “Reflected Brownian Motion”. In: *Wiley Encyclopedia of Operations Research and Management Science*. John Wiley & Sons, Inc., 2010.
- [68] EA Dinic. “An algorithm for the solution of the max-flow problem with the polynomial estimation”. In: *Doklady Akademii Nauk SSSR* 194.4 (1970), pp. 1277–1280.
- [69] Zvi Drezner et al. “Facility location”. In: Springer, 2002. Chap. The Weber problem, pp. 1–36.
- [70] Petros Drineas, Ravi Kannan, and Michael W. Mahoney. “Fast Monte Carlo algorithms for matrices I: Approximating matrix multiplication”. In: *SIAM J. Comput.* 36.1 (2006). Proceedings of the 42nd Annual IEEE Symposium on Foundations of Computer Science (FOCS), pp. 132–157.

- [71] Petros Drineas, Ravi Kannan, and Michael W. Mahoney. “Fast Monte Carlo Algorithms for Matrices II: Computing a Low-Rank Approximation to a Matrix”. In: *SIAM J. Comput.* 36.1 (2006), pp. 158–183.
- [72] Petros Drineas, Ravi Kannan, and Michael W. Mahoney. “Fast Monte Carlo Algorithms for Matrices III: Computing a Compressed Approximate Matrix Decomposition”. In: *SIAM J. Comput.* 36.1 (2006). Proceedings of the 14th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA), pp. 184–206.
- [73] Petros Drineas and Michael W Mahoney. “Effective resistances, statistical leverage, and applications to linear equation solving”. In: *arXiv preprint arXiv:1005.3097* (2010).
- [74] Petros Drineas, Michael W. Mahoney, and S. Muthukrishnan. “Relative-Error CUR Matrix Decompositions”. In: *SIAM J. Matrix Anal. Appl.* 30.2 (2008), pp. 844–881.
- [75] Petros Drineas, Michael W. Mahoney, and S. Muthukrishnan. “Sampling Algorithms for  $\ell_2$  Regression and Applications”. In: *Proceedings of the 17th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 2006, pp. 1127–1136.
- [76] Petros Drineas, Michael W. Mahoney, and S. Muthukrishnan. “Subspace Sampling and Relative-Error Matrix Approximation: Column-Based Methods”. In: *Proceedings of the 10th International Workshop on Randomization and Computation (RANDOM)*. 2006, pp. 316–326.
- [77] A Dutt, M Gu, and V Rokhlin. “Fast algorithms for polynomial interpolation, integration, and differentiation”. In: *SIAM Journal on Numerical Analysis* 33.5 (1996), pp. 1689–1711.
- [78] M. E. Dyer and A. M. Frieze. “Computing the volume of a convex body: a case where randomness provably helps”. In: *Proc. of AMS Symposium on Probabilistic Combinatorics and Its Applications*. 1991, pp. 123–170.
- [79] M. E. Dyer, A. M. Frieze, and R. Kannan. “A random polynomial-time algorithm for approximating the volume of convex bodies”. In: *J. ACM* 38.1 (1991), pp. 1–17.
- [80] Jack Edmonds. “Matroid intersection”. In: *Annals of discrete Mathematics* 4 (1979), pp. 39–49.
- [81] Jack Edmonds. “Matroid partition”. In: *Mathematics of the Decision Sciences* 11 (1968), pp. 335–345.
- [82] Jack Edmonds. “Submodular functions, matroids, and certain polyhedra”. In: *Edited by G. Goos, J. Hartmanis, and J. van Leeuwen* (1970), p. 11.
- [83] Jack Edmonds and Rick Giles. “A min-max relation for submodular functions on graphs”. In: *Studies in Integer Programming (PL Hammer, EL Johnson and BH Korte, eds.)*, *Ann. Discrete Math* 1 (1977), pp. 185–204.
- [84] Jack Edmonds and Richard M Karp. “Theoretical improvements in algorithmic efficiency for network flow problems”. In: *Journal of the ACM (JACM)* 19.2 (1972), pp. 248–264.
- [85] Shimon Even and R. Endre Tarjan. “Network Flow and Testing Graph Connectivity”. In: *SIAM Journal on Computing* 4.4 (Dec. 1975), pp. 507–518.
- [86] Dan Feldman and Michael Langberg. “A unified framework for approximating and clustering data”. In: *Proceedings of the forty-third annual ACM symposium on Theory of computing*. ACM. 2011, pp. 569–578.
- [87] Lisa Fleischer and Satoru Iwata. “A push-relabel framework for submodular function minimization and applications to parametric optimization”. In: *Discrete Applied Mathematics* 131.2 (2003), pp. 311–322.

- [88] Lisa Fleischer and Satoru Iwata. “Improved algorithms for submodular function minimization and submodular flow”. In: *Proceedings of the thirty-second annual ACM symposium on Theory of computing*. ACM. 2000, pp. 107–116.
- [89] Lisa Fleischer, Satoru Iwata, and S Thomas McCormick. “A faster capacity scaling algorithm for minimum cost submodular flow”. In: *Mathematical Programming* 92.1 (2002), pp. 119–139.
- [90] András Frank. “A weighted matroid intersection algorithm”. In: *Journal of Algorithms* 2.4 (1981), pp. 328–336.
