

Graduate School Introductory Computational Simulation Course Pedagogy

by

Laura L. Proctor

Submitted to the Department of Computation for Design and
Optimization
in partial fulfillment of the requirements for the degree of
Master of Science in Computation for Design and Optimization
at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

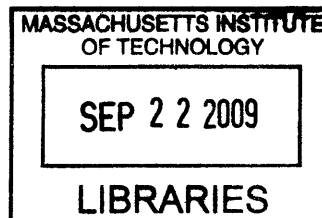
June 2009

© Massachusetts Institute of Technology 2009. All rights reserved.

Author
Department of Computation for Design and Optimization
January 9, 2009

Certified by
Jacob K. White
Cecil H. Green Professor of Electrical Engineering & Computer
Science
Thesis Supervisor

Accepted by
Jaime Peraire
Professor of Aeronautics and Astronautics
Co-Director, Department of Computation for Design & Optimization



ARCHIVES

Graduate School Introductory Computational Simulation

Course Pedagogy

by

Laura L. Proctor

Submitted to the Department of Computation for Design and Optimization
on January 12, 2009, in partial fulfillment of the
requirements for the degree of
Master of Science in Computation for Design and Optimization

Abstract

Numerical methods and algorithms have developed and matured vastly over the past three decades now that computational analysis can be performed on almost any personal computer. There is a need to be able to teach and present this material in a manner that is easy for the reader to understand and be able to go forward and use. Three popular courses at MIT were without lecture notes; in this thesis the lecture notes are presented. The first chapter covers material taught in Numerical Methods for Partial Differential Equations (2.097/6.339/16.920) specifically the Integral Equation Methods section of this course, chapter two shows the notes for the course Introduction to Numerical Simulation (2.096/6.336/16.910), and chapter three contains the notes for the class Foundations of Algorithms and Computational Techniques in Systems Biology (6.581/20.482). These course notes give a broad overview of many algorithms and numerical methods that one can use to solve many problems that span many fields - from biology to aerospace to electronics to mechanics.

Thesis Supervisor: Jacob K. White
Title: Cecil H. Green Professor

Acknowledgments

I would like to express my unending gratitude to Professor Jacob White for guiding me through this time. He has shown me a compassion that I can only hope to emulate, a patience that is rare, and a massive knowledge for numerical methods and integral equation methods. All of this, and a great sense of humor that puts one at ease.

A great thanks to the other professors with whom I've had the opportunity to work – notably Professor Luca Daniel for whom I was a teaching assistant in the Introduction to Numerical Simulation course. He is a great teacher, very enthusiastic, and easy to approach. Also, Professor Jaime Peraire, for whom I was a teaching assistant in the Numerical Methods for Partial Differential Equations course. He is also a great teacher, an understanding person, straightforward, and respectable.

I also want to thank my colleagues who have been wonderful in their own rite, and always willing to lend a hand: Bo Kim, JungHoon Lee, Brad Bond, Tarek Moselhy, Lei Zhang, Kin Sou, Dima Vasilyev, Jay Bardhan, Homer Reid, David Joe Willis, and Carlos Pinto Coelho.

In addition, I want to thank my wonderful fiance, Kevin Flaherty, who has been there for me throughout my long graduate student career.

Contents

1	Integral Equation Methods	13
1.1	Discretization of Boundary Integral Equations	13
1.2	Numerical Quadrature	42
1.3	First and Second Kind Equations	63
1.4	Radiation Conditions and Formulations	85
1.5	First and Second Kind Theory, part 2	105
1.6	Fast Algorithms for Integral Equation Methods	115
2	Foundations of Algorithms and Computational Techniques in Sys-	
	tems Biology	139
2.1	Motivation/Overview	140
2.2	Models of Proteins	142
2.3	Discrete Conformational Search	145
2.4	Binding and Docking	147
2.5	Binding and Docking - Molecular Dynamics Simulation	148
2.6	Molecular Dynamics and Electrostatics	150
2.7	Continuum Electrostatic Modeling I	152
2.8	Continuum Electrostatic Modeling II	154
2.9	Electrostatic Contributions to Binding and Design	156
2.10	Electrostatics Modeling	158
2.11	Statistical Mechanics	160
2.12	Statistical Mechanics	162
2.13	Formulating Models	163

2.14	Nonlinear Dynamics and Stability	165
2.15	Steady-State Problems	167
2.16	Parameter Fitting and Estimation	169
2.17	Parameter Estimation; Robustness, Fragility, Control	171
2.18	2-D and 3-D Light Microscopy; Image Reconstruction	172
2.19	Deconvolution	173
2.20	Deconvolution II	175
2.21	Blind Deconvolution	177
2.22	Optical Flow	178
2.23	High-Throughput Data and Analysis	179
2.24	Inference and Statistics	180
3	Introduction to Numerical Simulation	183
3.1	Example Problems and Basic Equations	184
3.2	Equation Formulation - Node Branch Stamping	194
3.3	Equation Formulation - Nodal Analysis	204
3.4	Linear Systems - LU	214
3.5	Linear Systems - Conditioning	222
3.6	Linear Systems - LU for Sparse Systems	236
3.7	Linear Systems - QR Factorization	243
3.8	Linear Systems - GCR Iterative Method	250
3.9	Linear Systems - GCR Convergence	266
3.10	Linear Systems - GCR Preconditioners	281
3.11	NonLinear Systems - Newton Method	288
3.12	NonLinear Systems - Multidimensional Newton	305
3.13	NonLinear Systems - A Case Study	316
3.14	ODE - Backward Euler, Forward Euler, Trapezoidal Rule	329
3.15	ODE - Multistep Methods I	341
3.16	ODE - Multistep Methods II	354
3.17	ODE - Periodic Steady State Analysis	362

3.18 PDE - Finite Difference Methods 1D	384
3.19 PDE - Finite Difference Methods 3D	396
3.20 PDE - Finite Element Methods & GCR Preconditioners	407
3.21 PDE - BEM Integral Equation Method I	429
3.22 PDE - BEM Integral Equation Method II	444
3.23 PDE - Fast Methods for Integral Equation Solvers	456
3.24 Model Order Reduction I	465
3.25 Model Order Reduction II	474
3.26 Model Order Reduction III	487

Introduction

Numerical methods and algorithms have developed and matured vastly over the past three decades now that computational analysis can be performed on almost any personal computer. There is a need to be able to teach and present this material in a manner that is easy for the reader to understand and be able to go forward and use. Three popular course at MIT were without lecture notes; in this thesis the lecture notes are presented. The first chapter covers material taught in Numerical Methods for Partial Differential Equations (2.097/6.339/16.920) specifically the Integral Equation Methods section of this course, chapter two shows the notes for the course Introduction to Numerical Simulation (2.096/6.336/16.910), and chapter three contains the notes for the class Foundations of Algorithms and Computational Techniques in Systems Biology (6.581/20.482). These course notes give a broad overview of many algorithms and numerical methods that one can use to solve many problems that span many fields - from biology to aerospace to electronics to mechanics.

Chapter 1

Integral Equation Methods

Numerical Methods for Partial Differential Equation is a course that covers several techniques for discretizing partial differential equations in order to solve them. Very often partial differential equations do not have analytic solutions and one needs to apply an appropriate method to solve them. In this thesis, the Integral Equation Methods used for solving Partial Differential Equations is covered.

1.1 Discretization of Boundary Integral Equations

Numerical Methods for PDEs

Boundary Element Methods, Lecture 1
Introduction to Discretization of Boundary Integral Equations

L. Proctor, S. De, C. Coelho, D. Willis, X. Wang, & J. White

November 23, 2008

1 Module Outline - 6 Lectures

Overview Integral Equation Methods

Applications: Aero, MEMS, IC Design

Informal Overview:

Formulation & Discretization

Quadrature and Cubature for computing integrals

1-D, 2-D, and dealing with Singularities

1st and 2nd Kind Theory - Part 1

Discretization

Formulating 3-D Integral Equations

Radiation Conditions, Ansatz & Green's Formulations

1st and 2nd Kind Theory - Part 2

Convergence theory

Fast Multipole and FFT-based methods

2 Outline for Today

Background

Exterior versus interior problems

Point source approach

Test Function Selection

Collocation Method

Galerkin Method

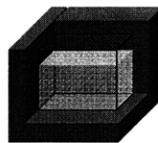
Some issues in 3D

Singular integrals

3 Background

3.1 Interior vs Exterior Problems

Interior



Exterior



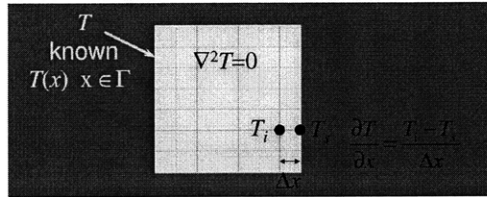
Temperature in a tank Ice cube in a bath

What is the heat distribution?

Heat flow = Thermal conductivity $\int_{surface} \frac{\partial T}{\partial n}$

3.1.1 The Interior Problem

Example: Heat Distribution in a Tank



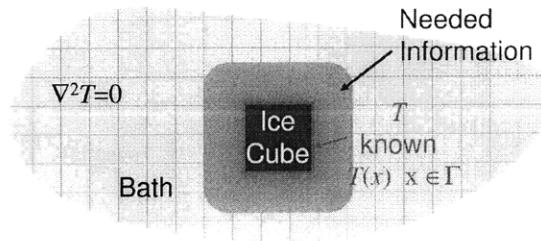
How does one determine the heat distribution using the finite difference method?

Use the above example of heat distribution in a tank, to see what is meant by an interior problem. This is a Dirichlet problem because $T(\vec{x})$, the temperature, is defined on the surface $\vec{x} \in \Gamma$. The steady state two dimensional heat flow equation is defined using the Laplace Equation, $\nabla^2 T(\vec{x}) = 0$. The domain of this problem is clearly the interior of the tank, $\vec{x} \in \Omega$.

3.1.2 The Exterior Problem

Up until this time, all problems that have been studied using finite-element and finite-difference methods have been *interior* problems, but now, we begin to wonder how to form an exterior problem. This poses some problems for these other methods such as generating the grid as well as mesh truncation.

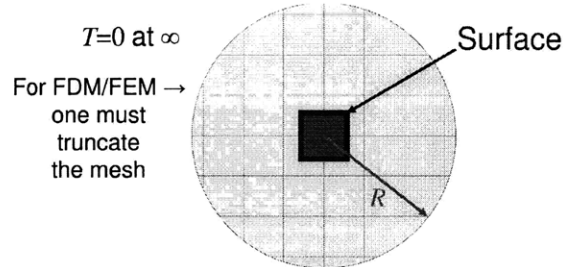
Example: Ice Cube in a Bath



Above is a problem showing an *ice cube in a bath* which is used to indicate an exterior problem. Again, this is a Dirichlet problem because $T(\vec{x})$ is defined on the surface $\vec{x} \in \Gamma$ and described by the steady state two dimensional heat flow equation defined using Laplace's Equation, $\nabla^2 T(\vec{x}) = 0$. The problem domain is the infinitely extending region exterior to the ice cube. A point that we will expand in a later lecture is that with exterior problems, an additional boundary condition is needed to specify what happens at a large distance away from our point source. Assuming there are no heat sources exterior to the cube will impose the following *radiation* boundary condition

$$\lim_{\|\vec{x}\| \rightarrow \infty} T(\vec{x}) \rightarrow 0.$$

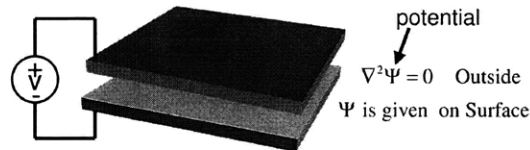
Suppose that for this specific problem, we are only interested in what occurs at the surface of the cube. It seems inefficient to use the finite-difference or finite-element methods where one needs to compute the temperature everywhere in Ω .



3.2 Examples

3.2.1 Computation of Capacitance

What is the capacitance?



$$\text{Capacitance} = \text{Dielectric Permittivity} \int \frac{\partial \Psi}{\partial n}$$

In the example in the slide, the yellow plates form a parallel-plate capacitor with an applied voltage V . In this 3-D electrostatics problem, the electrostatic potential Ψ satisfies Laplace's equation $\nabla^2 \Psi(x) = 0$ in the region exterior to the plates, and the potential is known on the surface of the plates (Dirichlet boundary condition). Furthermore, far from the plates,

$$\lim_{\|\vec{x}\| \rightarrow \infty} \Psi(\vec{x}) \rightarrow 0.$$

(Exterior Radiation Boundary Condition to be studied further in a future lecture). The value of interest is the capacitance, C , which satisfies

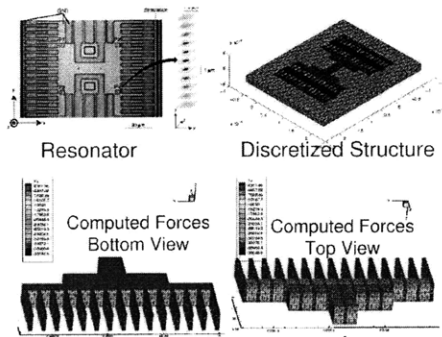
$$q = CV$$

where q , the net charge on one of the plates, is given by the surface normal of the potential integrated over one plate and scaled by a dielectric permittivity.

Note 1**Example 1: Capacitance problem**

This is a typical application example, determining the charge density on the surface of conducting plates given an applied voltage. In this particular example, the top plate potential is $\Psi = 0.5V$ and the bottom plate potential is $\Psi = -0.5V$, where V is the voltage noted in the figure.

For this exterior Dirichlet problem, one can write an integral equation that relates the surface charge density on the plates σ to the potential on the plates. This integral equation, $\Psi(\vec{x}) = \int_{\Gamma} \frac{1}{\|\vec{x}-\vec{x}'\|} \sigma(\vec{x}') dS'$, is often referred to by physicists as the superposition integral. In the integral equation, x is any point on the plate surfaces and the surface being integrated over is the union of the top and bottom plate surfaces. Note that the integration surface is not a connected domain, but this presents no difficulties.

3.2.2 Drag Force in a Microresonator

Note 2**Example 2: Drag force in a MEMS device**

The example in the slide is a microresonator, it is a structure that can be made to vibrate using electrostatic forces. The changing character of those vibrations can be used to sense rotation. The particulars of how the microresonator operates is not directly relevant to our discussion of integral equations, except for one point. In order to determine how much energy is needed to keep the microresonator vibrating, it is necessary to determine the fluid drag force on comb structures shown in the bottom part of the slide. The fluid is the air surrounding the structure, and at the micron-scale of these devices, air satisfies the incompressible Stokes equation,

$$\begin{aligned} \nabla^2 u(x) &= \nabla p(x) \\ \nabla \cdot u(x) &= 0 \end{aligned} \tag{1}$$

where u is the fluid velocity and p is the pressure. By specifying the comb velocity, and then computing the surface pressure and the normal derivative of velocities tangent to the surface, one can determine the net drag force on the comb. Once again, this is a problem in which the known quantities (the

comb velocity) and the quantities of interest (the derivative of the tangential components of fluid velocity) are on the surface.

3.2.3 Aircraft Drag



Discretization for F-18 pressure simulation (no lift)
Inviscid, Irrotational, Steady Flow
Potential flow: $\nabla^2 u(x) = 0$ $\nabla u = \text{velocity}$

Note 3

Example 3: Aircraft Drag

The potential flow model for aircraft drag computation will be discussed in more detail in subsequent lectures, so we only give a brief description here. In order to compute the drag on the wing of an aircraft, one must determine the difference between the wing velocity and the velocity of the air very close to the wing. If the air can be assumed inviscid, irrotational, and incompressible, the velocity is given by the gradient of a scalar potential which satisfies Laplace's equation. The boundary conditions for the Laplace's equation are given as a velocity boundary condition on the aircraft surface, equivalently a Neumann condition on the potential, and it is usually assumed that the potential approaches zero at infinity. The boundary condition at infinity is more subtle than it may seem, as we shall see in later lectures. Finally, it is common to introduce an artificial boundary in the domain, and specify a condition on that boundary to introduce rotational effects.

3.2.4 Capacitance of Microprocessor Signal Lines



Note 4 **Example 4: Capacitance of microprocessor signal lines**

This last example in the above slide is a picture of the wiring on a microprocessor integrated circuit. A typical microprocessor has millions of wires, so we are only looking at a small piece of a processor. The critical problem in this example is determining how long signals take to get from the output of a logical gate to the input of the next gate. To compute that delay, one must determine the capacitance on each of the wires given in the slide picture. To do so requires computing charges given electrostatic potentials as noted above.

3.3 Advantages of Integral Equation Method

3.3.1 What is common about these examples?

Exterior Problems

MEMS device - fluid (air) creates drag

Aircraft Design - exterior air flow

Signal Line - Exterior fields.

Quantities of interest are on surface

MEMS device - Just want surface traction force

Aircraft Design - Just want surface tangent velocities

Signal Line - Just want surface charge.

Exterior problem is linear and space-invariant

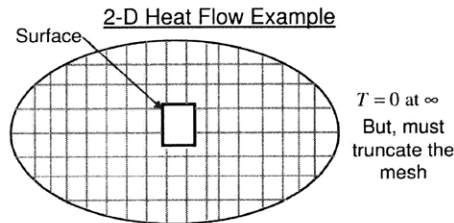
MEMS device - Exterior Stoke's flow equation (linear)

Aircraft Design - Laplace's equation, plus wakes.

Signal line - Laplace's equation in free spce (linear)

But problems are geometrically very complex

3.3.2 Why not use FDM / FEM?



Only need $\frac{\partial T}{\partial n}$ on the surface, but T is computed everywhere.

Must truncate the mesh, $\Rightarrow T(\infty) = 0$ becomes $T(R) = 0$.

Consider the two dimensional exterior heat conduction problem in the above figure in which the temperature is known on the surface of the square. Suppose the quantity of interest is the total heat flow out of the square.

The temperature T satisfies

$$\nabla^2 T(x) = 0 \quad x \in \Omega$$

$$T(x) \text{ given } x \in \Gamma \tag{2}$$

$$\lim_{\|x\| \rightarrow \infty} T(x) = 0$$

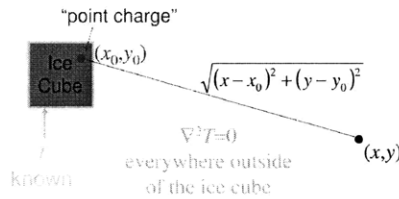
where Ω is the infinite domain outside the square and Γ is the region formed by the edges of the square.

Using finite-element or finite-difference methods to solve this problem requires introducing an additional approximation beyond discretization error. It is not possible to discretize all of Ω , as it is infinite, and therefore the domain must be truncated with an artificial finite boundary. In the slide, the artificial boundary is a large ellipse on which we assume the temperature is zero. Clearly, as the radius of the ellipse increases, the truncated problem more accurately represents the domain problem, but the number of unknowns in the discretization increases.

3.4 Point Source Approach

3.4.1 Green's Function

Heat Distribution in 2-D



$$\text{Green's Function: } T = \log \left(\sqrt{(x - x_0)^2 + (y - y_0)^2} \right)$$

From basic electrostatics, one knows that in 3-D, the potential field produced by a point charge decays inversely with the distance to the point charge. Since, roughly, one can represent any charge distribution using a sum of point charges, one can express the potential due to a charge density as a sum of point charge potentials. Therefore, point charge potentials play a special role, and are often referred to as Greens' functions for the problem.

In 2-D The potential due to a point charge is:

$$u = \log \left(\sqrt{(x - x_0)^2 + (y - y_0)^2} \right) \tag{3}$$

$$\forall (x, y) \neq (x_0, y_0)$$

In 3D The potential due to a point charge is:

$$u = \frac{1}{\sqrt{(x - x_0)^2 + (y - y_0)^2 + (z - z_0)^2}} \tag{4}$$

$$\forall (x, y, z) \neq (x_0, y_0, z_0)$$

In 2D If $u = \log \left(\sqrt{(x - x_0)^2 + (y - y_0)^2} \right)$
then $\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = 0 \quad \forall \quad (x, y) \neq (x_0, y_0)$

In 3D If $u = \frac{1}{\sqrt{(x-x_0)^2+(y-y_0)^2+(z-z_0)^2}}$
then $\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} + \frac{\partial^2 u}{\partial z^2} = 0 \quad \forall \quad (x, y, z) \neq (x_0, y_0, z_0)$

Proof: Just differentiate and see!

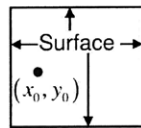
In the next few slides, we will use an informal semi-numerical approach to derive the integral form of Laplace's equation. We do this in part because such a derivation lends insight to the subsequent numerical procedures.

To start, recall from basic physics that the potential due to a point charge is related only to the distance between the point charge and the evaluation point. In 2-D the potential is given by the log of the distance, equation (3), and in 3-D the potential is inversely proportional to the distance, equation (4). These functions are sometimes referred to as Green's functions for Laplace's equation, but have the physical interpretation as the potential due to a point charge. We will be studying Green's functions in more depth later on.