- [91] András Frank and Éva Tardos. “An application of simultaneous Diophantine approximation in combinatorial optimization”. In: *Combinatorica* 7.1 (1987), pp. 49–65.
- [92] S. Fujishige. “Algorithms for solving the independent-flow problems”. In: *Journal of the Operations Research Society of Japan* (1978).
- [93] Satoru Fujishige. “An out-of-kilter method for submodular flows”. In: *Discrete applied mathematics* 17.1 (1987), pp. 3–16.
- [94] Satoru Fujishige and Satoru Iwata. “Algorithms for submodular flows”. In: *IEICE TRANSACTIONS on Information and Systems* 83.3 (2000), pp. 322–329.
- [95] Satoru Fujishige, Hans Röck, and Uwe Zimmermann. “A strongly polynomial algorithm for minimum cost submodular flow problems”. In: *Mathematics of Operations Research* 14.1 (1989), pp. 60–69.
- [96] Satoru Fujishige and Zhang Xiaodong. “An efficient cost scaling algorithm for the independent assignment problem”. In: *Journal of the Operations Research Society of Japan* 38.1 (1995), pp. 124–136.
- [97] Mituhiro Fukuda et al. “Exploiting sparsity in semidefinite programming via matrix completion I: General framework”. In: *SIAM Journal on Optimization* 11.3 (2001), pp. 647–674.
- [98] Zvi Galil and Éva Tardos. “An  $O(n^2(m+n \log n) \log n)$  min-cost flow algorithm”. In: *Journal of the ACM (JACM)* 35.2 (1988), pp. 374–386.
- [99] François Le Gall. “Powers of tensors and fast matrix multiplication”. In: *arXiv preprint arXiv:1401.7714* (2014).
- [100] François Le Gall. “Powers of tensors and fast matrix multiplication”. In: *arXiv:1401.7714* (2014).
- [101] Jean-Louis Goffin, Zhi-Quan Luo, and Yinyu Ye. “Complexity analysis of an interior cutting plane method for convex feasibility problems”. In: *SIAM Journal on Optimization* 6.3 (1996), pp. 638–652.
- [102] Jean-Louis Goffin and Jean-Philippe Vial. “Convex nondifferentiable optimization: A survey focused on the analytic center cutting plane method”. In: *Optimization Methods and Software* 17.5 (2002), pp. 805–867.
- [103] Jean-Louis Goffin and Jean-Philippe Vial. “Shallow, deep and very deep cuts in the analytic center cutting plane method”. In: *Mathematical Programming* 84.1 (1999), pp. 89–103.
- [104] Andrew V. Goldberg and Satish Rao. “Beyond the flow decomposition barrier”. In: *J. ACM* 45.5 (1998), pp. 783–797.
- [105] Andrew V Goldberg and Robert E Tarjan. “A new approach to the maximum-flow problem”. In: *Journal of the ACM (JACM)* 35.4 (1988), pp. 921–940.



- [106] Andrew V Goldberg and Robert E Tarjan. “Finding minimum-cost circulations by successive approximation”. In: *Mathematics of Operations Research* 15.3 (1990), pp. 430–466.
- [107] Clovis C Gonzaga. “Path-following methods for linear programming”. In: *SIAM review* 34.2 (1992), pp. 167–224.
- [108] Clóvis C Gonzaga and Elizabeth W Karas. “Fine tuning Nesterov’s steepest descent algorithm for differentiable convex programming”. In: *Mathematical Programming* 138.1-2 (2013), pp. 141–166.
- [109] Martin Grötschel, László Lovász, and Alexander Schrijver. “Geometric Algorithms and Combinatorial Optimization”. In: Springer, 1988.
- [110] Martin Grötschel, László Lovász, and Alexander Schrijver. “The ellipsoid method and its consequences in combinatorial optimization”. In: *Combinatorica* 1.2 (1981), pp. 169–197.
- [111] B. Grünbaum. “Partitions of mass-distributions and of convex bodies by hyperplanes”. In: *Pacific J. Math* 10.4 (1960), pp. 1257–1261.
- [112] N. Halko, P. G. Martinsson, and J. A. Tropp. “Finding Structure with Randomness: Probabilistic Algorithms for Constructing Approximate Matrix Decompositions”. In: *SIAM Review* 53.2 (2011), pp. 217–288.
- [113] Sariel Har-Peled and Akash Kushal. “Smaller coresets for k-median and k-means clustering”. In: *Proceedings of the twenty-first annual symposium on Computational geometry*. ACM. 2005, pp. 126–134.
- [114] J. Michael Harrison. *Brownian motion and stochastic flow systems*. Wiley series in probability and mathematical statistics. New York: Wiley, 1985.
- [115] Christoph Helmberg and Franz Rendl. “A spectral bundle method for semidefinite programming”. In: *SIAM Journal on Optimization* 10.3 (2000), pp. 673–696.
- [116] P. Indyk and Stanford University. Computer Science Dept. *High-dimensional computational geometry*. Stanford University, 2000.
- [117] Satoru Iwata. “A capacity scaling algorithm for convex cost submodular flows”. In: *Mathematical programming* 76.2 (1997), pp. 299–308.
- [118] Satoru Iwata. “A faster scaling algorithm for minimizing submodular functions”. In: *SIAM Journal on Computing* 32.4 (2003), pp. 833–840.
- [119] Satoru Iwata. “A fully combinatorial algorithm for submodular function minimization”. In: *Journal of Combinatorial Theory, Series B* 84.2 (2002), pp. 203–212.
- [120] Satoru Iwata. “Submodular function minimization”. In: *Mathematical Programming* 112.1 (2008), pp. 45–64.
- [121] Satoru Iwata, Lisa Fleischer, and Satoru Fujishige. “A combinatorial strongly polynomial algorithm for minimizing submodular functions”. In: *Journal of the ACM (JACM)* 48.4 (2001), pp. 761–777.
- [122] Satoru Iwata, S Thomas McCormick, and Maiko Shigeno. “A fast cost scaling algorithm for submodular flow”. In: *Information Processing Letters* 74.3 (2000), pp. 123–128.
- [123] Satoru Iwata, S Thomas McCormick, and Maiko Shigeno. “A Faster Algorithm for Minimum Cost Submodular Flows.” In: *SODA*. 1998, pp. 167–174.
- [124] Satoru Iwata, S Thomas McCormick, and Maiko Shigeno. “A strongly polynomial cut canceling algorithm for the submodular flow problem”. In: *Integer Programming and Combinatorial Optimization*. Springer, 1999, pp. 259–272.

- [125] Satoru Iwata and James B Orlin. “A simple combinatorial algorithm for submodular function minimization”. In: *Proceedings of the twentieth Annual ACM-SIAM Symposium on Discrete Algorithms*. Society for Industrial and Applied Mathematics. 2009, pp. 1230–1237.
- [126] Rahul Jain and Penghui Yao. “A parallel approximation algorithm for positive semidefinite programming”. In: *Foundations of Computer Science (FOCS), 2011 IEEE 52nd Annual Symposium on*. IEEE. 2011, pp. 463–471.
- [127] Klaus Jansen. “Approximate strong separation with application in fractional graph coloring and preemptive scheduling”. In: *Theoretical Computer Science* 302.1 (2003), pp. 239–256.
- [128] Fritz John. “Extremum problems with inequalities as side conditions”. In: *Studies and Essays presented to R. Courant on his 60th Birthday (1948)*.
- [129] R. Johnson and T. Zhang. “Accelerating Stochastic Gradient Descent using Predictive Variance Reduction”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2013.
- [130] William Johnson and Joram Lindenstrauss. “Extensions of Lipschitz mappings into a Hilbert space”. In: *Conference on modern analysis and probability*. Vol. 26. Contemporary Mathematics. 1984, pp. 189–206.
- [131] R. Kannan, L. Lovász, and M. Simonovits. “Random walks and an  $O^*(n^5)$  volume algorithm for convex bodies”. In: *Random Structures and Algorithms* 11 (1997), pp. 1–50.
- [132] R. Kannan and H. Narayanan. “Random walks on polytopes and an affine interior point method for linear programming”. In: *STOC*. 2009, pp. 561–570.
- [133] Ravindran Kannan and Hariharan Narayanan. “Random walks on polytopes and an affine interior point method for linear programming”. In: *Mathematics of Operations Research* 37.1 (2012), pp. 1–20.
- [134] David R Karger. “Better random sampling algorithms for flows in undirected graphs”. In: *Proceedings of the ninth annual ACM-SIAM symposium on Discrete algorithms*. Society for Industrial and Applied Mathematics. 1998, pp. 490–499.
- [135] David Karger and Matthew Levine. “Random sampling in residual graphs”. In: *Proceedings of the thirty-fourth annual ACM symposium on Theory of computing*. ACM. 2002, pp. 63–66.
- [136] Narendra Karmarkar. “A new polynomial-time algorithm for linear programming”. In: *Proceedings of the sixteenth annual ACM symposium on Theory of computing*. ACM. 1984, pp. 302–311.
- [137] Narendra Karmarkar. “Riemannian geometry underlying interior-point methods for linear programming”. In: *Contemporary Mathematics* 114 (1990), pp. 51–75.
- [138] Richard M Karp and Christos H Papadimitriou. “On linear characterizations of combinatorial optimization problems”. In: *SIAM Journal on Computing* 11.4 (1982), pp. 620–632.
- [139] Jonathan A. Kelner and Petar Maymounkov. “Electric routing and concurrent flow cutting”. In: *CoRR* abs/0909.2859 (2009).
- [140] Jonathan A. Kelner, Gary L. Miller, and Richard Peng. “Faster approximate multicommodity flow using quadratically coupled flows”. In: *Proceedings of the 44th symposium on Theory of Computing*. STOC '12. New York, New York, USA: ACM, 2012, pp. 1–18.
- [141] Jonathan A. Kelner et al. “A simple, combinatorial algorithm for solving SDD systems in nearly-linear time”. In: *Proceedings of the 45th Annual ACM Symposium on Theory of Computing (STOC)*. 2013, pp. 911–920.

- [142] Leonid G Khachiyan. “Polynomial algorithms in linear programming”. In: *USSR Computational Mathematics and Mathematical Physics* 20.1 (1980), pp. 53–72.
- [143] Leonid G Khachiyan. “Rounding of polytopes in the real number model of computation”. In: *Mathematics of Operations Research* 21.2 (1996), pp. 307–320.
- [144] LG Khachiyan, SP Tarasov, and II Erlikh. “The method of inscribed ellipsoids”. In: *Soviet Math. Dokl.* Vol. 37. 1. 1988, pp. 226–230.