▷ **Exercise 1** Show by direct differentiation that the functions in equations (3) and (4) satisfy $\nabla^2 u = 0$, in the appropriate dimension almost everywhere. ■

3.4.2 Scaling Green's Function

u is given on surface



$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = 0 \quad \text{outside}$$

$$\text{Let } u = \log \left(\sqrt{(x-x_0)^2 + (y-y_0)^2} \right)$$

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = 0 \quad \text{outside} \quad \text{Problem Solved}$$

Boundary conditions are not satisfied!

A simple idea for computing the solution of Laplace's equation outside the square is to let

$$u(x, y) = \alpha \log \sqrt{(x - x_0)^2 + (y - y_0)^2}$$

where (x_0, y_0) is a point inside the square. Clearly u will always satisfy $\nabla^2 u = 0$ outside the square, but u may not match the boundary conditions. By adjusting

α , it is possible to make sure to match the boundary conditions at at least one point.

This concept is applied to a circle as a simple example of how to match the boundary conditions.

$T = \log(\sqrt{x^2 + y^2})$
 $T(\bar{x} \in \Gamma) = \log R = 1?$
 Multiply by a constant:
 $T = \alpha \log(\sqrt{x^2 + y^2})$
 charge strength

▷ **Exercise 2** Suppose the potential on the surface of the square is a constant. Can you match that constant potential everywhere on the perimeter of the square by judiciously selecting α ? ■

u is given on surface

 $\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = 0$ outside

$$u = \sum_{i=1}^n \alpha_i \log(\sqrt{(x-x_i)^2 + (y-y_i)^2}) = \sum_{i=1}^n \alpha_i G(x-x_i, y-y_i)$$

Pick the α_i 's to match the boundary conditions!

To construct a potential that satisfies Laplace's equation and matches the boundary conditions at more points, let u be represented by the potential due to a sum of n weighted point charges in the square's interior. As shown in the slide, we can think of the potential due to a sum of charges as a sum of Green's functions. Of course, we have to determine the weights on the n point charges, and the weight on the i^{th} charge is denoted hereby α_i .

Source Strengths selected to give correct potential at test points.

$$\begin{bmatrix} G(x_{t_1} - x_1, y_{t_1} - y_1) & \cdots & G(x_{t_1} - x_n, y_{t_1} - y_n) \\ \vdots & \ddots & \vdots \\ G(x_{t_n} - x_1, y_{t_n} - y_1) & \cdots & G(x_{t_n} - x_n, y_{t_n} - y_n) \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \begin{bmatrix} \Psi(x_{t_1}, y_{t_1}) \\ \vdots \\ \Psi(x_{t_n}, y_{t_n}) \end{bmatrix}$$

To determine a system of n equations for the n α_i 's, consider selecting a set of n test points, as shown in the slide above. Then, by superposition, for each test point (x_{t_i}, y_{t_i}) ,

$$u(x_{t_i}, y_{t_i}) = \sum_{i=1}^n \alpha_i \log \sqrt{(x_{t_i} - x_0)^2 + (y_{t_i} - y_0)^2} = \sum_{i=1}^n \alpha_i G(x_{t_i} - x_0, y_{t_i} - y_0). \quad (5)$$

Writing an equation like (5) for each test point yields the matrix equation

$$\begin{bmatrix} G(x_{t_1} - x_1, y_{t_1} - y_1) & \cdots & G(x_{t_1} - x_n, y_{t_1} - y_n) \\ \vdots & \ddots & \vdots \\ G(x_{t_n} - x_1, y_{t_n} - y_1) & \cdots & G(x_{t_n} - x_n, y_{t_n} - y_n) \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \begin{bmatrix} \Psi(x_{t_1}, y_{t_1}) \\ \vdots \\ \Psi(x_{t_n}, y_{t_n}) \end{bmatrix} \quad (6)$$

The matrix A in equation (6) has some properties worth noting:

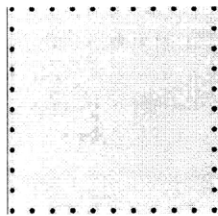
- A is dense, that is $A_{i,j}$ never equals zero. This is because every charge contributes to every potential.
- If the test points and the charge points are ordered so that the i^{th} test point is nearest the i^{th} charge, then $A_{i,i}$ will be larger than $A_{i,j}$ for all j .

The 2nd item above seems to suggest that A is diagonally dominant, but this is *not* the case. Diagonal dominance requires that the absolute *sum* of the off-diagonal entries is smaller than the magnitude of the diagonal. The matrix above easily violates that condition.

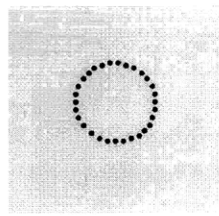
▷ **Exercise 3** Determine a set of test points and charge locations for the 2-D square problem that generates an A matrix where the magnitude of the diagonals are bigger than the absolute value of the off-diagonals, but the magnitude of the diagonal is smaller than the absolute sum of the off-diagonals. ■

3.4.3 Source Point Locations

Where should the sources be located?

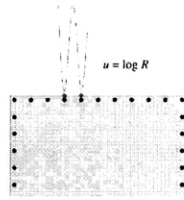


Close to the boundary

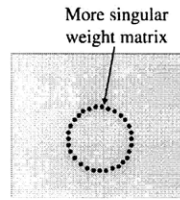


Clustered in the center

Problems with these placements:

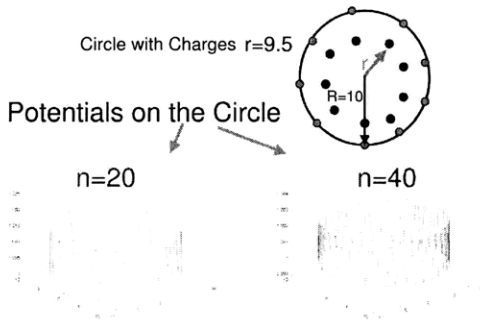


Close to the boundary



Clustered in the center

3.4.4 Computational Results



It is possible to construct a numerical scheme for solving exterior Laplace problems by adding progressively more point charges so as to match more boundary conditions. In the above graph, we show an example of using such a method to compute the potential exterior to a circle of radius 10, where the potential on the circle is given to be unity. In the example, charges are placed uniformly on a circle of radius 9.5, and test points are placed uniformly on the radius 10 circle. If 20 point charges are placed in a circle of radius 9.5, then the potential produced will be exactly one only at the 20 test points on the radius 10 circle. The potential produced by the twenty point charges on the radius 10 circle is plotted in the lower left corner of the slide above. As might be expected, the potential produced on the radius 10 circle is exactly one at the 20 test points, but then oscillates between 1 and 1.2 on the radius 10 circle. If 40 charges and test points are used, the situation improves. The potential on the circle still oscillates, as shown in the lower right hand corner, but now the amplitude is only between 1 and 1.004.

3.5 Charge Density

Want to smear point charges to the surface

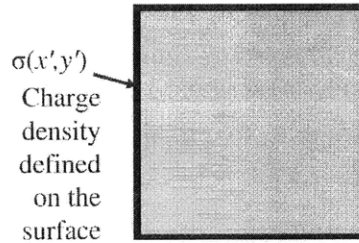


Results in an integral equation

$$\Psi(\vec{x}) = \int_{\Gamma} G(\vec{x}, \vec{x}') \sigma(\vec{x}') dS' \quad (7)$$

How do we solve the integral equation?

In equation (7) for which variable are we trying to solve?



$$\Psi(x, y) = \int_{\Gamma} \sigma(x', y') \log \sqrt{(x - x')^2 + (y - y')^2} dx' dy'$$

Single Layer Potential

The oscillating potential produced by the point charge method is due to the rapid change in potential as the separation between evaluation point and point charge shrinks. If the point charges could be smeared out, so that the produced potential did not rise to infinity with decreasing separation, then the resulting computed potential would not have the oscillation noted on the previous slide. In addition, it makes the most sense to smear the point charges onto the surface, as then the charge density and the known potential have the same associated geometry. The result is the integral equation (7), where now the unknown is a charge *density* on the surface and the potential due to that charge density is given by the well-known superposition integral. In the case of two or three dimensional Laplace problems, $G(\vec{x}, \vec{x}')$ can be written as $\hat{G}(\vec{x} - \vec{x}')$, as the potential is only a function of distance to the charge density and not a function of absolute position. For such a Green's function, this equation is,

$$\Psi(\vec{x}) = \int_{\Gamma} \hat{G}(\vec{x} - \vec{x}') \sigma(\vec{x}') dS', \quad (8)$$

which one may recognize from system theory as a convolution integral. This connection is quite precise. A space-invariant system has an impulse response, which is usually referred to as a Green's function. The output, in this case, the potential, is a convolution of the impulse response with the input, in this case, the charge density. Such an integral form of the potential is referred to as a *single layer potential*.

The single layer potential is an example of a class of integral equations known as “Fredholm integral equation of the First Kind”. A Fredholm integral equation of the Second Kind results when the unknown charge density exists not only under the integral sign but also outside it. An example of such an equation is

$$\Psi(\vec{x}) = \sigma(\vec{x}) + \int_{\Gamma} K(\vec{x} - \vec{x}')\sigma(\vec{x}')dS'. \tag{9}$$

Fredholm integral equations, in which the domain of integration is fixed, usually arise out of boundary value problems. Initial value problems typically give rise to the so-called Volterra integral equations, where the domain of integration depends on the output of interest. For example, consider the initial value problem

$$\begin{aligned} \frac{dx(t)}{dt} &= tx(t); \quad t \in [0, T], \quad T > 0. \\ x(t = 0) &= x_0 \end{aligned}$$

The “solution” of this equation is the following Volterra integral equation:

$$x(t) = x_0 + \int_0^t \xi x(\xi)d\xi.$$

4 Basis Functions

4.1 Basic Idea

Represent $\sigma(x) = \sum_{i=1}^n \alpha_i \underbrace{\varphi_i(x)}_{\text{Basis Functions}}$

Example Basis
 Represent circle with straight lines
 Assume σ is constant along each line



The basis functions are “on” the surface

Basis Functions can be used to approximate the surface charge density in a similar way in which they approximate geometry for finite elements.

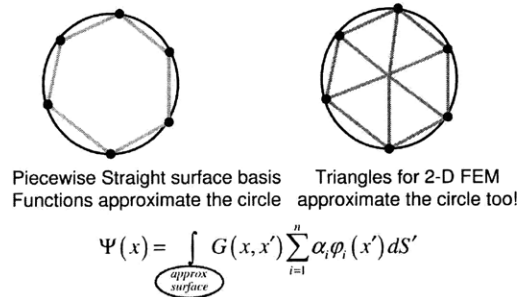
Numerical solution of the single layer potential

As we have studied extensively in the finite-element section of the course, one approach to numerically computing solutions to partial differential equations is to represent the solution approximately as a weighted sum of basis functions. Then, the original problem is replaced with the problem of determining the basis function weights. In finite-element methods, the basis functions exist in a volume, for integral equations they typically exist on a surface. For 2-D problems that means the basis functions are restricted to curves and in 3-D the

basic functions are on physical surfaces.

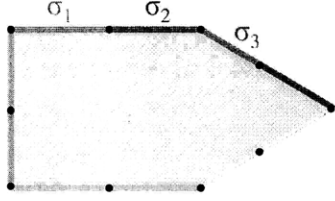
As an example, consider the circle in the above figure. One could try to represent the charge density on the circle by breaking the circle into n sub-arcs, and then assume the charge density is a constant on each sub-arc. Such an approach is not commonly used. Instead the geometry is usually approximated along with the charge density. In this example case, shown in the center right of the slide, the sub-arcs of the circle are replaced with straight sections, thus forming a polygon. The charge density is assumed constant on each edge of the polygon. The result is a piecewise constant representation of the charge density on a polygon.

4.2 Geometric Approximation

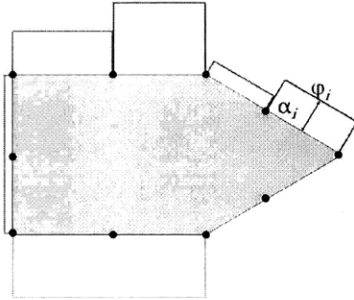


The idea that both the geometry and the unknown charge density has been approximated is not actually a new issue. As shown in the figure in the above slide, if FEM methods are used to solve an interior problem, and triangular elements are used, then the circle is approximated to exactly the same degree as when straight sections replace the sub-arcs for the surface integral equation. As shown at the bottom of the above slide, we can substitute the basis function representation into the integral equation, but then we should also note that the integral is now over the approximated geometry. It is common, but not mathematically justified, to ignore the errors generated by the geometry approximation. We will also ignore the error in the geometric approximation in our analyses, just for simplicity. In the case of polygonal geometries, there is no geometric approximation, so there is at least one case where the assumption is precise. It should be noted, however, that there are often analytic results only for smooth geometries, and then before making comparisons to such analytic results, it is necessary to examine the effect of the approximated geometry.

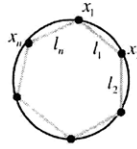
If the original problem is a polygon



there is no geometric approximation



4.2.1 Piecewise Constant Straight Sections



- 1) Pick a set of n Points on the surface
- 2) Define a new surface by connecting points with n lines.
- 3) Define $\varphi_i(x) = 1$ if x is on line l_i , otherwise, $\varphi_i(x) = 0$

$$\Psi(\vec{x}) = \int_{\text{approx surface}} G(\vec{x}, \vec{x}') \sum_{i=1}^n \alpha_i \varphi_i(\vec{x}') dS' = \sum_{i=1}^n \alpha_i \int_{\text{line } l_i} G(\vec{x}, \vec{x}') dS'$$

How do we determine the α_i 's?

We complete the description of using constant charge densities on straight sections as the basis. If we substitute this example basis function into the integral equation, as is done above, the result is to replace the original integration of the product of the Green's function and the density with a weighted sum of integrals over straight lines of just the Green's function. The next step is then to develop an approach for determining the weights, denoted here by α_i 's.

5 Test Points

5.1 Residuals

5.1.1 Definition and Minimization

$$R(\vec{x}) = \Psi(\vec{x}) - \int_{\text{approx surface}} G(\vec{x}, \vec{x}') \sum_{i=1}^n \alpha_i \varphi_i(\vec{x}') dS' \quad (10)$$

Pick the α_i 's to minimize $R(\vec{x})$

General Approach: Pick a set of test functions ϕ_1, \dots, ϕ_n and force $R(\vec{x})$ to be orthogonal to the set:

$$\int \phi_i(\vec{x}) R(\vec{x}) dS = 0 \quad \forall i \quad (11)$$

One way of assessing the accuracy of the basis function based approximation of the charge density is to examine how well the approximation satisfies the integral equation. To be more precise, we define the residual associated with the integral equation and an approximate solution, equation (10). Note that $R(\vec{x})$ is just the difference between the given potential on the surface and the potential produced by the approximated charge density. Note also that the equation is now over the approximate geometry and therefore \vec{x} and \vec{x}' are both on the approximated surface.

If the representation satisfies the integral equation exactly, then the residual $R(\vec{x})$ will be zero for all \vec{x} and the approximate solution is equal to the exact solution (provided the integral equation has a unique exact solution ... more on this later). In general, though, this is not possible, and instead we will try to pick the basis function weights, the α_i 's, to somehow minimize $R(\vec{x})$. One approach to minimizing $R(\vec{x})$ is to make it orthogonal to a collection of test functions, which may or may not be related to the basis functions, equation (11). Enforcing orthogonality in this case means ensuring that the integral of the product of $R(\vec{x})$ and $\phi(\vec{x})$ over the surface is zero.

5.1.2 Residual Minimization Using Test Functions

$$\boxed{\int \phi_i(\vec{x}) R(\vec{x}) dS = 0} \Rightarrow$$

$$\int \phi_i(\vec{x}) \Psi(\vec{x}) dS - \iint_{\text{approx surface}} \phi_i(\vec{x}) G(\vec{x}, \vec{x}') \sum_{j=1}^n \alpha_j \varphi_j(\vec{x}') dS' dS = 0 \quad (12)$$

We will generate different methods by choosing the ϕ_1, \dots, ϕ_n

Collocation : $\phi_i(\vec{x}) = \delta(\vec{x} - \vec{x}_{t_i})$ (point matching)

Galerkin Method : $\phi_i(\vec{x}) = \varphi_i(\vec{x})$ (basis = test)

Weighted Residual Method : $\phi_i(\vec{x}) = 1$ if $\varphi_i(\vec{x}) \neq 0$ (averages)

As noted in the equation (12), by substituting the definition of the residual into the equation (11), it is possible to generate n equations, one for each test function. The generated equation has two integrals. The first is a surface integral of the product of the given potential with the test function. The second integral is a double integral over the surface. The integrand of the double integral is a product of the test function, the Green's function, and the charge density representation.

Three different numerical techniques can be derived by altering the test functions.

5.2 Collocation

Collocation: $\phi_i(\vec{x}) = \delta(\vec{x} - \vec{x}_{t_i})$ (point matching)

$$\int \delta(\vec{x} - \vec{x}_{t_i}) R(\vec{x}) dS = R(\vec{x}_{t_i}) = 0 \Rightarrow$$

$$\sum_{j=1}^n \alpha_j \int_{\text{approx surface}} \overbrace{G(\vec{x}_{t_i}, \vec{x}') \varphi_j(\vec{x}') dS'}^{A_{i,j}} = \Psi(\vec{x}_{t_i}) \Rightarrow$$

$$\begin{bmatrix} A_{1,1} & \cdots & A_{1,n} \\ \vdots & \ddots & \vdots \\ A_{n,1} & \cdots & A_{n,n} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \begin{bmatrix} \Psi(\vec{x}_{t_1}) \\ \vdots \\ \Psi(\vec{x}_{t_n}) \end{bmatrix}$$

The collocation method, described in the above slide, uses shifted impulse functions as test functions, $\phi_i(\vec{x}) = \delta(\vec{x} - \vec{x}_i)$. Impulse functions, also called “delta” functions, have a *sifting* property when integrated with a smooth function $f(\vec{x})$,

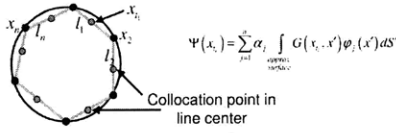
$$\int f(\vec{x}) \delta(\vec{x} - \vec{x}_i) dx = f(\vec{x}_i).$$

Impulse functions are also referred to as generalized functions, and they are specified only by their behavior when integrated with a smooth function. In the case of the impulse function, one can think of the function as being zero except for a very narrow interval around \vec{x}_i , and then being so large in that narrow interval that $\int \delta(\vec{x} - \vec{x}_i) dx = 1$.

As the summation equation in the middle of the above slide indicates, testing with impulse functions is equivalent to insisting that $R(\vec{x}_i) = 0$, or in words, that the potential produced by the approximated charge density should match the given potential at n test points. That the potentials match at the test points gives rise to the method's name, the point where the potential is exactly matched is “co-located” with a set of test points.

The $n \times n$ matrix equation at the bottom of the above slide has as its right-hand side the potentials at the test points. The unknowns are the basis function weights. The j^{th} matrix element for the i^{th} row is the potential produced at test point x_i by a charge density equal to basis function φ_j .

5.2.1 Centroid Collocation for Piecewise Constant Bases



$$\Psi(x) = \sum_{j=1}^n \alpha_j \int_{\text{line } j} G(r, x') \phi(x) dS'$$

$$\begin{bmatrix} A_{1,1} & \cdots & A_{1,n} \\ \vdots & \ddots & \vdots \\ A_{n,1} & \cdots & A_{n,n} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \begin{bmatrix} \Psi(\vec{x}_{t_1}) \\ \vdots \\ \Psi(\vec{x}_{t_n}) \end{bmatrix} \Rightarrow \Psi(\vec{x}_{t_i}) = \sum_{j=1}^n \alpha_j \underbrace{\int_{\text{line } j} G(\vec{x}_{t_i}, \vec{x}') dS'}_{A_{i,j}}$$

In the above slide, a specific collocation algorithm is described. First, the basis being used is the constant charge density on n straight sections or lines, as described above. Note that therefore the geometry is being approximated. Second, the collocation points being selected are the centroids of the basis functions, in this case just the center of each straight line. Note that the collocation point is on the approximated geometry, not the original geometry. So, one can think of the problem as having been restated to be on a polygon instead of the original circle. One could also have selected the collocation points on the original circle, but then the replacement interpretation does not hold.

In collocation, or point-matching, the charge densities on each of the straight lines are selected so that the resulting potential at the line centers matches the given potential. As the equations on this slide make clear, the matrix element $A_{i,j}$ is the potential at the center of line i due to a unit charge density along line j .