- [145] Adam R Klivans and Daniel Spielman. “Randomness efficient identity testing of multivariate polynomials”. In: *Proceedings of the thirty-third annual ACM symposium on Theory of computing*. ACM. 2001, pp. 216–223.
- [146] Ioannis Koutis, Gary L. Miller, and Richard Peng. “A Nearly- $m \log n$  Time Solver for SDD Linear Systems”. In: *Foundations of Computer Science (FOCS), 2011 IEEE 52nd Annual Symposium on*. Oct. 2011, pp. 590–598.
- [147] Ioannis Koutis, Gary L. Miller, and Richard Peng. “Approaching Optimality for Solving SDD Linear Systems”. In: *Proceedings of the 51st Annual IEEE Symposium on Foundations of Computer Science (FOCS)*. 2010, pp. 235–244.
- [148] Jakob Krarup and Steven Vajda. “On Torricelli’s geometrical solution to a problem of Fermat”. In: *IMA Journal of Management Mathematics* 8.3 (1997), pp. 215–224.
- [149] Andreas Krause. <http://submodularity.org/>.
- [150] Dilip Krishnan, Raanan Fattal, and Richard Szeliski. “Efficient preconditioning of laplacian matrices for computer graphics”. In: *ACM Transactions on Graphics (TOG)* 32.4 (2013), p. 142.
- [151] Kartik Krishnan and John E Mitchell. “A unifying framework for several cutting plane methods for semidefinite programming”. In: *Optimization methods and software* 21.1 (2006), pp. 57–74.
- [152] Kartik Krishnan and John E Mitchell. “Properties of a cutting plane method for semidefinite programming”. In: *Pacific Journal of Optimization* 8.4 (2012), pp. 779–802.
- [153] Kartik Krishnan and Tamás Terlaky. “Interior point and semidefinite approaches in combinatorial optimization”. In: *Graph theory and combinatorial optimization*. Springer, 2005, pp. 101–157.
- [154] Richard A. Kronmal and Arthur V. Peterson. “The Alias and Alias-rejection-mixture Methods for Generating Random Variables from Probability Distributions”. In: *Proceedings of the 11th Conference on Winter Simulation - Volume 1*. WSC ’79. San Diego, California, USA: IEEE Press, 1979, pp. 269–280.
- [155] Harold W. Kuhn. “A note on Fermat’s problem”. English. In: *Mathematical Programming* 4.1 (1973), pp. 98–107.
- [156] Sanjiv Kumar, Mehryar Mohri, and Ameet Talwalkar. “Sampling Methods for the Nystrom Method”. In: *Journal of Machine Learning Research* 13.1 (2012), pp. 981–1006.
- [157] Eugene L Lawler. “Matroid intersection algorithms”. In: *Mathematical programming* 9.1 (1975), pp. 31–56.
- [158] Gregory Lawler and Hariharan Narayanan. “Mixing times and lp bounds for Oblivious Routing”. In: *Workshop on Analytic Algorithms and Combinatorics, (ANALCO 09)*. 2009.

- [159] N. Le Roux, M. Schmidt, and F. Bach. “A Stochastic Gradient Method with an Exponential Convergence Rate for Strongly-Convex Optimization with Finite Training Sets”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2012.
- [160] Yin Tat Lee, Richard Peng, and Daniel A Spielman. “Sparsified Cholesky Solvers for SDD linear systems”. In: *arXiv preprint arXiv:1506.08204* (2015).
- [161] Yin Tat Lee, Satish Rao, and Nikhil Srivastava. “A New Approach to Computing Maximum Flows using Electrical Flows”. In: *Proceedings of the 45th symposium on Theory of Computing - STOC '13* (2013).
- [162] Yin Tat Lee and Aaron Sidford. “Efficient Accelerated Coordinate Descent Methods and Faster Algorithms for Solving Linear Systems.” In: *The 54th Annual Symposium on Foundations of Computer Science (FOCS)*. 2013.
- [163] Yin Tat Lee and Aaron Sidford. “Efficient Accelerated Coordinate Descent Methods and Faster Algorithms for Solving Linear Systems”. In: *Proceedings of the 54th Annual IEEE Symposium on Foundations of Computer Science*. 2013.
- [164] Yin Tat Lee and Aaron Sidford. “Efficient Inverse Maintenance and Faster Algorithms for Linear Programming”. In: *arXiv:1503.01752* (2015).
- [165] Yin Tat Lee and Aaron Sidford. “Path Finding I: Solving Linear Programs with  $\tilde{O}(\sqrt{\text{rank}})$  Linear System Solves”. In: *arXiv preprint arXiv:1312.6677* (2013).
- [166] Yin Tat Lee and Aaron Sidford. “Path Finding II: An  $\tilde{O}(m \sqrt{n})$  Algorithm for the Minimum Cost Flow Problem”. In: *arXiv preprint arXiv:1312.6713* (2013).
- [167] F. Thomson Leighton and Ankur Moitra. “Extensions and limits to vertex sparsification”. In: *Proceedings of the 42nd ACM symposium on Theory of computing*. STOC '10. New York, NY, USA: ACM, 2010, pp. 47–56.
- [168] A. Yu Levin. “On an algorithm for the minimization of convex functions”. In: *Soviet Math. Doklady* (1965).
- [169] D Lewis. “Finite dimensional subspaces of  $L_p$ ”. In: *Studia Mathematica* 63.2 (1978), pp. 207–212.
- [170] Mu Li, Gary L Miller, and Richard Peng. “Iterative row sampling”. In: *Proceedings of the 54th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*. 2013, pp. 127–136.
- [171] E Lieb and W Thirring. “Inequalities for the moments of the eigenvalues of the Schrödinger equation and their relation to Sobolev inequalities”. In: *Studies in Mathematical Physics: Essays in honor of Valentine Bargman, Lieb, E., Simon, B., Wightman, AS (eds.)* (1976), pp. 269–303.
- [172] Hendrik P. Lopuhaa and Peter J. Rousseeuw. “Breakdown Points of Affine Equivariant Estimators of Multivariate Location and Covariance Matrices”. In: *Ann. Statist.* 19.1 (Mar. 1991), pp. 229–248.
- [173] L. Lovász. “Hit-and-Run mixes fast”. In: *Math. Prog* 86 (1998), pp. 443–461.
- [174] L. Lovász. “How to compute the volume?” In: *Jber. d. Dt. Math.-Verein, Jubiläumstagung 1990* (1990), pp. 138–151.
- [175] L. Lovász and M. Simonovits. “Mixing rate of Markov chains, an isoperimetric inequality, and computing the volume”. In: *ROCS*. 1990, pp. 482–491.
- [176] L. Lovász and M. Simonovits. “On the randomized complexity of volume and diameter”. In: *Proc. 33<sup>rd</sup> IEEE Annual Symp. on Found. of Comp. Sci.* 1992, pp. 482–491.

- [177] L. Lovász and M. Simonovits. “Random walks in a convex body and an improved volume algorithm”. In: *Random Structures and Alg.* Vol. 4. 1993, pp. 359–412.
- [178] L. Lovász and S. Vempala. “Fast Algorithms for Logconcave Functions: sampling, Rounding, Integration and Optimization”. In: *FOCS*. 2006, pp. 57–68.
- [179] L. Lovász and S. Vempala. “Hit-and-run from a corner”. In: *SIAM J. Computing* 35 (4 2006), pp. 985–1005.
- [180] László Lovász and Santosh Vempala. “Simulated annealing in convex bodies and an  $O^*(n^4)$  volume algorithm”. In: *J. Comput. Syst. Sci.* 72.2 (2006), pp. 392–417.
- [181] A. Lubotzky, R. Phillips, and P. Sarnak. “Ramanujan Graphs”. In: *Combinatorica* 8.3 (1988), pp. 261–277.
- [182] Aleksander Madry. “Fast approximation algorithms for cut-based problems in undirected graphs”. In: *Proceedings of the 51st Annual Symposium on Foundations of Computer Science*. 2010.
- [183] Aleksander Madry. “Navigating Central Path with Electrical Flows: from Flows to Matchings, and Back”. In: *Proceedings of the 54th Annual Symposium on Foundations of Computer Science*. 2013.
- [184] Michael W. Mahoney. “Randomized Algorithms for Matrices and Data”. In: *Foundations and Trends in Machine Learning* 3.2 (2011), pp. 123–224.
- [185] Michael W Mahoney and Xiangrui Meng. “Low-distortion subspace embeddings in input-sparsity time and applications to robust linear regression”. In: *Proceedings of the 45th Annual ACM Symposium on Theory of Computing (STOC)*. 2013, pp. 91–100.
- [186] Stefano Marchesini, Yu-Chao Tu, and Hau-tieng Wu. “Alternating projection, ptychographic imaging and phase synchronization”. In: *arXiv preprint arXiv:1402.0550* (2014).
- [187] Adam W Marcus, Nikhil Srivastava, and Daniel A Spielman. “Interlacing Families IV: Bipartite Ramanujan Graphs of All Sizes”. In: *arXiv preprint arXiv:1505.08010* (2015). to appear in FOCS 2015.
- [188] G. A. Margulis. “Explicit group theoretical constructions of combinatorial schemes and their application to the design of expanders and concentrators”. In: *Problems of Information Transmission* 24.1 (July 1988), pp. 39–46.
- [189] S McCormick. *Submodular Function Minimization*. 2013.
- [190] S Thomas and McCormick. “Canceling most helpful total submodular cuts for submodular flow.” In: *IPCO*. 1993, pp. 343–353.
- [191] J. A. Meijerink and H. A. van der Vorst. “An Iterative Solution Method for Linear Systems of Which the Coefficient Matrix is a Symmetric  $M$ -Matrix”. In: *Mathematics of Computation* 31.137 (1977), pp. 148–162.
- [192] Xiangrui Meng, Michael A. Saunders, and Michael W. Mahoney. “LSRN: A Parallel Iterative Solver for Strongly Over- or Under-Determined Systems”. In: *SIAM J. Scientific Computing* 36.2 (2014).
- [193] Carl D Meyer Jr. “Generalized inversion of modified matrices”. In: *SIAM J. Appl. Math* 24.3 (1973), pp. 315–323.
- [194] Gary L. Miller and Richard Peng. “Approximate Maximum Flow on Separable Undirected Graphs”. In: *Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms*. New Orleans, Louisiana, USA: SIAM, 2013, pp. 1151–1170.