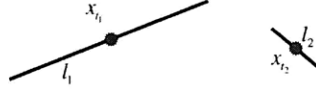
It should be noted that the matrix A is dense, the charge on line j contributes to the potential everywhere. Also note that if line j is far away from line i , then

$$A_{i,j} \approx \text{length}(\text{line}_j) \times G(\vec{x}_{t_i}, \vec{x}_{t_j}) \quad (13)$$

▷ **Exercise 4** Suppose we are using piecewise constant centroid collocation to solve a 2-D Laplace problem, so $G(x, y, x', y') = \log \sqrt{(x - x')^2 + (y - y')^2}$. Roughly how far apart do line sections i and j have to be for equation (13) to be accurate to within one percent? Assume line j has length of one. Does your answer depend on the orientation of line j ? Does your answer depend on the orientation of line i ? (You should answer yes to one of these and no to the other, do you see why?) ■

5.2.2 Centroid Collocation Generates Nonsymmetric A

$$\Psi(\vec{x}_{t_i}) = \sum_{j=1}^n \alpha_j \int_{\text{line } j} \overbrace{G(\vec{x}_{t_i}, \vec{x}') dS'}^{A_{i,j}}$$



$$A_{1,2} = \int_{\text{line } 2} G(\vec{x}_{t_1}, \vec{x}') dS' \neq \int_{\text{line } 1} G(\vec{x}_{t_2}, \vec{x}') dS' = A_{2,1} \quad (14)$$

Consider the two line sections, l_1 and l_2 given in the above figure. For Laplace problems, $G(\vec{x}, \vec{x}') = G(\vec{x}', \vec{x})$, which suggests a symmetry in the underlying integral equation that is not represented in the collocation discretization. This asymmetry is shown in equation (14) by noting that $A_{1,2} \neq A_{2,1}$. That is, the potential at the center of l_2 due to a unit charge density on l_1 is not equal to the potential at the center of l_1 due to a unit charge on l_2 .

It is possible to scale the variables to improve the symmetry, consider a change of variables

$$\hat{\alpha}_i = \alpha_i \times \text{length}(\text{line}_j).$$

In this change of variables, the unknowns $\hat{\alpha}_i$ are now the net line charges rather than the line charge densities. In this new system, $\hat{A}\hat{\alpha} = \Psi$, where the elements of the matrix \hat{A} are given by

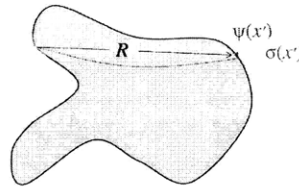
$$\hat{A}_{i,j} = \frac{1}{\text{length}(\text{line}_j)} \int_{\text{line}_j} G(\vec{x}_{t_i}, \vec{x}') dS'.$$

Under the change of variables, if line j is far away from line i , then

$$\hat{A}_{i,j} \approx G(\vec{x}_{t_i}, \vec{x}_{t_j}) \approx \hat{A}_{j,i}. \quad (15)$$

In other words, the elements of \hat{A} corresponding to distant terms are approximately symmetric.

Is $\Psi(\vec{x})$ due to $\sigma(\vec{x}')$ the same as $\Psi(\vec{x}')$ due to $\sigma(\vec{x})$?



Green's Function is due to $\log R$

▷ **Exercise 5** Give an example which shows that the scaled entries of \hat{A} can be far from symmetric. Assume we are using piecewise constant straight sections with centroid collocation and the 2-D Laplace's equation Green's function. ■

5.3 Galerkin

Galerkin: $\phi_i(x) = \varphi_i(x)$ (test=basis)

$$\int \varphi_i(x) R(x) dS = \int \varphi_i(x) \Psi(x) dS - \int \varphi_i(x) G(x, x') \sum_{j=1}^n \alpha_j \varphi_j(x') dS' dS = 0$$

$$\underbrace{\int \varphi_i(x) \Psi(x) dS}_{b_i} = \sum_{j=1}^n \alpha_j \underbrace{\int \int G(x, x') \varphi_i(x) \varphi_j(x') dS' dS}_{A_{i,j}}$$

$$\begin{bmatrix} A_{1,1} & \cdots & A_{1,n} \\ \vdots & \ddots & \vdots \\ A_{n,1} & \cdots & A_{n,n} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$

If $G(\vec{x}, \vec{x}') = G(\vec{x}', \vec{x})$ then $A_{i,j} = A_{j,i} \Rightarrow \mathbf{A}$ is symmetric

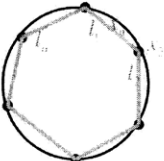
In the Galerkin method, the test functions are equal to the basis functions. In particular, one generates n equations for the basis function weights by insisting that $R(\vec{x})$ is orthogonal to each of the basis functions. Enforcing orthogonality corresponds to setting

$$\int \varphi(\vec{x}) R(\vec{x}) dS = 0$$

and substituting the definition of $R(\vec{x})$ into the orthogonality condition yields the equation in the center of the above slide.

Note that the Galerkin method yields a system of n equations, one for each orthogonality condition, and n unknowns, one for each basis function weight. Also, the system does not have the potential explicitly as the right hand side. Instead, the i^{th} right-hand side entry is the average of the product of the potential and the i^{th} basis function.

5.3.1 Galerkin for Piecewise Constant Bases



$$\underbrace{\int \Psi(x) dS}_{b_i} = \sum_{j=1}^n \alpha_j \underbrace{\int \int G(x, x') dS' dS}_{A_{i,j}}$$

$$\begin{bmatrix} A_{1,1} & \cdots & A_{1,n} \\ \vdots & \ddots & \vdots \\ A_{n,1} & \cdots & A_{n,n} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$

In the Galerkin method, the basis has constant charge density on n straight sections or lines. We will think of the problem as having been restated to be on a polygon instead of the original circle. The charge densities on each of the straight lines are selected so that the resulting line averaged potential matches the line averaged given potential. As the equations on the above slide make clear, the matrix element $A_{i,j}$ is the average potential over line i , scaled by the length of line i , due to a unit charge density along line j .

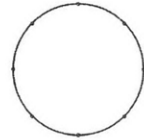
As with the collocation method, the matrix A is dense because the the charge on line j contributes to the averaged potentials everywhere. Also note that if line j is far away from line i , then

$$A_{i,j} \approx \text{length}(\text{line}_j) \times \text{length}(\text{line}_i) \times G(\vec{x}_i, \vec{x}_j) \quad (16)$$

▷ **Exercise 6** Suppose we are using piecewise constant centroid collocation to solve a 2-D Laplace problem, so $G(x, y, x', y') = \log \sqrt{(x - x')^2 + (y - y')^2}$. Roughly how far apart do line sections i and j have to be for equation (16) to be accurate to within one percent? Assume line j has length of one. Does your answer depend on the orientation of line j ? Does your answer depend on the orientation of line i ? (Your answer should be different than the answer you gave for the collocation method. Do you see why?) ■

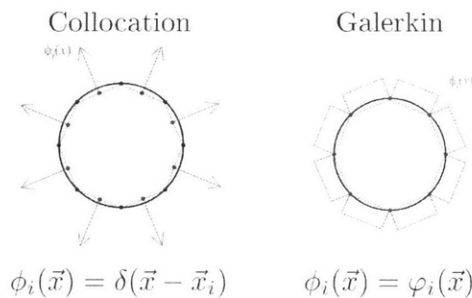
5.4 Summary

Compare the Collocation and Galerkin methods on a two-dimensional circle.

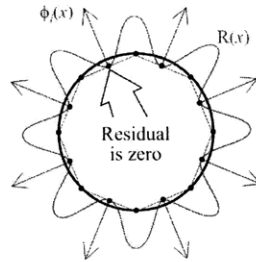


What do the test functions look like?
What do the Residuals look like?

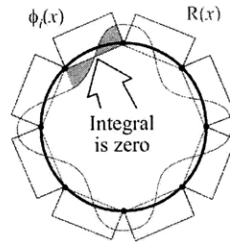
Test Functions



Collocation Method



Galerkin Method



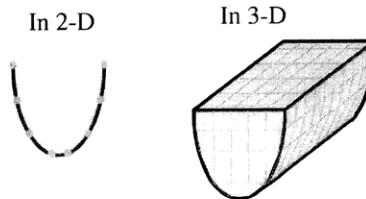
6 Issues in 3D

6.1 Geometric Representation

6.1.1 Introduction

Example: Ship's Hull

More errors are introduced with expansion of dimensions



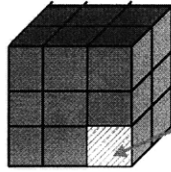
Note 6

"Leaky Panels"

Many papers in the literature on solving integral equations refer to "panel methods". The name is derived from the idea of breaking a surface into flat panels. In the application area of analyzing ocean wave forces on ship hulls, panel methods are commonly used. However, it is not possible to represent a curved hull with quadrilateral flat panels. Researchers in the area often create a best fit panelled surface in which there are gaps between the edges of the panels. Such a discretization technique is often referred to as using "leaky panels", a very compelling image.

Integral Equation : $\Psi(x) = \int_{surface} \frac{1}{\|x-x'\|} \sigma(x') dS'$

Discretize Surface into Panels



Represent $\sigma(x) \approx \sum_{i=1}^n \alpha_i \underbrace{\varphi_i(x)}_{\text{Basis Functions}}$

Panel j $\varphi_j(x) = 1$ if x is on panel j
 $\varphi_j(x) = 0$ otherwise

Consider solving the integral equation where the surface is the surface of the cube shown. The first step, as we have mentioned in previous lectures, is to develop a basis in which to represent the surface charge density σ .

The cube pictured in the slide has had its surface divided into panels, and a basis is derived from the panels. In particular, one can associate a basis function φ_j with each panel j by assigning $\varphi_j(\vec{x})$ the value one when \vec{x} is a point on panel j , and setting $\varphi_j(\vec{x}) = 0$ otherwise. If σ is approximated by a weighted combination of these basis functions, then the approximation is a piecewise constant representation of the charge density on the surface of the cube.

A few aspects of this basis set should be noted.

- The basis functions are orthogonal, that is if $i \neq j$,

$$\int \varphi_j(\vec{x}) \varphi_i(\vec{x}) dx = 0.$$

- These basis functions are normalized with respect to l_∞ , not l_2 . That is, $\|\varphi\|_\infty = 1$ but

$$\|\varphi_j\|_2^2 = \int \varphi_i(\vec{x}) \varphi_i(\vec{x}) dx = \text{panel area.}$$

6.2 Centroid Collocation

Put collocation points at panel centroids

$$\Psi(x_i) = \sum_{j=1}^n \alpha_j \int_{\text{panel } j} \frac{1}{\|x_i - x'\|} dS'$$

$A_{i,j}$

Collocation point x_i

$$\begin{bmatrix} A_{1,1} & \cdots & \cdots & A_{1,n} \\ \vdots & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ A_{n,1} & \cdots & \cdots & A_{n,n} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \begin{bmatrix} \Psi(x_1) \\ \vdots \\ \Psi(x_n) \end{bmatrix}$$

After one has decided on a basis with which to approximately represent the surface charge density, the next step is to develop a system of equations from which to determine the basis weights, denoted as the α_i 's. The most commonly used approach to forming such a system is to use collocation, though Galerkin methods are also quite widely used. Recall that in collocation, the basis function weights are determined by ensuring the the integral equation is exactly satisfied at a collection of “collocation” points. For panel methods, the most common choice for the position of the collocation points are the panel centroids, as shown in the cube diagram above.

The equation in the top of the above slide relates the potential at collocation point \vec{x}_{c_i} to the weights for the panel-based basis functions. To see how the equation was derived, consider evaluating the potential at the i^{th} collocation point using the original integral equation

$$\Phi(\vec{x}_{c_i}) = \int_{surface} \frac{1}{\|\vec{x}_{c_i} - \vec{x}'\|} \sigma(\vec{x}') dS', \quad (17)$$

where Φ is the know potential on the problem surface and σ is the unknown charge density. Substituting the approximate representation for σ ,

$$\sigma(\vec{x}) \approx \sum_{j=1}^n \alpha_j \varphi_j(\vec{x})$$

into the integral equation results in

$$\Phi(\vec{x}_{c_i}) = \int_{surface} G(\vec{x}_{c_i}, \vec{x}') \sum_{j=1}^n \alpha_j \varphi_j(\vec{x}') dS', \quad (18)$$

where $G(\vec{x}, \vec{x}') \equiv \frac{1}{\|\vec{x} - \vec{x}'\|}$ is used to simplify the formula. Exploiting the fact that $\varphi_j(\vec{x}) = 1$ if \vec{x} is on panel j , and zero otherwise, results in the formula at the top of the above slide.

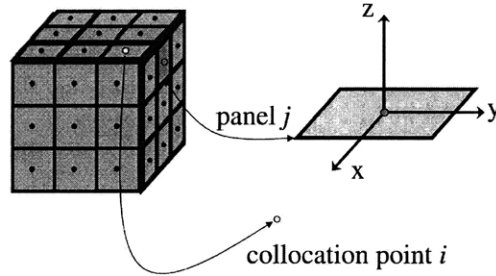
The system of equations from which to determine the basis function weights is given in the lower corner of the slide. The right hand side of the system is the vector of known potentials at the collocation points. The i, j^{th} element of the matrix A is the potential produced at collocation point i due to a unit charge density on panel j . The vector of α 's are the unknown panel charge densities.

▷ **Exercise 7** Determine a scaling of the α 's ($\hat{\alpha}_i = c_i \alpha_i$) such that the scaled matrix \hat{A} has the property

$$A_{i,j} \approx \frac{1}{\|\vec{x}_{c_i} - \vec{x}_{c_j}\|}$$

when $\|\vec{x}_{c_i} - \vec{x}_{c_j}\|$ is much larger than a panel diameter. ■

6.2.1 Calculating Matrix Elements



The diagram shows the integral formula for the matrix element $A_{i,j}$ and its approximations. The integral is given by:

$$A_{i,j} = \int_{\text{panel } j} \frac{1}{\|x_{c_i} - x'\|} dS'$$

The diagram also shows two approximations:

One point quadrature Approximation: $A_{i,j} \approx \frac{\text{Panel Area}}{\|x_{c_i} - x_{\text{centroid}_j}\|}$

Four point quadrature Approximation: $A_{i,j} \approx \sum_{j=1}^4 \frac{0.25 * \text{Area}}{\|x_{c_i} - x_{\text{point}_j}\|}$

In order to calculate the matrix entries for the system of equations described in the previous slide, recall that $A_{i,j}$ is the potential produced at collocation point i due to a unit charge density on panel j . The formula for $A_{i,j}$ is given on the top right of the above slide.

The figure on the left of the above slide is a diagram of how one typically computes the panel integral given on the top right. First, consider a shift and rotation of the coordinate system so that the panel lies in the x - y plane at $z = 0$, with the panel's center at $x = 0$, $y = 0$. The figure in the top left shows the panel in the shifted and rotated coordinate system. Note that the collocation point must also be placed in the new coordinate system.

If panel j is reasonably well separated from collocation point i , it is possible to approximate the integral given in the top right by a single point quadrature. More specifically, one could approximate the integral of $\frac{1}{\|\vec{x}_{c_i} - \vec{x}'\|}$ by a product of $\frac{1}{\|\vec{x}_{c_i} - \vec{x}_{\text{centroid}_j}\|}$ and the panel area. As show in the middle figure, a single point quadrature is like treating the panel as if it were a point charge at the panel's centroid, where the point charge's strength is equal to the panel area.

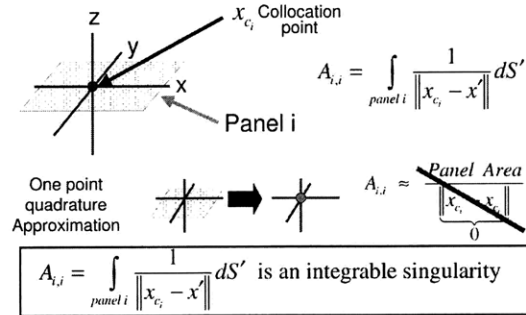
If the collocation point is close to the panel, then a single point quadrature will be insufficiently accurate. Instead, a more accurate four point quadrature scheme would be to break the panel into four subpanels, and then treat each of the subpanels as point charges at their respective centers. This simple idea is shown in the figure at the bottom of the above slide. This four point scheme is

equivalent to

$$\int_{\text{panel}_j} \frac{1}{\|\vec{x}_{c_i} - \vec{x}'\|} dS' \approx \sum_{j=1}^4 \frac{0.25 * \text{Area}}{\|\vec{x}_{c_i} - \vec{x}_{\text{point}_j}\|}$$

If the panel is a unit square in the x-y plane whose center is at the coordinate system origin, then the four \vec{x}_{point_j} 's are $(x, y, z) = (0.25, 0.25, 0)$, $(x, y, z) = (-0.25, 0.25, 0)$, $(x, y, z) = (-0.25, -0.25, 0)$, and $(x, y, z) = (0.25, -0.25, 0)$.

6.2.2 Calculating "Self Term"

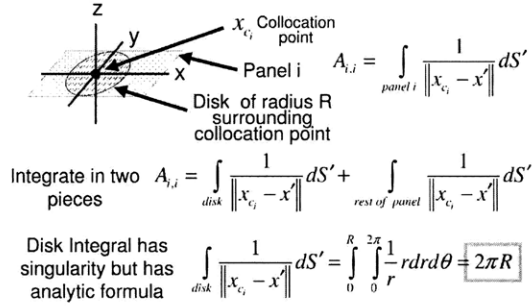


The diagonal terms $A_{i,i}$ can not be computed using the quadrature approximation given on the previous slide. To see this, consider the figure at the top left of the above slide, where a panel has been shifted and rotated into the x-y plane, and the collocation point is the center of the panel. The integral that must be computed is given on the right side of the top of the above slide.

As shown in the middle of the slide, using a single point quadrature scheme will fail, because the distance between the point charge approximation to the panel and the collocation point will be zero. Therefore, the single point formula will require computing the reciprocal of zero, which is infinite. The problem is that the integrand in

$$\int_{\text{panel}_i} \frac{1}{\|x_{c_i} - x'\|} dS' \tag{19}$$

is singular. That is, the integrand approaches infinity at a point x' which is in the domain of integration. What is not so obvious is that (19) is an integrable singularity. Therefore, even though the integrand approaches infinity at some point, the "area under the curve" is finite.



In the above slide, we both show that

$$\int_{\text{panel}_i} \frac{1}{\|\vec{x}_{c_i} - \vec{x}'\|} dS'$$

is integrable, and also give an idea about how to compute the integral.

As shown in the slide, first rotate and shift the coordinate system so that the panel is in the x - y plane at $z = 0$, and so that the collocation point (or equivalently the panel centroid) is at the origin. In this new coordinate system, the integral can be written as

$$A_{i,i} = \int_{\text{panel}(rs)_i} \frac{1}{\|\vec{x}'\|} dS'$$

where the notation $\text{panel}(rs)$ is used to indicate that the integral is over the rotated and shifted panel.

On the top left of the above slide, a circular disk of radius R and center at the collocation point is inscribed in the rotated panel. In the equations that follow the figure, it is noted that the panel integral can be recast as the sum of an integral over the disk plus an integral over the rest of the panel. The integrand in the integral over the rest of the panel is no longer singular, but the integrand in the integral over the disk is still singular.

The integral over the disk can be computed analytically by using a change of variables. After rotating and shifting the panel, the disk is in the x - y plane and its center, equal to the collocation point, is at zero. Therefore,

$$\int_{\text{disk}} \frac{1}{\|\vec{x}_{c_i} - \vec{x}'\|} dS' = \int_{\text{disk}} \frac{1}{\|\vec{x}'\|} dS'$$

Apply a change of variables as transformations of two-dimensional regions. Recall this mapping from an earlier lecture in the finite difference method. Suppose that a region $\hat{\Omega}$ in the r - θ plane is transformed one-to-one into the region Ω by differentiable equations of the form

$$x = r \cos \theta, \quad y = r \sin \theta.$$

Any function $f(x, y)$ defined on Ω can be thought of as a function $f(x(r, \theta), y(r, \theta))$ on $\hat{\Omega}$. The integrals of these functions are related by

$$\iint_{\Omega} f(x, y) dx dy = \iint_{\hat{\Omega}} f(x(r, \theta), y(r, \theta)) |J(r, \theta)| dr d\theta$$

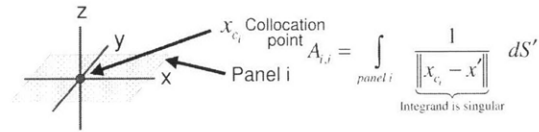
where $J(r, \theta)$ is the Jacobian determinant of the coordinate transformation where

$$J(r, \theta) = \begin{vmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} \end{vmatrix} = \begin{vmatrix} \cos\theta & -r\sin\theta \\ \sin\theta & r\cos\theta \end{vmatrix} = r(\cos^2\theta + \sin^2\theta) = r$$

So, now the integral may be put into the transformed coordinates using this transformation

$$\int_{disk} \frac{1}{\|\vec{x}_{c_i} - \vec{x}'\|} dS' = \int_0^R \int_0^{2\pi} \frac{1}{r} r d\theta dr$$

The integral over the disk is easily seen to be $2\pi R$.