- [195] Mehryar Mohri and Ameet Talwalkar. “Can matrix coherence be efficiently and accurately estimated?” In: *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2011, pp. 534–542.
- [196] Renato DC Monteiro. “First-and second-order methods for semidefinite programming”. In: *Mathematical Programming* 97.1-2 (2003), pp. 209–244.
- [197] Kazuhide Nakata et al. “Exploiting sparsity in semidefinite programming via matrix completion II: Implementation and numerical results”. In: *Mathematical Programming* 95.2 (2003), pp. 303–327.
- [198] Hariharan Narayanan. “Randomized interior point methods for sampling and optimization”. In: *Annals of Applied Probability* 26.1 (Feb. 2016), pp. 597–641.
- [199] Hariharan Narayanan and Alexander Rakhlin. “Random walk approach to regret minimization”. In: *Advances in Neural Information Processing Systems*. 2010, pp. 1777–1785.
- [200] Jelani Nelson and Huy L. Nguyen. “OSNAP: Faster Numerical Linear Algebra Algorithms via Sparser Subspace Embeddings”. In: *Proceedings of the 54th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*. 2013, pp. 117–126.
- [201] A. Nemirovski and D. Yudin. *Problem Complexity and Method Efficiency in Optimization*. Wiley Interscience, 1983.
- [202] Arkadi Nemirovski. “Efficient Methods in Convex Programming”. In: (1994).
- [203] Yu Nesterov. “A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ ”. In: *Soviet Mathematics Doklady* 27.2 (1983), pp. 372–376.
- [204] Yu Nesterov. “Complexity estimates of some cutting plane methods based on the analytic barrier”. In: *Mathematical Programming* 69.1-3 (1995), pp. 149–176.
- [205] Yu Nesterov. “Efficiency of coordinate descent methods on huge-scale optimization problems”. In: *Core discussion papers* 2 (2010), p. 2010.
- [206] Yu Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Vol. I. 2003.
- [207] Yu Nesterov. “Rounding of convex sets and efficient gradient methods for linear programming problems”. In: *Available at SSRN 965658* (2004).
- [208] Yu Nesterov. “Smooth minimization of non-smooth functions”. In: *Mathematical Programming* 103.1 (2005), pp. 127–152.
- [209] Yu E Nesterov and Michael J Todd. “Self-scaled barriers and interior-point methods for convex programming”. In: *Mathematics of Operations research* 22.1 (1997), pp. 1–42.
- [210] Yu Nesterov and Arkadi Nemirovski. “Conic formulation of a convex programming problem and duality”. In: *Optimization Methods and Software* 1.2 (1992), pp. 95–115.
- [211] Yu Nesterov and Arkadi Nemirovski. *Interior-Point Polynomial Algorithms in Convex Programming*. Studies in Applied Mathematics. SIAM, 1994.
- [212] Yu Nesterov and Arkadi Nemirovski. “Primal Central Paths and Riemannian Distances for Convex Sets”. In: *Foundations of Computational Mathematics* 8.5 (2008), pp. 533–560.
- [213] Yu Nesterov and Arkadi Nemirovski. *Self-concordant functions and polynomial-time methods in convex programming*. USSR Academy of Sciences, Central Economic & Mathematic Institute, 1989.

- [214] Yurii E Nesterov, Michael J Todd, et al. "On the Riemannian geometry defined by self-concordant barriers and interior-point methods". In: *Foundations of Computational Mathematics* 2.4 (2002), pp. 333–361.
- [215] James B Orlin. "A faster strongly polynomial minimum cost flow algorithm". In: *Operations research* 41.2 (1993), pp. 338–350.
- [216] James B Orlin. "A faster strongly polynomial time algorithm for submodular function minimization". In: *Mathematical Programming* 118.2 (2009), pp. 237–251.
- [217] James B Orlin. "Genuinely Polynomial Simplex and Non-Simplex Algorithms for the Minimum Cost Flow Problem". In: (1984).
- [218] James B Orlin. "Max flows in  $O(nm)$  time, or better". In: *Proceedings of the 45th annual ACM symposium on Symposium on theory of computing*. ACM, 2013, pp. 765–774.
- [219] James B Orlin, John VandeVate, et al. "On a " primal" matroid intersection algorithm". In: (1983).
- [220] Lawrence M. Ostresh. "On the Convergence of a Class of Iterative Methods for Solving the Weber Location Problem". In: *Operations Research* 26.4 (1978), pp. 597–609.
- [221] Onur Ozyesil, Amit Singer, and Ronen Basri. "Stable camera motion estimation using convex programming". In: *SIAM Journal on Imaging Sciences* 8.2 (2015), pp. 1220–1262.
- [222] Pablo A. Parrilo and Bernd Sturmfels. "Minimizing Polynomial Functions". In: *DIMACS Workshop on Algorithmic and Quantitative Aspects of Real Algebraic Geometry in Mathematics and Computer Science, March 12-16, 2001, DIMACS Center, Rutgers University, Piscataway, NJ, USA*. 2001, pp. 83–100.
- [223] Richard Peng and Daniel A Spielman. "An Efficient Parallel Solver for SDD Linear Systems". In: *arXiv preprint arXiv:1311.3286* (2013).