1. If panel is a flat polygon, analytical formulas exist.
2. Curved panels can be handled with projection.

7 Summary

Integral Equation Methods

Exterior versus interior problems

Start with using point sources

Standard Solution Methods

Collocation Method

Galerkin Method

Integrals for 3D Problems

Singular Integrals

1.2 Numerical Quadrature

Numerical Methods for PDEs

Methods, Lecture 2
Numerical Quadrature

Notes by L. Proctor, S. De and J. White

November 26, 2008

1 Outline

Gaussian Quadrature

- Convergence properties
- Essential role of orthogonal polynomials
- Multidimensional Integrals

Techniques for singular kernels

- Adaptation and variable transformation
- Singular quadrature.

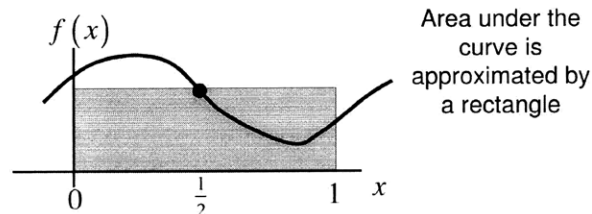
2 Introduction

Numerical Quadrature is employed as an approximation used to evaluate integrals. We seek an appropriate numerical procedure applied to a definite integral, $I\{f\} \equiv \int_a^b f(x)dx$, where the approximation is essentially of the form

$I_n\{f\} \equiv \sum_{i=1}^n \alpha_i f(x_i)$. The n distinct points, x_i are the quadrature nodes we have chosen and the quadrature coefficients, or weights, are the α_j terms. In general, we would like to have the smallest possible quadrature error, $E_n\{f\} \equiv I\{f\} - I_n\{f\}$.

2.1 Simple Quadrature Example

$$\int_0^1 f(x)dx \simeq f\left(\frac{1}{2}\right)$$

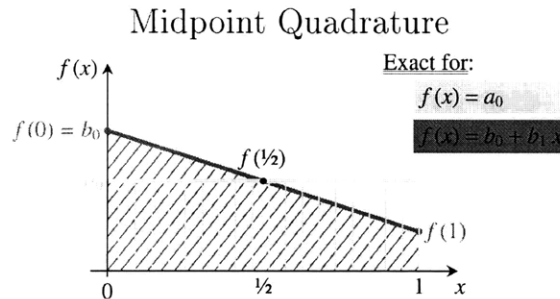


To simplify notation, consider the more generic problem of developing a good numerical technique for evaluating the integral of a function $f(x)$ on the domain $[0, 1]$. We assume that the integrand is a “smooth” function, though we will examine this assumption later. First we have developed a naive approach for obtaining a good approximation of the integral, one we call a **simple quadrature scheme**.

The simplest approach is to replace the integral with the the product of the interval (in this case one) and the integrand evaluated at a point inside the interval. If the selected evaluation point is the center of the interval, $x = 0.5$, we call the scheme **midpoint quadrature**.

A midpoint quadrature scheme replaces the area under the curve $f(x)$ by a

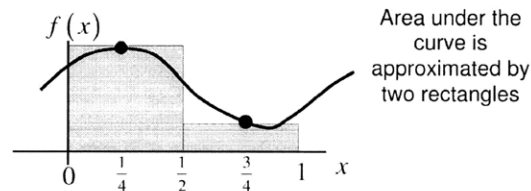
rectangle whose height is the function $f(x)$ evaluated at the midpoint $x = 0.5$. The scheme is exact when $f(x)$ is a constant. However, what is less obvious is that the scheme is exact when $f(x)$ is a line (an affine function of x) as well. The most obvious way of seeing this is by realizing that when $f(x)$ is a straight line, the area under it is a trapezoid. This trapezoid has exactly the same area as the rectangle which this scheme uses to approximate the integral (can you see why?).



▷ **Exercise 1** Suppose endpoint quadrature (in which the area under the curve is replaced by a rectangle whose height is $f(x)$ evaluated at $x = 0$ or $x = 1$) is used instead of midpoint quadrature. For what class of functions is endpoint quadrature exact? ■

2.2 Improving the Accuracy

$$\int_0^1 f(x) dx \simeq \frac{1}{2} f\left(\frac{1}{4}\right) + \frac{1}{2} f\left(\frac{3}{4}\right)$$



One way of improving the midpoint quadrature scheme is to divide the interval $[0, 1]$ into subintervals $[0, 0.5]$ and $[0.5, 1]$, then write the integral

$$\int_0^1 f(x) dx = \int_0^{0.5} f(x) dx + \int_{0.5}^1 f(x) dx,$$

and finally apply a midpoint rule to the integral in each subinterval. We obtain the scheme shown in the slide. The factor $\frac{1}{2}$ appearing in front of $f(\frac{1}{2})$ and $f(\frac{3}{4})$ are just the domain lengths.

▷ **Exercise 1** Can you come up with an expression for the error in this case? How much does the accuracy improve? ■

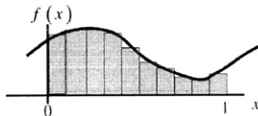
Dividing the interval into two reduces the error, now consider using n subintervals and repeating the midpoint quadrature rule on each subinterval. We obtain the scheme

$$\int_0^1 f(x)dx = \sum_{i=1}^n \underbrace{\frac{1}{n}}_{\text{subinterval length}} f(x_{c_i})$$

where the centroid of the i^{th} subinterval is $x_{c_i} = \frac{1}{2}(\frac{i-1}{n} + \frac{i}{n}) = \frac{i-\frac{1}{2}}{n}$. There is no doubt that the accuracy improves, but the key question is by how much? How does this error decrease with the number of subintervals used? And finally, are there clever ways of obtaining better accuracy with less effort?

2.3 General n -point formula

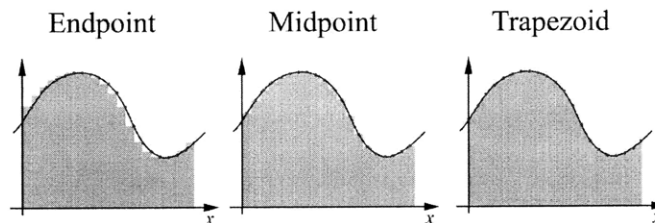
$$\int_0^1 f(x)dx \simeq \sum_{i=1}^n \underbrace{w_i}_{\text{weight } \frac{1}{n}} \underbrace{f(x_i)}_{\text{test point } \frac{i-\frac{1}{2}}{n}}$$



Key questions about the method:
How fast do the errors decay with n ?
Are there better methods?

2.4 Different Geometric Approximations

Which geometry is the most accurate?



3 General Quadrature Formula

3.1 General 1-D Formula

$$\int_0^1 f(x)dx \simeq \sum_{i=1}^n \underbrace{w_i}_{\text{Weight}} \underbrace{f(x_i)}_{\text{Evaluation Point}}$$

Free to pick the **evaluation points**.
Free to pick the **weights** for each point.

An n -point formula has $2n$ degrees of freedom!

After all the hard work we did dividing the domain into subintervals, we realize that we cannot even integrate a quadratic function exactly on the domain. There must be something that we can do to improve this scheme. We go back and look at the general form of the quadrature approximation scheme. All we are doing is approximating an integral by a weighted sum of function evaluations. So far we have been choosing these weights as the subinterval lengths. We have also been choosing the evaluation points as the center of the interval, in the midpoint quadrature scheme. The weights are just some normalizing factors which we have taken to be the fraction of the domain over which we are evaluating. The equality of areas of trapezoids and rectangles that we previously discussed gives us the extra polynomial accuracy of being able to obtain the area under a straight line exactly. So, what would happen if we were to choose both the integration points and the weights intelligently? For an n -point formula we have n weights and n evaluation points to choose. That gives us $2n$ degrees of freedom. Hence we must be able to exactly integrate a polynomial of degree at most $(2n - 1)$. This idea gives rise to the **Gaussian quadrature** scheme.

3.2 Evaluation Points & Weights Selection

Can make the result exact if $f(x)$ is a polynomial

$$f(x) = a_0 + a_1x + a_2x^2 + \dots + a_lx^l = p_l(x)$$

Select x_i 's and w_i 's such that

$$\int_0^1 p_l(x)dx = \sum_{i=1}^n w_i p_l(x_i)$$

for ANY polynomial up to (and including) l^{th} order
With $2n$ degrees of freedom, $l = 2n - 1$

Let $p_l(x)$ denote a polynomial of degree l in the variable x ($a_l \neq 0$). We want to select the weights and integration points such that the formula

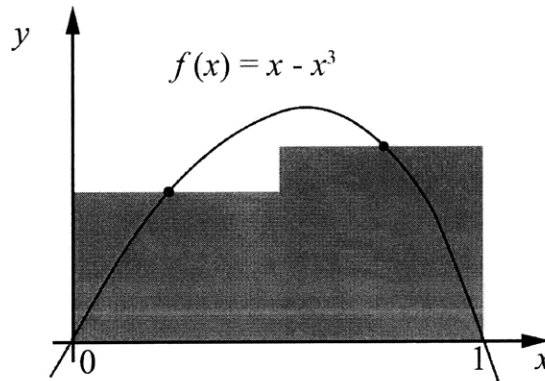
$$\int_0^1 f(x)dx = \int_0^1 p_l(x)dx = \sum_{i=1}^n w_i p_l(x_i)$$

is exact for all polynomials of degree up to (and including) l . Obviously, with $2n$ degrees of freedom, the best we can do is $l = 2n - 1$.

Note 1

Example: Third Order Polynomial

As an example, consider integratin the function $f(x) = x - x^3$. from $x = 0$ to $x = 1$. The exact solution is, $I\{f\} = \int_0^1 (x - x^3)dx = \frac{1}{4}$. It is stated above that since we have a polynomial of degree three ($l = 3$), then we will be able to find an exact solution using two point quadrature ($n = 2$). But, one must note, that finding the exact solution is not simply a matter of applying midpoint quadrature, as we have done previously. The solution using midpoint quadrature is 0.2813, not very close to 0.25. Using Gaussian Quadrature, the method of which will be studied in further detail later on, will provide the exact solution.



Assuming the weights, w_i , remain bounded, and the derivatives of $f(x)$ are bounded on $[0, 1]$,

$$\left| \int_0^1 f(x)dx - \sum_{i=1}^n w_i f(x_i) \right| \leq \frac{K}{(2n)!}$$

Gaussian quadrature converges **very** quickly!!

4 Error Analysis

4.1 Taylor Series Expansion

To derive the error of the midpoint quadrature scheme analytically, consider the interval $[0, h]$, $h > 0$ and Taylor expand $f(x)$ about the center of this interval, $\bar{x} = \frac{h}{2}$,

$$f(x) = f(\bar{x}) + \Delta(x) \frac{df(\bar{x})}{dx} + \frac{\Delta(x)^2}{2!} \frac{d^2 f(\xi)}{dx^2} \quad \text{for some } \xi \in [0, h],$$

where $\Delta(x) = x - \bar{x}$. The last term in the expansion is the Taylor series remainder. Integrating this expansion over the interval $[0, h]$

$$\int_0^h f(x) dx = hf(\bar{x}) + \frac{h^3}{24} \frac{d^2 f(\xi)}{dx^2}.$$

Hence the error in the midpoint quadrature approximation is

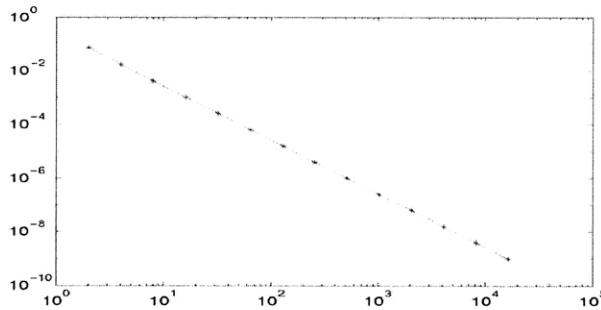
$$E = \int_0^h f(x) dx - hf(\bar{x}) = \frac{h^3}{24} \frac{d^2 f(\xi)}{dx^2}.$$

A function that is first-order polynomial in x would have zero as the second derivative, and therefore the above expression tells us that the error of midpoint quadrature for such functions is identically zero. In addition, the above expression tells us that the error scales as the cube of the domain length.

▷ **Exercise 2** If the midpoint quadrature scheme uses rectangular geometry to approximate the area under the curve, then why is the first derivative needed in the Taylor Series expansion? How does one represent the expansion of the trapezoidal approximation? ■

4.2 Example - Error vs. n

$$\int_0^1 \sin(x) dx \simeq \sum_{i=1}^n \frac{1}{n} \sin\left(\frac{i - \frac{1}{2}}{n}\right)$$



Note 2

Above is the example of integrating $f(x) = \sin(x)$ on the domain $[0, 1]$. We obtain progressively better answers to the integration by increasing the number of subintervals n . The error in evaluating the integral is plotted as a function of the number of subintervals (n) above. The error appears to be going down as $\mathcal{O}\left(\frac{1}{n^2}\right)$.

From what we have just seen, the error inside the i^{th} subinterval (of length $h = \frac{1}{n}$) is $\frac{h^3}{24} \frac{d^2 f(\xi_i)}{dx^2}$ for some $\xi_i \in [\frac{i-1}{n}, \frac{i}{n}]$. Hence, for the entire interval $[0, 1]$ we can sum these errors and obtain the error, E_n for an approximation using n subintervals as

$$E_n = \frac{nh^3}{24} \underbrace{\left(\frac{1}{n} \sum_{i=1}^n \frac{d^2 f(\xi_i)}{dx^2} \right)}_M$$

It is easy to see that if $f(x)$ is a continuous function, M (being the mean) must be bounded by the maximum and the minimum of $f(x)$ on the interval $[0, 1]$ and hence, there must exist some $\xi \in [0, 1]$ such that $M = d^2 f(\xi)/dx^2$. Hence we obtain the estimate

$$E_n = \frac{nh^3}{24} \frac{d^2 f(\xi)}{dx^2} = \frac{1}{24n^2} \frac{d^2 f(\xi)}{dx^2}$$

since $h = 1/n$. This error estimate tells us that the scheme is again exact for constants and linear functions on the domain (no higher order polynomials!) and, for a smooth function, the error decays **algebraically**.

4.3 The Exactness Criteria

Consider the Taylor series for $f(x)$ expanded about $x = 0$

$$f(x) = f(0) + \frac{\partial f(0)}{\partial x} x + \cdots + \frac{1}{l!} \frac{\partial^l f(0)}{\partial x^l} x^l + R_{l+1}$$

R_{l+1} is the **remainder**

$$R_{l+1} = \frac{1}{(l+1)!} \frac{\partial^{l+1} f(\tilde{x})}{\partial x^{l+1}} x^{l+1}$$

where $\tilde{x} \in [0, x]$

Note 3

Of all functions, why are we interested in integrating polynomials? The reason comes from the structure of Taylor's series expansion. The Taylor expansion of a function in a local neighborhood of a point (here this point is chosen as 0 without loss of generality) is nothing but a power series expansion. The higher

the order of polynomials that our scheme can integrate means a higher order of the remainder term in the expansion. The integral of the remainder over the domain is precisely the error in numerical integration.

The exactness condition requires

$$\int_0^1 p_l(x) dx = \int_0^1 (a_0 + a_1 x + a_2 x^2 + \dots + a_l x^l) dx = \sum_{i=1}^n w_i p_l(x_i)$$

for any set of $l + 1$ coefficients a_0, a_1, \dots, a_l

Equivalently

$$\int_0^1 a_0 dx + \int_0^1 a_1 x dx + \int_0^1 a_2 x^2 dx + \dots + \int_0^1 a_l x^l dx = \sum_{i=1}^n w_i p_l(x_i)$$

This slide needs little clarification. Our exactness criterion is

$$\int_0^1 p_l(x) dx = \int_0^1 (a_0 + a_1 x + a_2 x^2 + \dots + a_l x^l) dx = \sum_{i=1}^n w_i p_l(x_i)$$

which is the same as

$$a_0 \int_0^1 dx + a_1 \int_0^1 x dx + \dots + a_l \int_0^1 x^l dx = a_0 \sum_{i=1}^n w_i + a_1 \sum_{i=1}^n w_i x_i + \dots + a_l \sum_{i=1}^n w_i x_i^l$$

For this to be an identity for the $(l + 1)$ arbitrary coefficients a_i , we must have the $(l + 1)$ conditions

$$\sum_{i=1}^n w_i x_i^j = \int_0^1 x^j dx \quad \text{for } j = 0, 1, \dots, l$$

Using the Taylor series results, the exactness criteria, and the innate linearity of the quadrature scheme

$$\int_0^1 f(x) dx - \sum_{i=1}^n w_i f(x_i) = \underbrace{\frac{1}{(l+1)!} \frac{\partial^{l+1} f(\tilde{x})}{\partial x^{l+1}} \left(\int_0^1 x^{l+1} dx - \sum_{i=1}^n w_i x_i^{l+1} \right)}_{\text{Remainder}}$$

Exactness condition will be satisfied if and only if

$$\begin{aligned} \int_0^1 dx &= \sum_{i=1}^n w_i \cdot 1 \\ \int_0^1 x dx &= \sum_{i=1}^n w_i \cdot x_i \\ &\vdots \\ \int_0^1 x^l dx &= \sum_{i=1}^n w_i \cdot x_i^l \end{aligned}$$

Reorganizing exactness equations

$$\begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \\ \vdots & \vdots & \ddots & \vdots \\ x_1^l & x_2^l & \cdots & x_n^l \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} - \begin{bmatrix} 1 \\ \vdots \\ \vdots \\ \int_0^1 x^l dx \end{bmatrix} = 0$$

Nonlinear, since x_i 's and w_i 's are unknowns

What is a practical way of computing the evaluation points and weights? The system of equations is not easy to solve since x_i 's and w_i 's are unknowns.

5 Computing the Points & Weights

5.1 Newton's Method

Could use **Newton's Method**

$$F(y) = 0 \Rightarrow J_F(y^k) (y^{k+1} - y^k) = -F(y^k)$$

The nonlinear function for Newton is then

$$F \begin{pmatrix} w \\ x \end{pmatrix} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \\ \vdots & \vdots & \ddots & \vdots \\ x_1^l & x_2^l & \cdots & x_n^l \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} - \begin{bmatrix} 1 \\ \vdots \\ \vdots \\ \int_0^1 x^l dx \end{bmatrix} = 0$$

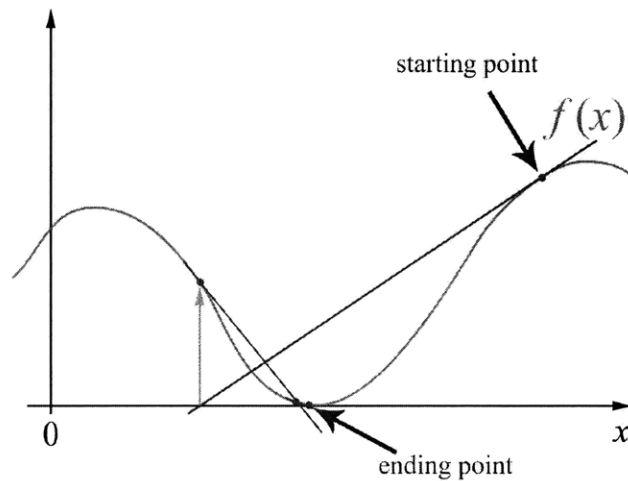
Note 4

Newton's method is an iterative technique for finding a value y such that $F(y) = 0$. The method is based on linearizing the problem about a guess at y , and then

updating the value of y by solving the linearized problem. In particular, the iterate y^{k+1} is determined from y^k by solving the linear system of equations

$$F(y^k) + J_F(y^k) (y^{k+1} - y^k) = 0$$

where $J_F(y^k)$ is the Jacobian (multidimensional derivative) of the nonlinear function $F(y)$. The iteration is continued until the updated y is sufficiently close to the exact solution, a criterion that can be difficult to verify. Newton's method does not always converge, a phenomenon that is more likely when $J_F(y)$ is nearly singular. For more about Newton's method, see the 6.336/16.910/2.096 course notes (available under open courseware).



Newton Method Jacobian reveals the problem

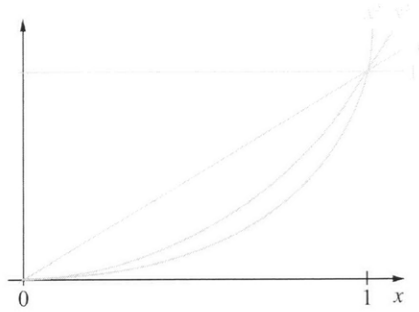
$$J_F \begin{pmatrix} w \\ x \end{pmatrix} = \begin{bmatrix} 1 & 1 & \dots & 1 & 0 & 0 & \dots & 0 \\ x_1 & x_2 & \dots & x_n & w_1 & w_2 & \dots & w_n \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ x_1^l & x_2^l & \dots & x_n^l & lw_1 x_1^{l-1} & \dots & \dots & lw_n x_n^{l-1} \end{bmatrix}$$

Columns become linearly dependent for high order

Looking at the Jacobian of the problem, we realize that the first n columns become increasingly linearly dependent for large l . This is bound to happen since we are looking into the space $\text{span}\{1, x, \dots, x^l\}$. This basis always becomes ill-conditioned with increasing l . The solution is to obtain a polynomial basis that is "normalized" in some sense so that it is properly conditioned.

5.2 Orthogonal Polynomials

5.2.1 Introduction



Exactness criteria will be satisfied if and only if

$$\left. \begin{array}{l} \int_0^1 c_0(x) dx = \sum_{i=1}^n w_i c_0(x_i) \\ \int_0^1 c_1(x) dx = \sum_{i=1}^n w_i c_1(x_i) \\ \vdots \\ \int_0^1 c_i(x) dx = \sum_{i=1}^n w_i c_i(x_i) \end{array} \right\} \begin{array}{l} \text{BUT} \\ \text{Each } c_i \text{ polynomial must} \\ \text{Contain an } x^i \text{ term} \\ \text{Be linearly independent} \end{array}$$

The only difference from the previous set of criteria is that these polynomials have better properties than the ones we chose before.

5.2.2 Orthogonality

For the normalized integral, two polynomials are said to be **orthogonal** if

$$\int_0^1 c_i(x)c_j(x)dx = 0 \quad \text{for } j \neq i$$

The above integral is often referred to as an inner product and ascribed the notation

$$(c_i, c_j) = \int_0^1 c_i(x)c_j(x)dx$$

The connection between polynomial inner products and vector inner products can be seen by sampling.

5.2.3 Exactness Criteria

Consider rewriting the exactness criteria

$$\begin{array}{cc}
 \int_0^1 c_0(x) dx = \sum_{i=1}^n w_i c_0(x_i) & \int_0^1 c_n(x) dx = \sum_{i=1}^n w_i c_n(x_i) \\
 \vdots & \vdots \\
 \underbrace{\int_0^1 c_{n-1}(x) dx = \sum_{i=1}^n w_i c_{n-1}(x_i)}_{\text{Low order terms}} & \underbrace{\int_0^1 c_{2n-1}(x) dx = \sum_{i=1}^n w_i c_{2n-1}(x_i)}_{\text{High Order Terms}}
 \end{array}$$

Recall that $l = 2n - 1$
 where $l = \text{degree of polynomial}$ & $n = \text{number of coefficients}$

Call the first $(n - 1)$ conditions the “lower order terms” and the last n conditions the “higher order terms.”

5.2.4 Higher Order Terms Contain Lower Order Terms

Write the higher order terms differently

$$\begin{array}{ccc}
 \int_0^1 c_n(x) dx = \sum_{i=1}^n w_i c_n(x_i) & \Rightarrow & \int_0^1 c_n(x) c_0(x) dx = \sum_{i=1}^n w_i c_n(x_i) c_0(x_i) \\
 & & \vdots \\
 \int_0^1 c_{2n-1}(x) dx = \sum_{i=1}^n w_i c_{2n-1}(x_i) & \Rightarrow & \int_0^1 c_n(x) c_{n-1}(x) dx = \sum_{i=1}^n w_i c_n(x_i) c_{n-1}(x_i)
 \end{array}$$

The products $c_n(x)c_j(x)$ are linearly independent.

In this slide we express the “higher order terms” as conditions involving “lower order terms.”

5.2.5 Using Orthogonality and Roots

Use orthogonal polynomials

$$\int_0^1 c_n(x) c_0(x) dx = \sum_{i=1}^n w_i c_n(x_i) c_0(x_i)$$

$$\int_0^1 c_n(x) c_{n-1}(x) dx = \sum_{i=1}^n w_i c_n(x_i) c_{n-1}(x_i)$$

Pick the x_i 's to be n roots of $c_n(x)$
 The higher order constraints are exactly satisfied!

This elegant step relies on polynomial orthogonality. If we choose the polynomial $c_n(x)$ such that it is orthogonal to all polynomials of inferior degree (i.e. $c_0(x), c_1(x), \dots, c_{n-1}(x)$) and the x_i 's are roots of this polynomial, then the higher order n conditions are automatically satisfied. Note that for this derivation we used polynomials which are orthogonal on the interval $[0, 1]$. Such polynomials are shifted and scaled versions of the classical Legendre polynomials, which are orthogonal on the interval $[-1, 1]$.

5.2.6 Satisfying Lower Order Constraints

An abbreviated exactness equation

$$\begin{bmatrix} c_0(x_1) & \cdots & c_0(x_n) \\ \vdots & \ddots & \vdots \\ c_{n-1}(x_1) & \cdots & c_{n-1}(x_n) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} = \begin{bmatrix} \int_0^1 c_0(x) dx \\ \vdots \\ \int_0^1 c_{n-1}(x) dx \end{bmatrix}$$

Now **linear**, x_i 's are known
 Rows are sampled orthogonal polynomials.

By using the roots of $c_n(x)$ for the x_i 's, the higher order constraints are automatically satisfied. Since the x_i 's are now known, only the weights are still unknown. The lower n constraints can be used to determine the weights, generating a linear system.

6 Gaussian Quadrature Algorithm

1. Construct $n + 1$ orthogonal polynomials

$$\int_0^1 c_i(x) c_j(x) dx = 0 \quad \text{for } j \neq i$$

2. Compute n **roots**, x_i , $i = 1, \dots, n$ of the n^{th} order orthogonal polynomial such that $c_n(x_i) = 0$
3. Solve a linear system for the **weights** w_i

4. Approximate the integral as a sum

$$\int_0^1 f(x)dx = \sum_{i=1}^n w_i f(x_i)$$

6.1 Example

Note 5

Example: Third Order Polynomial

Recall the example of the third-order polynomial, $f(x) = x - x^3$. We would like to determine w_1 , w_2 , x_1 , and x_2 so that the integration formula,

$$I_2\{f\} = w_1 f(x_1) + w_2 f(x_2)$$

gives the exact result.

Generalize this problem to any third-order polynomial of the form

$$f(x) = a_0 + a_1 x + a_2 x^2 + a_3 x^3.$$

Given the exactness condition, the integral can be written:

$$\int_0^1 (a_0 + a_1 x + a_2 x^2 + a_3 x^3) dx = w_1 (a_0 + a_1 x_1 + a_2 x_1^2 + a_3 x_1^3) + w_2 (a_0 + a_1 x_2 + a_2 x_2^2 + a_3 x_2^3)$$

Gather the like terms:

$$a_0 : w_1 + w_2 = \int_0^1 dx = 1$$

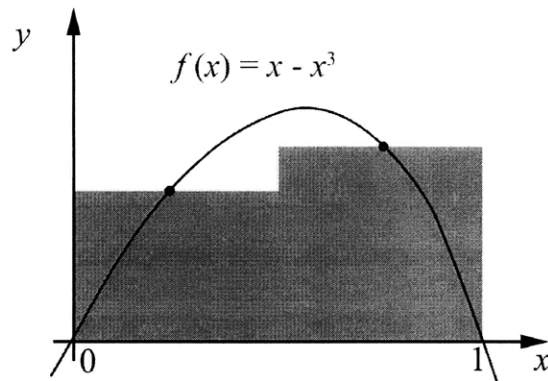
$$a_1 : w_1 x_1 + w_2 x_2 = \int_0^1 x dx = \frac{1}{2}$$

$$a_2 : w_1 x_1^2 + w_2 x_2^2 = \int_0^1 x^2 dx = \frac{1}{3}$$

$$a_3 : w_1 x_1^3 + w_2 x_2^3 = \int_0^1 x^3 dx = \frac{1}{4}$$

There are four equations above and four unknowns, so the system can be easily solved to give, $w_1 = \frac{1}{2}$ $w_2 = \frac{1}{2}$ $x_1 = \frac{3-\sqrt{3}}{6}$ $x_2 = \frac{3+\sqrt{3}}{6}$.

Plugging in these values, for the equation above, gives the exact solution of 0.25.



6.2 Summary

This slide summarizes the technique of finding weights and integration points for Gauss quadrature.

6.2.1 Accuracy Result

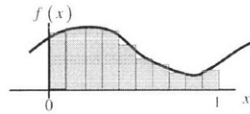
$$\int_0^1 f(x)dx \simeq \sum_{i=1}^n w_i f(x_i)$$

Key properties of the method

- An n -point Gauss quadrature rule is **exact** for polynomials of order $2n - 1$
- Error is proportional to $\frac{1}{(2n)!}$ (like $\frac{1}{n^n}$)

6.2.2 Simple vs. Gauss Quadrature

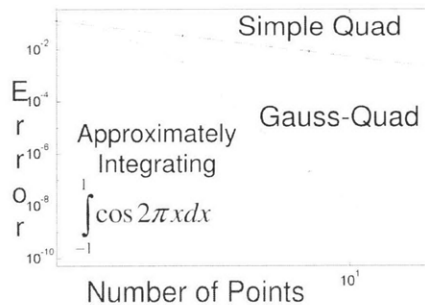
$$\int_0^1 f(x)dx \simeq \sum_{i=1}^n \frac{1}{n} f\left(\frac{i - \frac{1}{2}}{n}\right)$$



Key property of the method

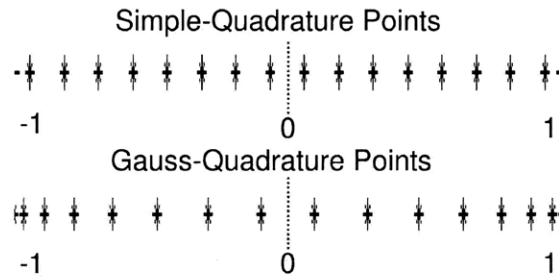
- Error is proportional to $\frac{1}{n^2}$

▷ **Exercise 3** Do you see that the simple quadrature scheme is a special case of Gauss quadrature? ■



Notice that for a smooth function $f(x) = \cos(2\pi x)$, which is infinitely differentiable, Gauss quadrature far outperforms the simple quadrature scheme

6.2.3 Evaluation Point Placement



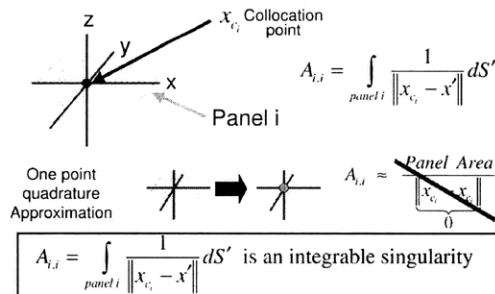
Notice the clustering at interval ends

In the Gauss quadrature scheme the evaluation points are roots of Legendre polynomials which are clustered at the ends of the interval.

7 The Singular Kernel Problem

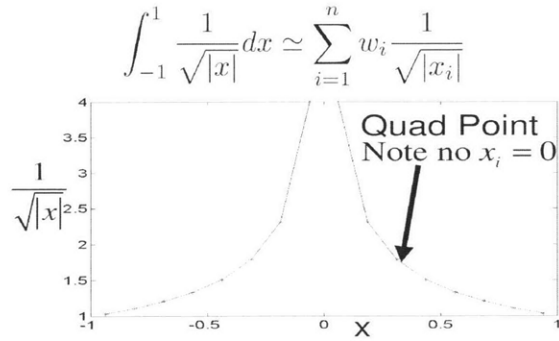
7.1 Calculating the “Self-Term”

Now lets go back to our problem of solving Laplace’s equation on a 3D domain using boundary integral representation. We realize that we now have some sophisticated tools to handle integrals of functions that are smooth. But what about the integral on the panel where the centroid x_{c_i} is located? The Green’s function blows up at the centroid. However, the function is integrable because the integrand blows up at a rate that is slower than the rate at which the surface measure goes to zero in the vicinity of the singularity. So we know that the integral exists and is finite, but is Gauss quadrature capable of performing well in the presence of this singularity?

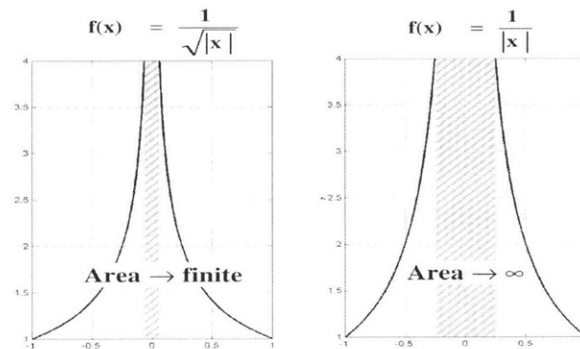


7.2 Symmetric 1-D Example

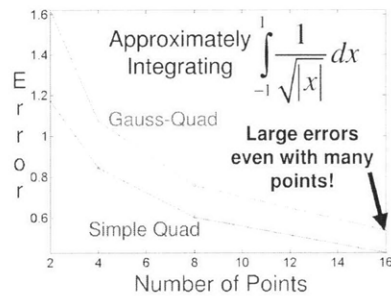
In 1D we look at a function $f(x) = \frac{1}{\sqrt{|x|}}$ which is integrable on $[-1, 1]$ but has a singularity at $x = 0$.



7.2.1 Integrable and Nonintegrable Singularities



7.2.2 Comparing Quadrature Schemes

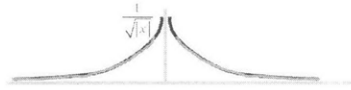


We observe that Gauss quadrature is not very good at integrating this function. The convergence is rather poor. As a matter of fact, it is more inaccurate than the simple quadrature scheme. In the next few slides we present several techniques of handling integrals with singularities (which are integrable, of course)

- Subinterval (adaptive) quadrature
- Change of variables of integration
- Singular (Gaussian) quadrature

7.3 Improved Techniques

7.3.1 Subinterval (Adaptive) Quadrature



Subdivide the integration interval

$$\int_{-1}^1 \frac{1}{\sqrt{|x|}} dx = \int_{-1}^{-0.1} \frac{1}{\sqrt{|x|}} dx + \int_{-0.1}^0 \frac{1}{\sqrt{|x|}} dx + \int_0^{0.1} \frac{1}{\sqrt{|x|}} dx + \int_{0.1}^1 \frac{1}{\sqrt{|x|}} dx$$

Use Gauss quadrature in each subinterval

Polynomials fit subintervals better

Expensive if many subintervals used.

7.3.2 Change of Variables - for Simple Cases

Change variables to eliminate singularity

$$y^2 = x \quad \Rightarrow \quad 2y dy = dx$$

$$\int_{-1}^1 \frac{1}{\sqrt{|x|}} dx = 2 \int_0^1 \frac{1}{\sqrt{|y^2|}} 2y dy = 2 \int_0^1 2 dy$$

Apply Gauss quadrature on desingularized integrand.

7.3.3 Singular Quadrature - Complicated Cases

Basic Concept

Integrand has known singularity $s(x)$

$$\int_{-1}^1 f(x)s(x) dx \text{ where } f(x) \text{ is smooth}$$

Develop a quadrature formula exact for

$$\int_{-1}^1 p_l(x)s(x) dx \text{ where } p_l(x) \text{ is polynomial of order } l$$

Calculate weights like Gauss quadrature

It is possible to generate Gaussian quadrature schemes of the form

$$\int_{-1}^1 s(x)f(x)dx = \sum_{i=1}^n w_i f(x_i)$$

for functions which have a known singularity $s(x) > 0$. The quadrature formula needs to be exact when $f(x)$ is a polynomial of order at most l . Not surprisingly, it turns out that the integration points are the n roots of a polynomial $c_n(x)$ of degree $n = (l + 1)/2$ which is orthogonal to all polynomials of inferior degree with respect to the weight $s(x)$, i.e.

$$\int_{-1}^1 s(x)c_n(x)c_j(x) = 0 \quad \text{for } j = 0, 1, \dots, (n - 1).$$

An example is the singular integral

$$I = \int_{-1}^1 \frac{f(x)}{\sqrt{1-x^2}}$$

Here, $s(x) = 1/\sqrt{1-x^2}$ and the corresponding orthogonal polynomials turn out to be the Chebyshev polynomials. The integration points are given in closed form by

$$x_i = \cos\left(\pi \frac{2i-1}{2n}\right)$$

and the corresponding weights are $w_i = \pi/n$.

Singular Quadrature Weights

$$\begin{bmatrix} c_0(x_1) & \cdots & c_0(x_n) \\ \vdots & \ddots & \vdots \\ c_{n-1}(x_1) & \cdots & c_{n-1}(x_n) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} = \begin{bmatrix} \int_{-1}^1 c_0(x)s(x)dx \\ \vdots \\ \int_{-1}^1 c_{n-1}(x)s(x)dx \end{bmatrix}$$

Need (analytic) formulas for integrals of $c(x)s(x)$

The lower order constraints can be used to compute the integration weights.

8 Summary

Easy technique for computing integrals

Piecewise constant approach

Gaussian quadrature

Faster convergence
Essential role of orthogonal polynomials
Techniques for singular kernels
Adaptation and Variable Transformation
Singular quadrature
What about multiple dimensions?

1.3 First and Second Kind Equations

Numerical Methods for PDEs

Integral Equation Methods, Lecture 3
First and Second Kind Equations

Notes by L. Proctor, S. De and J. White

December 1, 2008

1 Outline

Convergence issues in 1D

- First and second kind integral equations
- Develop some intuition about the difficulties

Convergence for second kind equations

- Consistency and stability issues

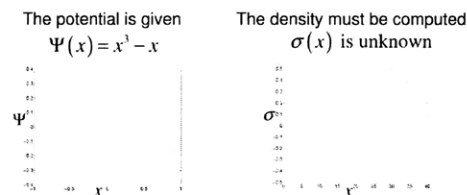
Nystrom Methods

- High order convergence

2 Example Problems in 1D

2.1 First Kind Equation

$$\Psi(x) = \int_{-1}^1 |x - x'| \sigma(x') dS' \quad x \in [-1, 1]$$



In the next several slides we will investigate the convergence properties of integral equation discretization methods. How these methods converge depends on what kind of integral equation is being solved. Examining this issue will introduce one of the subtle points about integral equations.

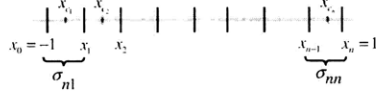
To begin, consider the example one-dimensional first-kind integral equation on the above slide. For this equation, we assume that the potential, $\Phi(x)$, is known and that the charge density $\sigma(x)$ is unknown. Here, x is in the interval $[-1, 1]$, and the integration is over that same interval. Note that for this example, the Green's function is given by $G(x, x') = |x - x'|$.

In the left plot below the equation, an example potential, $x^3 - x$, is given and plotted as a function of x . On the right is a plot of charge density as a function of x which is a possible solution to the integral equation. As we will see shortly, the solution for the charge density is not so easy to find.

2.1.1 Collocation Discretization

$$\Psi(x) = \int_{-1}^1 |x - x'| \sigma(x') dS' \quad x \in [-1, 1] \quad (1)$$

Piecewise-Constant Centroid-Collocation



$$\Psi(x_{c_i}) = \sum_{j=1}^n \sigma_{nj} \int_{x_{j-1}}^{x_j} |x_{c_i} - x'| dS' \quad (2)$$

To compute the numerical solution to this one-dimensional problem, consider solving integral equation (1) using a piecewise-constant collocation scheme. In such a scheme, we first select $n + 1$ points on the interval, in this case $[-1, 1]$. We denote those points as $\{x_0, x_1, \dots, x_n\}$, as shown in the above figure. For this example, $x_0 = -1$ and $x_n = 1$. Over the subintervals define a set of basis functions, $\{\varphi_1(x), \varphi_2(x), \dots, \varphi_n(x)\}$, where

$$\varphi_i(x) = \left\{ \begin{array}{ll} 1 & x \in [x_{i-1}, x_i] \\ 0 & \text{otherwise} \end{array} \right\}. \quad (3)$$

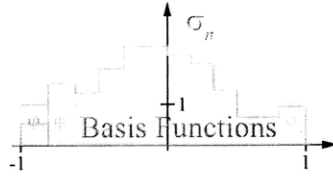
The charge density σ can then be represented approximately as

$$\sigma(x) \approx \sigma_n(x) \equiv \sum_{i=1}^n \sigma_{ni} \varphi_i(x), \quad (4)$$

where σ_{ni} is the weight associated with the i^{th} basis function. It may seem odd that we used the same letter to represent the density and the basis function weights, but there is a reason. The above basis set is such that only one basis function is nonzero for a given x , and basis functions only take on the value zero or one. Therefore, σ_{ni} will be equal to the approximate charge density when $x \in [x_{i-1}, x_i]$.

Charge Density Representation

$$\sigma_n(x) \equiv \sum_{i=1}^n \sigma_{ni} \varphi_i(x)$$



Plugging the basis function representation of the charge density, equation (4), into equation (1) yields

$$\Psi(x) = \int_{-1}^1 |x - x'| \sum_{i=1}^n \sigma_{ni} \varphi_i(x') dS',$$

which can be simplified by exploiting equation (3). Recall that the residual, $R(x)$, is defined as how well the weighted combination of basis functions satisfies the integral equation. In this centroid collocation case,

$$R(x) = \Psi(x) - \sum_{j=1}^n \sigma_{nj} \int_{x_{j-1}}^{x_j} |x - x'| dS'.$$

If collocation is used to solve this equation, then $R(x_{c_i}) = 0$ for all i , where x_{c_i} is the i^{th} collocation point. The collocation points shown in the figure are the subinterval center points, $x_{c_i} = \frac{1}{2}(x_{i-1} + x_i)$. Note that there are other choices for collocation points, such as $x_{c_i} = x_i$.