- [224] Frank Plastria and Mohamed Elosmani. "On the convergence of the Weiszfeld algorithm for continuous single facility location allocation problems". English. In: *TOP* 16.2 (2008), pp. 388–406.
- [225] Harald Räcke. "Optimal hierarchical decompositions for congestion minimization in networks". In: *Proceedings of the 40th annual ACM symposium on Theory of computing*. STOC '08. Victoria, British Columbia, Canada: ACM, 2008, pp. 255–264.
- [226] Srinivasan Ramaswamy and John E Mitchell. "A long step cutting plane algorithm that uses the volumetric barrier". In: *Department of Mathematical Science, RPI, Troy, NY* (1995).
- [227] James Renegar. "A polynomial-time algorithm, based on Newton's method, for linear programming". In: *Mathematical Programming* 40.1-3 (1988), pp. 59–93.
- [228] Vladimir Rokhlin and Mark Tygert. "A fast randomized algorithm for overdetermined linear least-squares regression". In: *Proc. Natl. Acad. Sci. USA* 105.36 (2008), pp. 13212–7.
- [229] Tamas Sarlos. "Improved approximation algorithms for large matrices via random projections". In: *Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*. 2006, pp. 143–152.
- [230] M Schmidt. "minFunc: unconstrained differentiable multivariate optimization in Matlab". In: *URL [http://www. di. ens. fr/mschmidt/Software/minFunc. html](http://www.di.ens.fr/mschmidt/Software/minFunc.html)* (2012).
- [231] Alexander Schrijver. "A combinatorial algorithm minimizing submodular functions in strongly polynomial time". In: *Journal of Combinatorial Theory, Series B* 80.2 (2000), pp. 346–355.

- [232] Alexander Schrijver. *Combinatorial optimization: polyhedra and efficiency*. Vol. 24. Springer, 2003.
- [233] Jonah Sherman. “Nearly Maximum Flows in Nearly Linear Time”. In: *Proceedings of the 54th Annual Symposium on Foundations of Computer Science*. 2013.
- [234] Maiko Shigeno and Satoru Iwata. “A dual approximation approach to weighted matroid intersection”. In: *Operations research letters* 18.3 (1995), pp. 153–156.
- [235] Yoel Shkolnisky and Amit Singer. “Viewing direction estimation in cryo-EM using synchronization”. In: *SIAM journal on imaging sciences* 5.3 (2012), pp. 1088–1110.
- [236] Naum Z Shor. “Cut-off method with space extension in convex programming problems”. In: *Cybernetics and systems analysis* 13.1 (1977), pp. 94–96.
- [237] Amit Singer and Yoel Shkolnisky. “Three-dimensional structure determination from common lines in cryo-EM by eigenvectors and semidefinite programming”. In: *SIAM journal on imaging sciences* 4.2 (2011), pp. 543–572.
- [238] Amit Singer and H-T Wu. “Vector diffusion maps and the connection Laplacian”. In: *Communications on pure and applied mathematics* 65.8 (2012).
- [239] Maurice Sion. “On general minimax theorems”. In: *Pacific J. Math* 8.1 (1958), pp. 171–176.
- [240] Daniel D. Sleator and Robert Endre Tarjan. “A data structure for dynamic trees”. In: *Proceedings of the thirteenth annual ACM symposium on Theory of computing*. STOC '81. Milwaukee, Wisconsin, USA: ACM, 1981, pp. 114–122.
- [241] R.L. Smith. “Efficient Monte-Carlo procedures for generating points uniformly distributed over bounded regions”. In: *Operations Res.* 32 (1984), pp. 1296–1308.
- [242] Daniel A. Spielman and Nikhil Srivastava. “Graph Sparsification by Effective Resistances”. In: *SIAM Journal on Computing* 40.6 (2011), pp. 1913–1926.
- [243] Daniel A. Spielman and Shang-Hua Teng. “Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems”. In: *Proceedings of the thirty-sixth annual ACM symposium on Theory of computing*. ACM. 2004, pp. 81–90.
- [244] Daniel A. Spielman and Shang-Hua Teng. “Spectral Sparsification of Graphs”. In: *CoRR* abs/0808.4134 (2008).
- [245] Michel Talagrand. “Embedding subspaces of  $L_1$  into  $L_1$ ”. In: *Proceedings of the American Mathematical Society* 108.2 (1990), pp. 363–369.
- [246] Michel Talagrand. “Embedding Subspaces of  $L_p$  in  $L_p$ ”. In: *Geometric Aspects of Functional Analysis*. Springer, 1995, pp. 311–326.
- [247] Éva Tardos. “A strongly polynomial minimum cost circulation algorithm”. In: *Combinatorica* 5.3 (1985), pp. 247–255.
- [248] Michael J Todd. “Semidefinite optimization”. In: *Acta Numerica 2001* 10 (2001), pp. 515–560.
- [249] Nobuaki Tomizawa and Masao Iri. “ALGORITHM FOR DETERMINING RANK OF A TRIPLE MATRIX PRODUCT  $AXB$  WITH APPLICATION TO PROBLEM OF DISCERNING EXISTENCE OF UNIQUE SOLUTION IN A NETWORK”. In: *ELECTRONICS & COMMUNICATIONS IN JAPAN* 57.11 (1974), pp. 50–57.
- [250] Burt Totaro. “The curvature of a Hessian metric”. In: *International Journal of Mathematics* 15.04 (2004), pp. 369–391.