Using the fact that $R(x_{c_i}) = 0$ leads to

$$R(x_{c_i}) = \Psi(x_{c_i}) - \sum_{j=1}^n \sigma_{nj} \int_{x_{j-1}}^{x_j} |x_{c_i} - x'| dS' = 0$$

which can be reorganized into equation (2).

The Matrix

We can now generate a system of equations that can be used to solve for the σ_{ni} 's, the piecewise constant charge densities for each of the subintervals.

One column for each density value

$$\begin{bmatrix} \int_{x_0}^{x_1} |x_{c_1} - x'| dS' & \cdots & \int_{x_{n-1}}^{x_n} |x_{c_1} - x'| dS' \\ \vdots & \ddots & \vdots \\ \int_{x_0}^{x_1} |x_{c_n} - x'| dS' & \cdots & \int_{x_{n-1}}^{x_n} |x_{c_n} - x'| dS' \end{bmatrix} \begin{bmatrix} \sigma_{n1} \\ \vdots \\ \sigma_{nn} \end{bmatrix} = \begin{bmatrix} \Psi(x_{c_1}) \\ \vdots \\ \Psi(x_{c_n}) \end{bmatrix}$$

One row for each collocation point

The right-hand side of this system of equations is a vector of known potentials at the interval centers (the collocation points). The i^{th} row of the matrix corresponds to unfolding the sum in the collocation equation

$$\Psi(x_{c_i}) = \sum_{j=1}^n \sigma_{nj} \int_{x_{j-1}}^{x_j} |x_{c_i} - x'| dS',$$

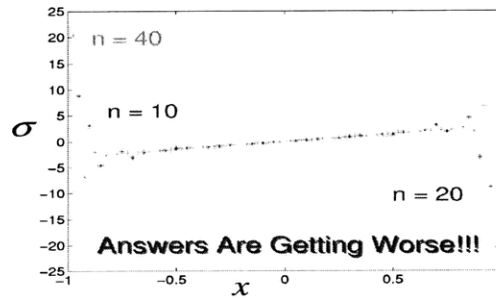
and the entries in the j^{th} column correspond to how much the charge on the j^{th} interval contributes to the i^{th} potential.

Note that the matrix is square and dense.

▷ **Exercise 1** Is the above matrix symmetric? If we used $x_{c_i} = x_i$, would the matrix still be symmetric? ■

2.1.2 Numerical Results with Increasing n

One usually believes that a discretization scheme should produce progressively more accurate answers as the discretization is refined. In this case, as we divide the interval into progressively finer subintervals, one might expect that the piecewise constant representation of the charge density given by $\sigma_n(x) \approx \sum_{i=1}^n \sigma_{ni} \varphi_i(x)$ would become more accurate as n increases.



As is clear from the above plot, the piecewise-constant centroid-collocation discretization of (1) is not converging. In the plot, which is hard to decipher without looking at a color version, the σ_{ni} 's produced using $n = 10$, $n = 20$ and $n = 40$ subintervals are shown. For each discretization, a point is plotted at σ_{ni} , x_i for $i = 1, \dots, n$, so there are ten points plotted for the coarsest discretization and forty points plotted for the finest discretization, but all sets of points span the interval $x \in [-1, 1]$.

What is clear from comparing the blue points ($n = 10$) to the red points ($n = 20$) and to the green points ($n = 40$), is that the charge density seems to be approaching infinity as the discretization is refined. The results are certainly not converging.

Why is this happening? Is the numerical technique at fault, or is the integral equation a problem?

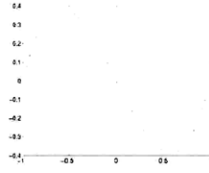
2.2 Second Kind Equation

We are going to postpone examining what went wrong with the discretization the first-kind integral equation, and instead examine a Second Kind integral equation example. As in the first-kind case, we are assuming the potential, $\Psi(x)$, is known and that the charge density $\sigma(x)$ is unknown. Also, x is in the interval $[-1, 1]$, and the integration is over that same interval. Once again, the Green's function is given by $G(x, x') = |x - x'|$.

$$\Psi(x) = \sigma(x) + \int_{-1}^1 |x - x'| \sigma(x') dx' \quad x \in [-1, 1]$$

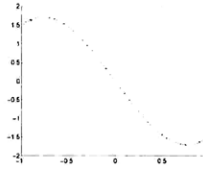
The potential is given

$$\Psi(x) = x^3 - x$$



The density must be computed

$$\sigma(x) \text{ is unknown}$$



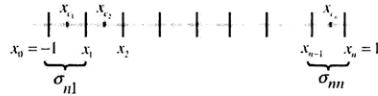
The above equation is second-kind instead of first-kind because the unknown charge density appears both inside *and* outside of the integral. In the first-kind equation, the density appeared only inside the integral. This seemingly small difference has enormous numerical ramifications.

In the left plot above, an example given potential, $\Psi(x) = x^3 - x$ is plotted as a function of x . On the right is a plot of a charge density as a function of x which satisfies this second kind integral equation. As we will see below, this equation is easily solved numerically.

2.2.1 Collocation Discretization

$$\Psi(x) = \sigma(x) + \int_{-1}^1 |x - x'| \sigma(x') dS' \quad x \in [-1, 1] \quad (5)$$

Centroid Collocated Piecewise Constant Scheme



$$\Psi(x_{c_i}) = \sigma_{ni} + \sum_{j=1}^n \sigma_{nj} \int_{x_{j-1}}^{x_j} |x_{c_i} - x'| dS' \quad (6)$$

To compute the numerical solution to the one-dimensional second-kind equation (5), once again consider using a piecewise-constant collocation scheme. Again, we select $n + 1$ points on the interval and denote those points as $\{x_0, x_1, \dots, x_n\}$, as shown in the figure above. For this example, $x_0 = -1$ and $x_n = 1$. The corresponding basis functions, $\{\varphi_1(x), \varphi_2(x), \dots, \varphi_n(x)\}$, are the same as in equation (3):

$$\varphi_i(x) = \left\{ \begin{array}{ll} 1 & x \in [x_{i-1}, x_i] \\ 0 & \text{otherwise} \end{array} \right\}. \quad (7)$$

The charge density σ is approximately represented by

$$\sigma(x) \simeq \sigma_n(x) \equiv \sum_{i=1}^n \sigma_{ni} \varphi_i(x), \quad (8)$$

where σ_{ni} is the weight associated with the i^{th} basis function. Plugging the basis function representation of the charge density, equation (8) into the second kind integral equation (5) gives:

$$\Psi(x) = \sum_{j=1}^n \sigma_{nj} \varphi_j(x) + \int_{-1}^1 |x - x'| \sum_{i=1}^n \sigma_{ni} \varphi_i(x') dS',$$

which can be simplified by exploiting the specific basis functions, equation (7) to

$$\Psi(x) = \sum_{j=1}^n \sigma_{nj} \varphi_j(x) + \sum_{j=1}^n \sigma_{nj} \int_{x_{j-1}}^{x_j} |x - x'| dS'. \quad (9)$$

As shown in the above figure, the collocation points are the subinterval center points, $x_{c_i} = \frac{1}{2}(x_{i-1} + x_i)$. When collocation is used, equation (9) must be satisfied exactly at the collocation points and therefore

$$\Psi(x_{c_i}) = \sum_{i=1}^n \sigma_{ni} \varphi_j(x_{c_i}) + \sum_{j=1}^n \sigma_{nj} \int_{x_{j-1}}^{x_j} |x_{c_i} - x'| dS'. \quad (10)$$

Note that $\varphi_j(x_{c_i}) = 0$ when $i \neq j$, and $\varphi_i(x_{c_i}) = 1$. Using this fact yields equation (6).

The Matrix

Just as in the discretized first-kind equation, we generate a system of equations that can be used to solve for the σ_{ni} 's, the piecewise constant charge densities for each of the subintervals.

$$\begin{bmatrix} \int_{x_0}^{x_1} |x_{c_1} - x'| dS' & \cdots & \int_{x_{n-1}}^{x_n} |x_{c_1} - x'| dS' \\ \vdots & \ddots & \vdots \\ \int_{x_0}^{x_1} |x_{c_n} - x'| dS' & \cdots & \int_{x_{n-1}}^{x_n} |x_{c_n} - x'| dS' \end{bmatrix} \begin{bmatrix} \sigma_{n1} \\ \vdots \\ \sigma_{nn} \end{bmatrix} = \begin{bmatrix} \Psi(x_{c_1}) \\ \vdots \\ \Psi(x_{c_n}) \end{bmatrix}$$

The right-hand side of this system of equations is a vector of known potentials at interval centers (the collocation points). The i^{th} row of the matrix corresponds to unfolding the sum in the collocation equation

$$\Psi(x_{c_i}) = \sigma_{ni} + \sum_{j=1}^n \sigma_{nj} \int_{x_{j-1}}^{x_j} |x_{c_i} - x'| dS'$$

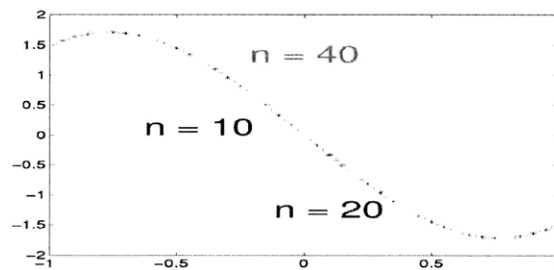
and the entries in the j^{th} column corresponds to how much the charge on the j^{th} interval contributes to the i^{th} potential.

The major difference between the matrix in this discretized second-kind example and the first-kind example is circled on the above slide. There is an additional one on the diagonal of the discretized second-kind equation that did not appear in the first-kind equation. In other words,

$$A_{second\ kind} = I + A_{first\ kind}.$$

2.2.2 Numerical Results with Increasing n

Unlike the results from discretizing the first kind equation, progressively refining the discretization of the second kind equation produces more accurate answers.



Answers Are Improving!!!

Once again, the plot is a little hard to decipher without looking at the color version. It shows the σ_{ni} 's produced using $n = 10$, $n = 20$ and $n = 40$ subintervals. For each discretization, a point is plotted at σ_{ni}, x_i for $i = 1, \dots, n$, so there are ten points plotted for the coarsest discretization and forty points plotted for the finest discretization, but all sets of points span the interval $x \in [-1, 1]$.

What is clear from comparing the blue points ($n = 10$) to the red points ($n = 20$) and to the green points ($n = 40$), is that the charge density seems to be approaching a smooth solution.

What is the essential difference between first and second kind equations?

Is it some aspect of the numerical technique or are these two equations really that different? In the next slides, we will try to answer this question.

2.3 Difficulty with the First Kind Equation

We will make use of operator-function notation for much of the next sections, both for clarity and brevity. For example, the charge density, σ , and the potential, Ψ , are functions of the independent variable x . When we mean the function, we use just use the function name, such as σ or Ψ . When we give an explicit formula in terms of x for the function, or are denoting the function's value for a particular \hat{x} , we use follow the function name by the value of the

independent variable in parentheses. For example, $\sigma(\hat{x})$ is the value of function σ when evaluated at $x = \hat{x}$.

The operator-function notation is a little less obvious in the case of operators that map functions to functions, like the integral operator. The integral operator takes a function, in this case σ , and produces another function that we might refer to as the potential. If we denote the integral operator from (5) as K , then $K\sigma$ is a function. If we wish to evaluate the function generated by applying K to σ at some \hat{x} , then we write $(K\sigma)(x)$. Note that $K\sigma(x)$ would *NOT BE CORRECT*. The operator K takes functions and $\sigma(x)$ is a value.

2.3.1 Singular Integral Operator

Denote the integral operator as K

$$(K\sigma)(x) \equiv \int_{-1}^1 |x - x'| \sigma(x') dS' = \Psi(x) \Rightarrow K\sigma = \Psi \quad (11)$$

The integral operator is **singular** : K has a null space

$$\sigma_0(x) = 0, x \neq 0, \sigma_0(0) = 1$$

$$(K\sigma_0)(x) = \int_{-1}^1 |x - x'| \sigma_0(x') dS' = 0 \text{ for all } x \quad (12)$$

$$\Rightarrow K\sigma_0 = 0 \quad (13)$$

$$\text{If } K\sigma^a = \Psi \text{ then } K(\sigma^a + \sigma_0) = \Psi$$

In equation (11), we introduce the abstract notion that

$$\int_{-1}^1 |x - x'| \sigma(x') dS'$$

is an operator on the function σ , which we denote with the symbol K . As shown on the top of the slide, this notation makes writing the integral equation look just like writing a matrix equation.

The key problem is that the operator K is singular. And if

$$K\sigma = \Psi$$

were a matrix equation with a singular K , one would not be surprised to discover the system of equations is hard, or impossible, to solve.

We will not try, in this lecture, to be formal about the concept of a singular operator. To do so, we would necessarily be examining details about certain types of function spaces. Instead, we will try to develop some intuition. In particular, we will draw an analogy to matrices and note that if an operator is

singular, it must have a null space.

To see that K does have a null space, consider the spike function $\sigma_0(x)$ depicted in equation (13). This spike function is one at $x = 0$ and zero otherwise. Note, this function is *not* an impulse function. Unlike the impulse function, the spike's value at $x = 0$ is finite and the area under its curve is obviously zero.

As noted in equation (13), $K\sigma_0 = 0$. To see this consider that since σ_0 is nonzero only at $x = 0$, and therefore

$$\int_{-1}^1 |x - x'| \sigma_0(x') dS' = |x| \int_{-1}^1 \sigma_0(x') dS'.$$

Since $\int_{-1}^1 \sigma_0(x') dS' = 0$, as the area under σ_0 's curve is zero, then $K\sigma_0 = 0$.

The statement "If $K\sigma^a = \Psi$ then $K(\sigma^a + \sigma_0) = \Psi$ " says that if K has a null space, and there exists a solution, then there exist infinitely many solutions.

One last comment should be made. The spike function we generated is not unique. Simply shifting the nonzero point would generate an infinite number of spike functions which would all be in the null space of K . That is, K has an incredibly rich null space.

2.3.2 Eigenvalues

Difficulty from the Matrix

Collocation generates a discrete form of K

$$K\sigma = \Psi \quad \rightarrow \quad \underline{K}_n \vec{\sigma}_n = \vec{\Psi}_n$$

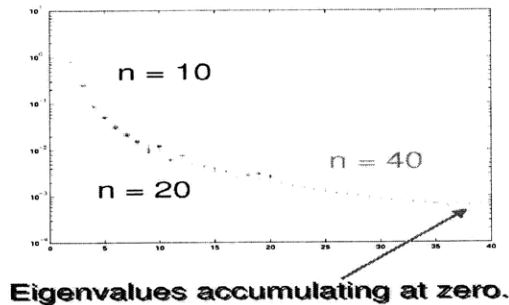
$$\underbrace{\begin{bmatrix} \int_{x_0}^{x_1} |x_1 - x'| dS' & \dots & \int_{x_{n-1}}^{x_n} |x_1 - x'| dS' \\ \vdots & \ddots & \vdots \\ \int_{x_0}^{x_n} |x_n - x'| dS' & \dots & \int_{x_{n-1}}^{x_n} |x_n - x'| dS' \end{bmatrix}}_{\underline{K}_n} \begin{bmatrix} \sigma_{n1} \\ \vdots \\ \sigma_{nn} \end{bmatrix} = \begin{bmatrix} \Psi(x_1) \\ \vdots \\ \Psi(x_n) \end{bmatrix}$$

K : functions to functions, \underline{K}_n :: vectors to vectors!

As shown above, discretizing the integral equation by combining a piecewise constant charge density representation with collocation at subinterval centers results in a system of equations which relates the subinterval σ_i 's to the collocation point potentials. From this perspective, the matrix above can be thought of as a discrete representation of the operator K . We denote the matrix with \underline{K}_n with an underline to indicate that it is a matrix, and was generated using a discretization with n basis functions. We also denoted the vectors $\vec{\sigma}_n$ and $\vec{\Psi}_n$ with arrows to avoid confusing vectors with functions. Later, we will need an operator version of the discretized representation of the operator K , but for the moment, the matrix is sufficient.

Numerical Results with Increasing n

If the operator K is singular, one might expect to see that reflected in the eigenvalues of a matrix generated by discretizing K . In particular, one would expect the matrix to have eigenvalues that are near zero.

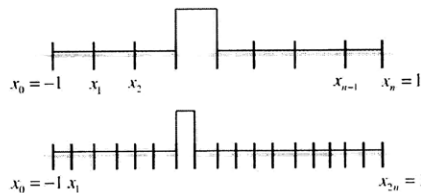


In the figure above, the eigenvalues of matrices generated by discretizing K for the 1-D problem are plotted. Discretizing using 10 subintervals generates a matrix with 10 eigenvalues plotted in blue. The blue eigenvalue closest to zero is ≈ 0.01 . As the discretization is refined to 20 subintervals, the minimum eigenvalue (plotted in red) drops to ≈ 0.003 , and with 40 subintervals the minimum eigenvalue (plotted in green) drops to ≈ 0.0009 . Examining this data suggests that as the discretization is refined, the generated matrix more accurately reflects the operator K , and therefore the matrix is becoming closer to being singular.

As the discretization is refined, the matrix is larger and has more eigenvalues. Notice that as the discretization is refined from $n = 10$ to $n = 20$ to $n = 40$, all the additional eigenvalues are closer to zero.

Intuition About Eigenvalues

As the discretization is refined, $\sigma_0(x)$ becomes better approximated



As the discretization is refined, K 's null space can be more accurately represented.

As an alternative view of why refining the discretization for the first kind equation produces a matrix with more and more smaller eigenvalues, consider the figures above. In the top plot, one of the basis functions is plotted for a coarse discretization. In the bottom plot, one of the basis functions is plotted for a finer discretization. Notice that as the discretization is refined, these basis functions look progressively more like the spike function mentioned previously. And since the spike function is in the null space of K , one would expect that

finer discretizations would generate “spikier” basis functions whose associated eigenvalues would be near zero.

2.4 Second Kind Equation Has Fewer Problems

Second Kind equation

$$\begin{aligned} ((I + K)\sigma)(x) &\equiv \sigma(x) + \int_{-1}^1 |x - x'| \sigma(x') dS' = \Psi(x) \\ \Rightarrow (I + K)\sigma &= \Psi \end{aligned} \tag{14}$$

$$(I + K)(\sigma_0 + \sigma) \neq (I + K)\sigma$$

$$\begin{aligned} (I+K) &\left(\begin{array}{c} \text{Graph of } \sigma(x) \text{ (smooth curve)} \\ \text{Graph of } \sigma_0(x) \text{ (spike at } x=0) \end{array} \right) \\ &= \text{Graph of } (\sigma + \sigma_0) \text{ (smooth curve with spike)} \neq \Psi \end{aligned}$$

$\sigma_0(x)=0, x \neq 0, \sigma_0(0)=1$

As shown in equation (14), the abstract operator for the second-kind equation is denoted by $I + K$, where I here is just the identity operator and K is the integral operator.

To see why the spike function, σ_0 , is not in the null space of the operator $I + K$, or equivalently that

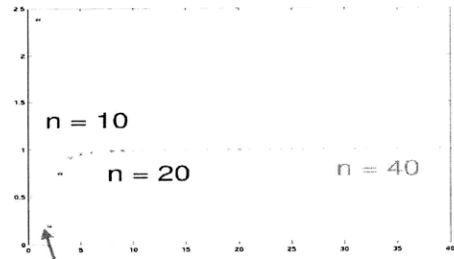
$$(I + K)(\sigma_0 + \sigma) \neq (I + K)(\sigma)$$

consider the figures above. If a spike is added to a smooth σ , the $(I + K)$ operator will preserve the spike. Another way to see this is to consider that since σ_0 is in the null space of K ,

$$(I + K)\sigma_0 = (I)\sigma_0 + K\sigma_0 = \sigma_0 \neq 0.$$

2.4.1 Eigenvalues

Numerical Results with Increasing n



Eigenvalues do not get closer to zero.

As we noted before, the matrix associated with discretizing the operator $I + K$ is identical to the sum of the identity matrix and the matrix associated with discretizing K alone. In the plot above, we once again present the eigenvalues generated by discretizing the 1-D example problem. Discretizing using 10 subintervals generates a matrix with 10 eigenvalues plotted in blue. The blue eigenvalue closest to zero is ≈ 0.2 . As the discretization is refined to 20 subintervals, the minimum eigenvalue (plotted in red) is still ≈ 0.2 , and with 40 subintervals the minimum eigenvalue (plotted in green) is still ≈ 0.2 . Examining this data suggests that as the discretization is refined, and the generated matrix more accurately reflects the operator $I + K$, the matrix is not becoming more singular. In fact, the eigenvalues are accumulating near one, an unsurprising result given that the eigenvalues of the discretized K operator were accumulating at zero.

▷ **Exercise 2** Estimate how many iterations will be needed for a Krylov-subspace based algorithm to converge for the 1-D discretized second-kind example. Will the number of iterations increase as the discretization is refined? ■

▷ **Exercise 3** Suppose the integral equation were changed to

$$\Psi(x) = \sigma(x) + \frac{1}{\lambda} \int_{-1}^1 |x - x'| \sigma(x') dS'.$$

For what value of λ would the solution no longer be unique. (you can answer this just by looking at the eigenplot above). ■

As the above exercise makes clear, a second-kind integral equation does not always have a unique solution. However, a first-kind equation almost never has a unique solution, the exception being when the Green's function is singular, as we will investigate in a subsequent lecture.

3 Theory of 2nd Kind Equations

The convergence theory for discretization methods applied to second-kind integral equations has an elegant simplicity, but only when examined using a care-

fully chosen abstraction. The theory is also surprisingly practical; the insights gained can be used to construct very high-order discretization schemes.

3.1 Comparison problem

General Second kind integral equation

$$\Psi(x) = \sigma(x) + \int G(x, x')\sigma(x')dx' \Rightarrow \Psi = (I + K)\sigma \quad (15)$$

Discrete matrix equivalent

$$\vec{\Psi}_n = (I + \underline{K}_n) \vec{\sigma}_n \quad (16)$$

How to compare function σ to vector $\vec{\sigma}_n$?

How to compare operator K to matrix \underline{K}_n ?

One approach to overcoming the comparison problems is to construct representations of the discretization that are functions and operators on functions. The most obvious approach to generating the functions associated with a discretization is by interpolation, but generating the operators associated with the discretization is a little more subtle.

3.1.1 Operator-Function Notation

General Second kind integral equation

$$\Psi(x) = \sigma(x) + \int G(x, x')\sigma(x')dx' \Rightarrow \Psi = (I + K)\sigma \quad (17)$$

Discretized operator-function equivalent

$$\Psi_n = (I + K_n) \sigma_n \quad (18)$$

σ_n, Ψ_n are functions of x (e.g. by interpolation)

K_n maps functions to functions like K (How constructed)?

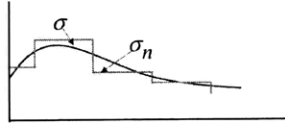
3.1.2 Orthogonal Galerkin

Representation $\sigma_n(x) = \sum_{i=1}^n \sigma_{ni} \varphi_i(x) \quad \int \varphi_i(x) \varphi_j(x) = \delta(i - j)$

Projection $\sigma_n = (P\sigma)(x)$

$$(P\sigma)(x) \equiv \sum_{i=1}^n \overbrace{\left(\int \sigma(x) \varphi_i(x) dx \right)}^{\sigma_{ni}} \varphi_i(x)$$

$$(P\sigma)(x) = \sum_{i=1}^n \sigma_{ni} \varphi_i(x)$$



If the density *sigma* is to be approximated by a weighted combination of n orthogonal basis functions, then the functional representation σ_n associated with the vector of weights $\vec{\sigma}_n$ is given. If the basis functions are orthonormal, then the σ_n associated with an arbitrary σ can be constructed by simple projection. It is worth noting that the projection operator, denoted P , has no effect on σ_n . That is $P\sigma_n = \sigma_n$.

▷ **Exercise 4** Why does the formula for the projection operator above require orthogonality of the basis functions? ■

If a Galerkin method is used to discretized the integral equation, then the associated operator is easy to construct, as shown below.

3.1.3 Ortho Galerkin Operator

$$\begin{aligned} (K\sigma_n)(x) &= (KP\sigma)(x) = \sum_{i=1}^n \sigma_{ni} \int G(x, x') \varphi_i(x') dx' \\ (PKP\sigma)(x) &= \sum_{j=1}^n \left(\int \varphi_j(x) KP\sigma(x) dx \right) \varphi_j(x) \\ &= \sum_{j=1}^n \left(\sum_{i=1}^n \sigma_{ni} \int \int \varphi_j(x) G(x, x') \varphi_i(x') dx dx' \right) \varphi_j(x) \\ P(I + KP)\sigma_n &= (I + PKP)\sigma_n = P\Psi \\ (I + K_n)\sigma_n &= \Psi_n \end{aligned}$$

The last equation on the above slide contains a subtle point. $P(I+KP)\sigma_n$ really equals $(P + PKP)\sigma_n$. However, P is equivalent to the identity operator when applied to σ_n , as projecting σ_n reproduces σ_n . So, we are free to conveniently chose to interpret $(P + PKP)$ as $(I + PKP)$ as a difference appears only when applying the operator to functions that are not weighted combinations of the n Galerkin basis functions.

For second-kind integral equations, one can prove a convergence theory for almost any reasonable discretization scheme, assuming that the equation has a

unique solution. As noted above, second-kind integral equations do not necessarily have unique solutions, but we will restrict ourselves to the unique solution case in analyzing convergence. In particular, we will assume that the second-kind integral equation operator has a bounded inverse.

Before beginning the derivation, let's readdress the notation definitions.

Let K denote the integral operator, and therefore the general form is

$$K\sigma = \int G(x, x')\sigma(x')dx'$$

Let σ_n denote a numerical approximation to σ on x based on using n basis functions. Note here that σ_n is a function of x and would typically be given by

$$\sigma_n(x) = \sum_{i=1}^n \sigma_{ni}\varphi_i(x).$$

Let K_n be the discrete representation of the integral operator. Note that K_n is *not* the matrix \underline{K}_n , but an operator that maps a function of x into another function of x . For example, if the discretization scheme uses a basis to approximate σ , and the basis weights were determined by a collocation scheme, a not necessarily unique associated K_n could be given by

$$(K_n\sigma)(x) = V \left(\int G(x, x')P\sigma(x')dx' \right)$$

where in the orthonormal basis set case

$$(P\sigma)(x) = \sum_{i=1}^n \left(\int \sigma(x')\varphi_i(x')dx' \right) \varphi_i(x), \quad (19)$$

and

$$(Vu)(x) = \sum_{i=1}^n u(x_{c_i})\varphi_i(x). \quad (20)$$

where u is a arbitrary function used to define the action of operator V .

Equations (19) and (20) deserve some explanation. The piecewise constant basis is orthonormal, so the formula in equation (19) is a simple projection of σ onto the basis. If centroid collocation is used, then the discrete potentials computed by evaluating the integral operator at the collocation points must be converted to a function of x by interpolation. In equation (20), the $\varphi_i(x)$'s act as interpolation functions.

With the examples of how Galerkin and centroid-collocation discretization schemes lead to function and operator representations, the second-kind integral equation convergence theory can be presented in a very transparent fashion, as will be show below. What the theorem demonstrates is that if a discretization scheme generates progressively more accurate representations of the integral operator as n increases, then the discretization method converges. That is,

$$\lim_{n \rightarrow \infty} \|\sigma - \sigma_n\| \rightarrow 0,$$

where the comparison between σ and σ_n is unambiguous as both are functions of x .

▷ **Exercise 5** Suppose a nonorthogonal basis is used to represent σ . How would the projection operator in equation (19) change? ■

3.2 Main Theorem

Given $(I + K)\sigma = \Psi$ & $\|(I + K)^{-1}\| < C$
 Means Equation uniquely solvable
 $(I + K_n)\sigma_n = \Psi_n$
 Remainder of Discrete Equivalent

Consistency:

If $\lim_{n \rightarrow \infty} \max_{\|\sigma\|=1} \|(K - K_n)\sigma\| \rightarrow 0$ and $\lim_{n \rightarrow \infty} \|\Psi - \Psi_n\| \rightarrow 0$

Then $\lim_{n \rightarrow \infty} \|\sigma - \sigma_n\| \rightarrow 0$

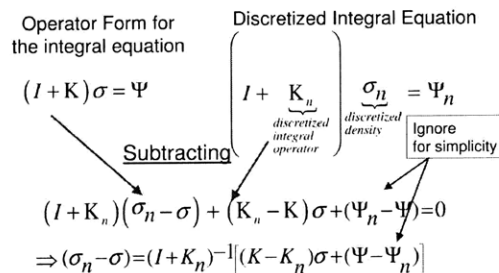
3.3 Rough Proof

To derive a relationship between the errors in the computed solution and the errors in the operator representation, we write the exact equation alongside the discrete equation.

Exact Equation	Discretized Equation
$\Psi = (I + K)\sigma$	$\Psi_n = (I + K_n)\sigma_n$

$$\begin{aligned} \Rightarrow \Psi - \Psi_n &= (I + K)\sigma - (I + K_n)\sigma_n \\ \Rightarrow (\Psi - \Psi_n) &= (\sigma - \sigma_n) + K\sigma - K_n\sigma_n \\ \Rightarrow (\Psi - \Psi_n) &= (\sigma - \sigma_n) + K\sigma - K_n\sigma + K_n\sigma - K_n\sigma_n \\ \Rightarrow (\Psi - \Psi_n) &= (I + K_n)(\sigma - \sigma_n) + (K - K_n)\sigma \\ \Rightarrow (\Psi - \Psi_n) - (K - K_n)\sigma &= (I + K_n)(\sigma - \sigma_n) \\ \Rightarrow (I + K_n)^{-1} [(\Psi - \Psi_n) - (K - K_n)\sigma] &= (\sigma - \sigma_n). \end{aligned}$$

The results on the slide below skip many of these intermediate steps but present the essential results.



The equation for the solution error (previous slide)

$$\underbrace{(\sigma_n - \sigma)}_{\text{solution error}} = (I + K_n)^{-1}(K - K_n)\sigma$$

Taking norms

$$\underbrace{\|\sigma_n - \sigma\|}_{\substack{\text{Error which} \\ \text{should go to} \\ \text{zero as } n \\ \text{increases}}} \leq \underbrace{\|(I + K_n)^{-1}\|}_{\substack{\text{Needs a} \\ \text{bound, that is} \\ \text{stability}}} \underbrace{\|(K - K_n)\sigma\|}_{\substack{\text{Goes to} \\ \text{zero} \\ \text{by consistency}}}$$

We complete deriving a relationship between the errors in the computed solution and the errors in the operator representation. In order to establish that consistency implies convergence, the inverse of the discretized operator must be bounded.

3.3.1 Stability Bound

Norm of solution error

$$\|(\sigma_n - \sigma)\| \leq \|(I + K_n)^{-1}\| \|(K - K_n)\sigma\|$$

Deriving the stability bound

$$(I + K_n)^{-1} = [I + K - (K - K_n)]^{-1} = [(I - (I + K)^{-1}(K - K_n))]^{-1}(I + K)^{-1}$$

Taking norms

$$\|(I + K_n)^{-1}\| \leq \underbrace{\|(I + K)^{-1}\|}_{\substack{\text{Bounded by } C \\ \text{by Assumption}}} \|(I - (I + K)^{-1}(K - K_n))^{-1}\|$$

Repeating from last slide

$$\|(I + K_n)^{-1}\| \leq \underbrace{\|(I + K)^{-1}\|}_{\substack{\text{Bounded by } C \\ \text{by Assumption}}} \|(I - (I + K)^{-1}(K - K_n))^{-1}\|$$

Bounding terms

$$\|(I + K_n)^{-1}\| \leq \frac{C}{1 - \underbrace{\|(I + K)^{-1}(K - K_n)\|}_{\text{Will be less than 0.5 for } n \text{ larger than some } n_0 \text{ by consistency and solvability}}} < 2C \text{ for}$$

$$n \geq n_0$$

Final result:

$$\lim_{n \rightarrow \infty} \|(K - K_n)\sigma\| = 0$$

IMPLIES

$$\lim_{n \rightarrow \infty} \|(\sigma_n - \sigma)\| = 0$$

What does this mean?

The discretization convergence of a second kind integral equation solver only depends on how well the integral is approximated.

The final result, noted on the above slide, is that the solution error is bounded by a constant multiplying the error in the integral operator representation. This suggests that any method which can accurately represent the integral operator can be used to discretize a second-kind integral equation.

4 Nystrom Method

4.1 1-D Second Kind Example

4.1.1 Collocation Discretization

Integral Equation

$$\Psi(x) = \sigma(x) + \int_{-1}^1 G(x, x')\sigma(x')dS' \quad x \in [-1, 1]$$

Apply **quadrature** to **Collocation equation**

$$\begin{aligned} \Psi(x_i) &= \sigma(x_i) + \int_{-1}^1 G(x_i, x')\sigma(x')dS' \\ \Rightarrow \Psi(x_i) &= \sigma(x_i) + \sum_{j=1}^n w_j G(x_i, x_j)\sigma(x_j) \end{aligned}$$

After applying quadrature to Collocation:

$$\Psi(x_i) = \sigma(x_i) + \sum_{j=1}^n w_j G(x_i, x_j)\sigma(x_j)$$

x_i is a collocation point

x_j 's are quadrature points

Now set **quadrature points = collocation points**

In Gaussian quadrature, described in the previous lecture, an integral is approximated using weighted combinations of the integrand. As a reminder, the Gaussian quadrature formula for integrating a function on the unit interval is

$$\int_0^1 f(x)dx \simeq \sum_{i=1}^n w_i f(x_i)$$

where the x_i 's are the evaluation points given by the zeros of an n^{th} -order orthogonal polynomial on the unit interval, and the w_i 's are the weights determined by solving exactness equations.

The key idea behind a Nystrom method for discretizing an integral equation is to use the Gaussian quadrature evaluation points as the test points in a collocation method for solving an integral equation. Then, the collocation method integrals can be approximated using the Gaussian quadrature scheme, resulting in a system of equations which only require evaluations of the integrand at the test=quadrature points. The second kind theory predicts that the error in such a scheme is proportional to the error in the quadrature scheme for computing the collocation integrals.

Set **quadrature points = collocation points**

$$\begin{aligned} \Psi(x_1) &= \sigma_{n1} + \sum_{j=1}^n w_j G(x_1, x_j) \sigma_{nj} \\ &\vdots \\ \Psi(x_n) &= \sigma_{nn} + \sum_{j=1}^n w_j G(x_n, x_j) \sigma_{nj} \end{aligned}$$

System of n equations in n unknowns

Collocation equation per quad/colloc point

Unknown density per quad/colloc point

4.1.2 Discretization-Matrix Comparison

Nystrom Matrix

$$\begin{bmatrix} 1 + w_1 G(x_1, x_1) & \cdots & w_n G(x_1, x_n) \\ \vdots & \ddots & \vdots \\ w_1 G(x_n, x_1) & \cdots & 1 + w_n G(x_n, x_n) \end{bmatrix} \begin{bmatrix} \sigma_{n1} \\ \vdots \\ \sigma_{nn} \end{bmatrix} = \begin{bmatrix} \Psi(x_1) \\ \vdots \\ \Psi(x_n) \end{bmatrix}$$

Piecewise Constant Collocation Matrix

$$\begin{bmatrix} 1 + \int_{x_1}^{x_2} G(x_1, x') dx' & \cdots & \int_{x_n}^{x_{n+1}} G(x_1, x') dx' \\ \vdots & \ddots & \vdots \\ \int_{x_1}^{x_2} G(x_n, x') dx' & \cdots & 1 + \int_{x_n}^{x_{n+1}} G(x_n, x') dx' \end{bmatrix} \begin{bmatrix} \sigma_{n1} \\ \vdots \\ \sigma_{nn} \end{bmatrix} = \begin{bmatrix} \Psi(x_1) \\ \vdots \\ \Psi(x_n) \end{bmatrix}$$

Nystrom Matrix

Just Green's function evaluations – No integrals
 Entries each have a quadrature weight
 Collocation points are quadrature points
 High order quadrature=faster convergence?

Piecewise Constant Collocation Matrix

Integrals of Green's function over line sections
 Distant terms equal Green's function
 Collocation points are basis function centroids
 Low order method always

4.2 K_n and Ψ_n for Nystrom Method

$$K_n \sigma = \sum_{i=1}^n \left(\sum_{j=1}^n w_j G(x_i, x_j) \sigma(x_j) \right) \varphi_i(x)$$

$$\Psi_n = \sum_{i=1}^n \Psi(x_i) \varphi_i(x)$$

4.3 Convergence

4.3.1 Theorem

In the limit as $n \rightarrow \infty$ (number of quad points $\rightarrow \infty$)

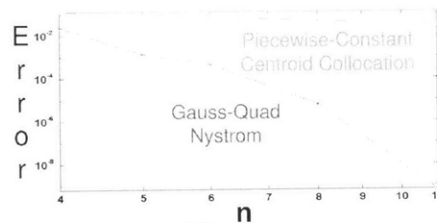
The discretization error = $\max_{\|\sigma\|=1} \|(K - K_n)\sigma\| \rightarrow 0$

AT THE SAME RATE as the underlying quadrature!!

Gauss Quadrature \Rightarrow Exponential Convergence!

4.3.2 Comparison

$$\cos 2\pi x = \sigma(x) + \int_{-1}^1 (x - x')^2 \sigma(x') dS'$$



4.3.3 Caveat

If Nystrom method can have exponential convergence, why use anything else?

Gaussian quadrature has exponential convergence for **nonsingular** kernels

Most physical problems of interest have **singular kernels** ($\frac{1}{r}$, $\frac{\exp ikr}{r}$, etc)

5 Summary

Convergence Issues in 1D

1st and 2nd kind integral equations, null spaces

Convergence for second kind equations

Show consistency and stability issues

Nystrom methods

High order convergence

Did not address singular integrands

1.4 Radiation Conditions and Formulations

Numerical Methods for PDEs

Integral Equation Methods, Lecture 4
Radiation Conditions and Formulations

Notes by L. Proctor, C. Coelho, and J. White

December 3, 2008

1 Outline

Laplace Problems

Exterior Radiation Condition

Potential Representations

Monopole and Dipole Densities

Principle Value

Ansatz and Green's Theorem

Dirichlet and Neumann problems

2 3D Problems

2.1 3D Laplace Equation

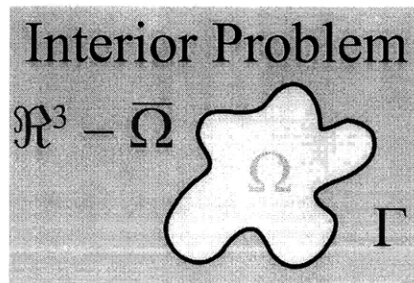
Laplace's equation in 3-D

$$\nabla^2 u(\vec{x}) = \frac{\partial^2 u(\vec{x})}{\partial x^2} + \frac{\partial^2 u(\vec{x})}{\partial y^2} + \frac{\partial^2 u(\vec{x})}{\partial z^2} = 0$$

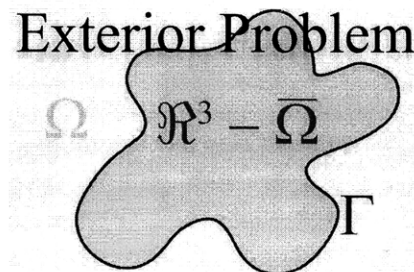
where

$$\vec{x} = (x, y, z) \in \Omega$$

and Ω is bounded by Γ .



The exterior problem is simply the region $\mathbb{R} - \bar{\Omega}$ of the interior problem.



One feature of using integral equation methods is that exterior problems can be solved using the same surface discretization required to solve an interior problem (assuming a linear space-invariant problem like Laplace's equation). This is true even though the exterior domain is infinite and the interior domain is finite. Exterior problems do introduce an additional complication, one must consider the boundary condition "at infinity" (later).

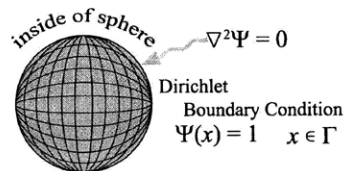
2.2 Boundary Conditions

2.2.1 Dirichlet

Dirichlet Condition

$$u(\vec{x}) = u_{\Gamma}(\vec{x}) \quad \vec{x} \in \Gamma$$

Interior Problem

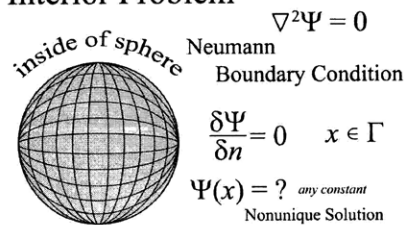


Can you determine the solution to Laplace's equation inside the sphere?

The solution of the interior Dirichlet problem is unique.

2.2.2 Neumann

Interior Problem



The solution of the interior Neumann problem is not unique.

For the solution of the exterior Neumann problem to be unique, it is sufficient to impose a *radiation condition*. In this case, a radiation condition would be a specification of how $u(\vec{x})$ approaches zero as $\vec{x} \rightarrow \infty$.