- [251] Joel A. Tropp. “User-Friendly Tail Bounds for Sums of Random Matrices”. In: *Foundations of Computational Mathematics* 12.4 (2012), pp. 389–434.
- [252] Ulrich Trottenberg, Cornelius W Oosterlee, and Anton Schuller. *Multigrid*. Academic press, 2000.
- [253] Pravin M. Vaidya. “A New Algorithm for Minimizing Convex Functions over Convex Sets”. In: *FOCS*. 1989, pp. 338–343.
- [254] Pravin M Vaidya. “An algorithm for linear programming which requires  $O(((m+n)n^2+(m+n)1.5n)L)$  arithmetic operations”. In: *Mathematical Programming* 47.1-3 (1990), pp. 175–201.
- [255] Pravin M. Vaidya. “Reducing the Parallel Complexity of Certain Linear Programming Problems”. In: *FOCS*. 1990, pp. 583–589.
- [256] Pravin M. Vaidya. “Solving linear equations with symmetric diagonally dominant matrices by constructing good preconditioners.” Unpublished manuscript UIUC 1990. A talk based on the manuscript was presented at the IMA Workshop on Graph Theory and Sparse Matrix Computation, October 1991, Minneapolis. 1990.
- [257] Pravin M Vaidya. “Speeding-up linear programming using fast matrix multiplication”. In: *Foundations of Computer Science, 1989., 30th Annual Symposium on*. IEEE. 1989, pp. 332–337.
- [258] Pravin M Vaidya and David S Atkinson. “A technique for bounding the number of iterations in path following algorithms”. In: *Complexity in Numerical Optimization* (1993), pp. 462–489.
- [259] Lieven Vandenbergh and Stephen Boyd. “Semidefinite programming”. In: *SIAM review* 38.1 (1996), pp. 49–95.
- [260] Yehuda Vardi and Cun-Hui Zhang. “The multivariate L1-median and associated data depth”. In: *Proceedings of the National Academy of Sciences* 97.4 (2000), pp. 1423–1426.
- [261] S. Vempala. “Geometric Random Walks: A Survey”. In: *MSRI Combinatorial and Computational Geometry* 52 (2005), pp. 573–612.
- [262] Santosh S Vempala. “Recent progress and open problems in algorithmic convex geometry”. In: *LIPICs-Leibniz International Proceedings in Informatics*. Vol. 8. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik. 2010.
- [263] Vincenzo Viviani. “De Maximis et Minimis Geometrica Divinatio Liber 2”. In: *De Maximis et Minimis Geometrica Divinatio* (1659).
- [264] Jens Vygen. “A note on Schrijver’s submodular function minimization algorithm”. In: *Journal of Combinatorial Theory, Series B* 88.2 (2003), pp. 399–402.
- [265] S. Fujishige W. Cui. “A primal algorithm for the submodular flow problem with minimum mean cycle selection”. In: *Journal of the Operations Research Society of Japan* (1988).
- [266] C Wallacher and Uwe T Zimmermann. “A polynomial cycle canceling algorithm for submodular flows”. In: *Mathematical programming* 86.1 (1999), pp. 1–15.
- [267] Alfred Weber. *The Theory of the Location of Industries*. Aber den I der Industrien. Chicago University Press, 1909.
- [268] E. Weiszfeld. “Sur le point pour lequel la somme des distances de n points donnees est minimum”. In: *Tohoku Mathematical Journal* (1937), pp. 355–386.

- [269] Virginia Vassilevska Williams. “Multiplying matrices faster than Coppersmith-Winograd”. In: *Proceedings of the forty-fourth annual ACM symposium on Theory of computing*. ACM. 2012, pp. 887–898.
- [270] Przemyslaw Wojtaszczyk. *Banach spaces for analysts*. Vol. 25. Cambridge University Press, 1996.
- [271] Guoliang Xue and Yinyu Ye. “An efficient algorithm for minimizing a sum of Euclidean norms with applications”. In: *SIAM Journal on Optimization* 7 (1997), pp. 1017–1036.
- [272] Yinyu Ye. “Complexity analysis of the analytic center cutting plane method that uses multiple cuts”. In: *Mathematical Programming* 78.1 (1996), pp. 85–104.
- [273] Yinyu Ye. *Interior point algorithms: theory and analysis*. Vol. 44. John Wiley & Sons, 2011.
- [274] David B Yudin and Arkadii S Nemirovski. “Evaluation of the information complexity of mathematical programming problems”. In: *Ekonomika i Matematicheskie Metody* 12 (1976), pp. 128–142.
- [275] Zhizhen Zhao and Amit Singer. “Rotationally invariant image representation for viewing direction classification in cryo-EM”. In: *Journal of structural biology* 186.1 (2014), pp. 153–166.
- [276] U Zimmermann. “Minimization on submodular flows”. In: *Discrete Applied Mathematics* 4.4 (1982), pp. 303–323.
- [277] Uwe Zimmermann. “Negative circuits for flows and submodular flows”. In: *Discrete applied mathematics* 36.2 (1992), pp. 179–189.
- [278] Anastasios Zouzias. “A Matrix Hyperbolic Cosine Algorithm and Applications”. In: *Proceedings of the 39th International Colloquium on Automata, Languages and Programming (ICALP)*. 2012, pp. 846–858.