2.2.3 Exterior

Dirichlet Boundary Condition

$$u(\vec{x}) = u_{\Gamma}(\vec{x}) \quad \vec{x} \in \Gamma$$

Neumann Boundary Condition

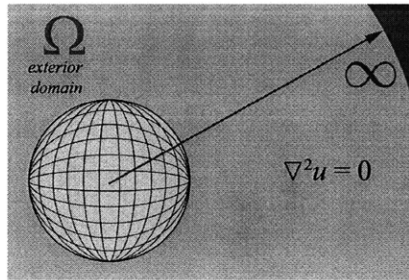
$$\frac{\partial u(\vec{x})}{\partial n_{\vec{x}}} = \frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}} \quad \vec{x} \in \Gamma$$

PLUS

A Radiation Condition

3 Radiation Condition

3.1 Condition at “Infinity”



3.2 Types of Conditions

A radiation condition of the form

$$\lim_{\|\vec{x}\| \rightarrow \infty} u(\vec{x}) \rightarrow 0$$

is not specific enough!
Need

$$\lim_{\|\vec{x}\| \rightarrow \infty} u(\vec{x}) \rightarrow O(\|\vec{x}\|^{-1})$$

OR

$$\lim_{\|\vec{x}\| \rightarrow \infty} u(\vec{x}) \rightarrow O(\|\vec{x}\|^{-2})$$

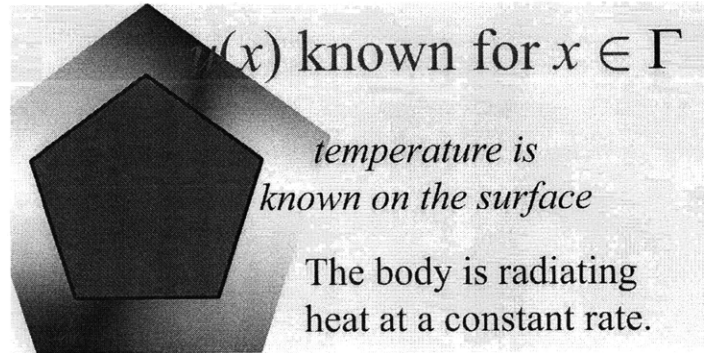
3.3 Examples

The criteria for choosing a radiation condition are best understood by considering several physical examples.

3.3.1 Radiated Heat

Problem Set-Up

$$\begin{aligned}\nabla^2 u(\vec{x}) &= 0 & \vec{x} \in \Omega \\ u(\vec{x}) &= u_0(\vec{x}) & \vec{x} \in \Gamma \\ \text{Radiated Heat} &= \int_{\Gamma} \frac{\partial u}{\partial n} dS\end{aligned}$$



Limit of Expanding Domain

- Infinite Problem
Limit as $R \rightarrow \infty$
- Heat Leaving Sphere
 $\int_{\text{Sphere}_R} \frac{\partial u}{\partial n} dS$ is Constant!

3.3.2 Rad Heat Case 1

$$\lim_{\|\vec{x}\| \rightarrow \infty} u(\vec{x}) \rightarrow O(\|\vec{x}\|^{-1}) \rightarrow \lim_{\|\vec{x}\| \rightarrow \infty} \frac{\partial u(\vec{x})}{\partial n} \rightarrow O(\|\vec{x}\|^{-2})$$

Since the surface of a sphere increases as R^2 :

$$\lim_{R \rightarrow \infty} \int_{\text{Sphere}_R} \frac{\partial u}{\partial n} dS \rightarrow \text{Constant}$$

Radiation condition models net heat loss.

3.3.3 Rad Heat Case 2

$$\lim_{\|\vec{x}\| \rightarrow \infty} u(\vec{x}) \rightarrow O(\|\vec{x}\|^{-2}) \rightarrow \lim_{\|\vec{x}\| \rightarrow \infty} \frac{\partial u(\vec{x})}{\partial n} \rightarrow O(\|\vec{x}\|^{-3})$$

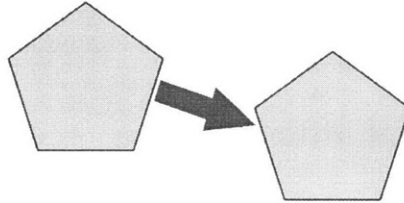
AND

$$\int_{\text{Sphere}_R} \frac{\partial u}{\partial n} dS \rightarrow 0$$

Can NOT model heat loss!

3.3.4 Heat Transfer

$$\nabla^2 u(\vec{x}) = 0 \quad \vec{x} \in \Omega \quad u(\vec{x}) = u_0(\vec{x}) \quad \vec{x} \in \Gamma$$



Heat flows from higher temperature object to lower temperature object, but no heat radiates out.

If

$$\lim_{\|\vec{x}\| \rightarrow \infty} u(\vec{x}) \rightarrow O(\|\vec{x}\|^{-2})$$

Then

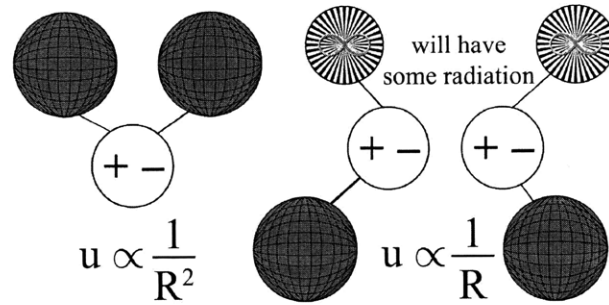
$$\lim_{\|\vec{x}\| \rightarrow \infty} \frac{\partial u(\vec{x})}{\partial n} \rightarrow O(\|\vec{x}\|^{-3})$$

And

$$\int_{\text{Sphere}_R} \frac{\partial u(\vec{x})}{\partial n} dS \rightarrow 0$$

This condition ensures all heat transferred.

3.3.5 Electrostatics



The above image is supposed to represent two scenarios, each scenario has two conducting bodies.

In the first scenario, on the left, there are two conductors that are treated as if there is a voltage source across them. There is a positive charge on one conductor, a negative charge on the other conductor, and the net charge should be zero. If the net charge on the two conductors is zero, then the integral of the normal electric field, $(\frac{\partial u(\vec{x})}{\partial n})$ over a bounding sphere, one that contains both spherical conductors, should be zero. This is just a statement of the well-known Gauss's theorem in electrostatics. In particular

$$\int_{Sphere_R} \frac{\partial u(\vec{x})}{\partial n} dS = 0$$

for any $Sphere_R$ containing the two conductors. In order for this integral to be zero regardless as the radius of $Sphere_R$ approaches ∞ ,

$$\lim_{\|\vec{x}\| \rightarrow \infty} \frac{\partial u(\vec{x})}{\partial n}$$

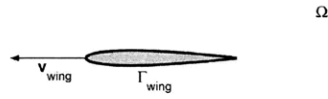
must decay at least as fast as $O(\|\vec{x}\|^{-3})$. Therefore, $u(\vec{x})$ must decay like $O(\|\vec{x}\|^{-2})$.

For the scenario on the right, the two spherical conductors are set to two different potentials with respect to the point at infinity. In such a case, it is unlikely that the sum of the charge on the two spheres will be zero,

$$\int_{Sphere_R} \frac{\partial u(\vec{x})}{\partial n} dS \neq 0$$

and therefore $u(\vec{x})$ should not be forced to decay any faster than $O(\|\vec{x}\|^{-2})$.

3.3.6 Potential Flow



Assumptions

Irrotational flow (velocity = potential gradient):

$$\mathbf{v}(\vec{x}) = \nabla u(\vec{x})$$

Air is incompressible (velocity divergence free):

$$\nabla \cdot \mathbf{v}(\vec{x}) = \nabla^2 u(\vec{x}) = 0.$$

Nonpenetrating wing boundary condition:

$$\nabla u(\vec{x}) \cdot \mathbf{n}(\vec{x}) = \mathbf{v}_{wing}(\vec{x}) \cdot \mathbf{n}(\vec{x}).$$

What is the right radiation condition?

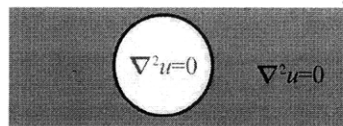
4 Formulations – Problem Types

4.1 Single Domain

	EXTERIOR	INTERIOR
DIRICHLET		
NEUMANN		

4.2 Coupled Domain

Example: Bimetallic Electrical Conductivity



Potential and Electric Current Continuity:

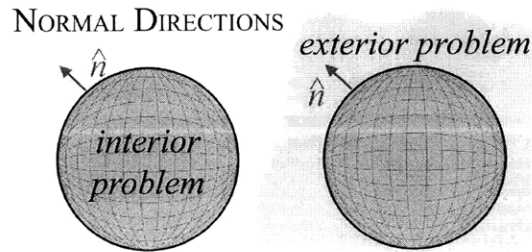
$$u(\vec{x}^+) = u(\vec{x}^-) \quad \vec{x} \in \Gamma$$

$$\alpha^+ \frac{\partial u(\vec{x}^+)}{\partial n_{\vec{x}}} = \alpha^- \frac{\partial u(\vec{x}^-)}{\partial n_{\vec{x}}} \quad \vec{x} \in \Gamma$$

An example of a coupled domain problem would be a conductivity problem involving multiple materials. To determine the electrical conductivity between two terminals of an object made of multiple materials, one would determine the ratio of the voltage across the object's terminals and the current flowing through the object. Electrical current density in an ideal linear conductor is a vector quantity given by the gradient of the potential, ∇u , scaled by a factor known as the conductivity of the material. In an ideal linear conductor there is no accumulation of charge at any interior point, implying that the current density has zero divergence. Therefore, the potential in an ideal linear conductor satisfies Laplace's equation, $\nabla^2 u = 0$. If an object is made of multiple materials with different electrical conductivities, then the boundary between materials satisfies interface conditions. At the boundary between materials, both the potential and the current density in the surface-normal direction are continuous. Since the conductivities of the two materials are different, continuity of the current density implies a jump in the gradient of the potential across the material boundary.

4.3 Normals

Normals usually point from Interior to Exterior.



Typically, the surface normal is assumed to point in the direction from the interior domain to the exterior domain. There are many situations where this typical practice is confusing or ambiguous, so it is often necessary to be explicit about the direction of the normal.

5 Surface Density Integrals

5.1 Monopole & Dipole

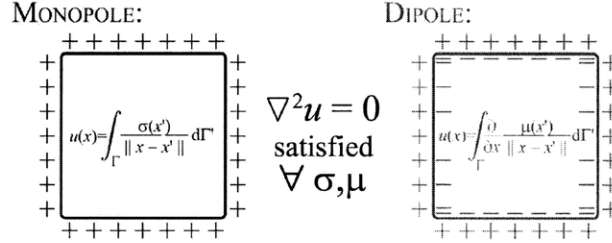
Potential due to a monopole density (σ):

$$u(\vec{x}) = \int_{\Gamma} \frac{1}{\|\vec{x} - \vec{x}'\|} \sigma(\vec{x}') d\Gamma'$$

Potential due to a dipole density (μ):

$$u(\vec{x}) = \int_{\Gamma} \frac{\partial}{\partial n_{\vec{x}'}} \frac{1}{\|\vec{x} - \vec{x}'\|} \mu(\vec{x}') d\Gamma'$$

where the normal points out of the domain Ω bounded by Γ .



Monopole or dipole densities can be used to generate potentials that satisfy $\nabla^2 u(\vec{x}) = 0$ for all $\vec{x} \in \Omega$. The monopole and dipole potentials differ in the radiation condition they satisfy. If the surface, γ , is finite in extent, then in the limit as $\|\vec{x}\| \rightarrow \infty$, the monopole potential decays like $\|\vec{x}\|^{-1}$, and the dipole potential decays like $\|\vec{x}\|^{-2}$.

Either representation can be used to derive surface integral equations, but care must be used when evaluating the associated potentials when $\vec{x} \in \Gamma$.

5.2 Surface Potentials

The monopole potential is continuous as x passes through Γ , so

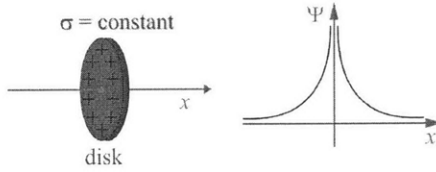
$$u_{\Gamma}(\vec{x}) = \int_{\Gamma} \frac{1}{\|\vec{x} - \vec{x}'\|} \sigma(\vec{x}') d\Gamma' \quad \vec{x} \in \Gamma$$

The dipole potential “jumps” as x passes through Γ , so the limit as $\vec{x} \rightarrow \Gamma$ of

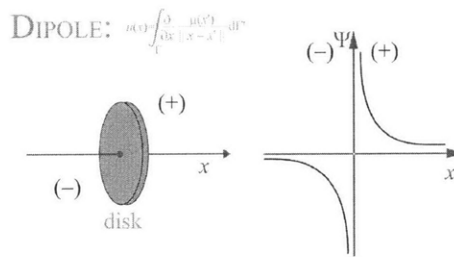
$$u(\vec{x}) = \int_{\Gamma} \frac{\partial}{\partial n_{\vec{x}'}} \frac{1}{\|\vec{x} - \vec{x}'\|} \mu(\vec{x}') d\Gamma'$$

depends on how Γ is approached.

MONOPOLE:
$$v(x) = \int_{\Gamma} \frac{\sigma(y)}{|y-x|} d\Gamma$$



Don't be put-off by the graph above. The monopole potential is continuous (it does not go off to infinity, as it may seem to in the above figure), but it is not continuously differentiable, there will be a discontinuity in the derivative at x_0 .



5.2.1 Principle Value Integral

If $f(y)$ is singular for some $y = x_0$, where $x_0 \in \Gamma$, then the principle value integral is

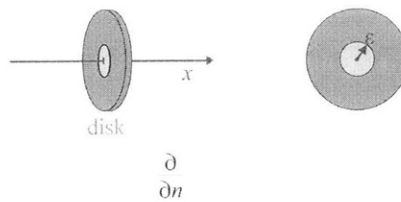
$$\int_{\Gamma}^{PV} f(\vec{y}) d\Gamma \equiv \lim_{\epsilon \rightarrow 0} \int_{\Gamma - B(x_0, \epsilon) \cap \Gamma} f(\vec{y}) d\Gamma$$

when $B(x_0, \epsilon)$ is the ϵ radius ball about x_0 .

The P.V. is a special kind of limit

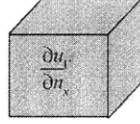
Limit of deleting and ever shrinking portion of the integration domain.
 NOT EQUIVALENT TO limiting processes on f !

CAUCHY PRINCIPLE VALUE INTEGRAL



5.2.2 Monopole Derivative (MD)

Consider a cube geometry:



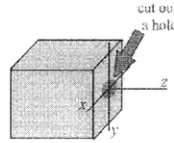
$$\frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}} = \lim_{\vec{x} \rightarrow \Gamma^+} \frac{\partial}{\partial n_{\vec{x}}} \int_{\Gamma} \frac{1}{\|\vec{x} - \vec{x}'\|} \sigma(\vec{x}') d\Gamma' \quad \vec{x} \in \Gamma$$

The plus (+) in Γ^+ indicates exterior approach.

In the above slide, we consider computing the normal derivative of the monopole potential just outside the boundary γ . As will be shown in the next few slides, the derivative can be represented as the sum of a principle value integral and an extra term.

5.2.3 MD Disk Removal

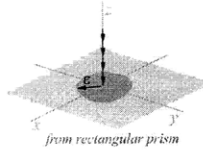
$$\frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}} = \lim_{\vec{x} \rightarrow \Gamma^+, \epsilon \rightarrow 0} \left[\frac{\partial}{\partial n_{\vec{x}}} \int_{\Gamma - B(x, \epsilon)} \frac{\sigma(\vec{x}')}{\|\vec{x} - \vec{x}'\|} d\Gamma' + \frac{\partial}{\partial n_{\vec{x}}} \int_{B(x, \epsilon)} \frac{\sigma(\vec{x}')}{\|\vec{x} - \vec{x}'\|} d\Gamma' \right]$$



Consider the entire side panel on the right of the cube in the above slide, and consider evaluating the normal derivative of the potential generated by a monopole distribution on the surface of the cube's right side. Specifically, assume that we wish to evaluate the derivative at a point \vec{x} in the center of the green disk on the cube's right side. The matter is complicated by the fact that the integrand goes to infinity when $\vec{x} = \vec{x}'$. Thus, we need to break up the integral into two pieces. One piece is the entire panel minus the green disk, and the other piece is just the green disk.

5.2.4 MD Disk Picture

$$\frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}} = \lim_{\vec{x} \rightarrow \Gamma^+, \epsilon \rightarrow 0} \left[\frac{\partial}{\partial n_{\vec{x}}} \int_{\Gamma - B(x, \epsilon)} \frac{\sigma(\vec{x}')}{\|\vec{x} - \vec{x}'\|} d\Gamma' + \frac{\partial}{\partial n_{\vec{x}}} \int_{B(x, \epsilon)} \frac{\sigma(\vec{x}')}{\|\vec{x} - \vec{x}'\|} d\Gamma' \right]$$

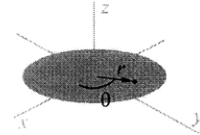


The first integral is the Principle Value Integral and the second integral is the integral of just the disk.

Given that the disk was extracted from right the surface of a cube, the disk is flat, and the normal is in the z -axis direction.

5.2.5 MD Disk Eval

$$\begin{aligned} \lim_{\vec{x} \rightarrow \Gamma^+} \frac{\partial}{\partial n_{\vec{x}}} \int_{B(x, \epsilon)} \frac{1}{\|\vec{x} - \vec{x}'\|} \sigma(\vec{x}') d\Gamma' \\ \approx \lim_{z \rightarrow 0^+} \frac{\partial}{\partial z} \int_0^{2\pi} \int_0^\epsilon \frac{\sigma(\vec{x})}{\sqrt{r^2 + z^2}} r dr d\theta \\ = \lim_{z \rightarrow 0^+} \frac{\partial}{\partial z} 2\pi \sigma(\vec{x}) \left[\sqrt{\epsilon^2 + z^2} - |z| \right] \\ = -2\pi \sigma(\vec{x}). \end{aligned}$$



Note 1

Disk Evaluation Math

For this problem, it is quite straightforward to see how one changes from cartesian to cylindrical coordinates. But, the algebra and calculus involved in solving this integral may not be as straightforward, herein is presented one method, broken-down into bite-size pieces:

$$\lim_{z \rightarrow 0^+} \frac{\partial}{\partial z} \int_0^{2\pi} \int_0^\epsilon \frac{\sigma(\vec{x})}{\sqrt{r^2 + z^2}} r dr d\theta.$$

Use trigonometric substitution to solve the integral with respect to r . Substitute $r = z \tan \alpha$ and $dr = z \sec^2 \alpha d\alpha$ and simplify to get the following integral:

$$= \lim_{z \rightarrow 0^+} \frac{\partial}{\partial z} \int_0^{2\pi} \int_0^{\alpha(\epsilon)} \frac{\sigma(\vec{x}) z \sin \alpha}{\cos^2 \alpha} d\alpha d\theta.$$

This integral is easily solved using direct substitution of $u = \cos \alpha$:

$$= \lim_{z \rightarrow 0^+} \frac{\partial}{\partial z} \int_0^{2\pi} \int_0^{u[\alpha(\epsilon)]} -z \frac{\sigma(\vec{x})}{u^2} du d\theta = \lim_{z \rightarrow 0^+} \frac{\partial}{\partial z} \int_0^{2\pi} z \frac{\sigma(\vec{x})}{u} \Big|_0^{u[\alpha(\epsilon)]} d\theta.$$

Plug back in for $u = \cos \alpha = \frac{z}{\sqrt{r^2 + z^2}} \Big|_{r=\epsilon}$

$$\begin{aligned} &= \lim_{z \rightarrow 0^+} \frac{\partial}{\partial z} \int_0^{2\pi} \sigma(\vec{x}) \sqrt{r^2 + z^2} \Big|_0^{\epsilon} d\theta \\ &= \lim_{z \rightarrow 0^+} \frac{\partial}{\partial z} \int_0^{2\pi} \sigma(\vec{x}) \left[\sqrt{\epsilon^2 + z^2} - \sqrt{z^2} \right] d\theta. \end{aligned}$$

Integrating the last part is quite simple,

$$= \lim_{z \rightarrow 0^+} \frac{\partial}{\partial z} \sigma(\vec{x}) \left(\sqrt{\epsilon^2 + z^2} - \sqrt{z^2} \right) \theta \Big|_0^{2\pi} = 2\pi \sigma(\vec{x}) \lim_{z \rightarrow 0^+} \frac{\partial}{\partial z} \left(\sqrt{\epsilon^2 + z^2} - \sqrt{z^2} \right).$$

Finally, take the derivatives with respect to z , the normal,

$$= 2\pi \sigma(\vec{x}) \lim_{z \rightarrow 0^+} \left(\frac{z}{\sqrt{\epsilon^2 + z^2}} - \frac{z}{\sqrt{z^2}} \right).$$

It can now be seen that, since $\lim_{z \rightarrow 0^+} \frac{z}{\sqrt{\epsilon^2 + z^2}} = 0$ and $\lim_{z \rightarrow 0^+} \frac{z}{\sqrt{z^2}} = \text{sign } z$ that

$$\lim_{z \rightarrow 0^+} \frac{\partial}{\partial z} \int_0^{2\pi} \int_0^\epsilon \frac{\sigma(\vec{x})}{\sqrt{r^2 + z^2}} r dr d\theta = -2\pi \sigma(\vec{x}).$$

5.2.6 MD Final

$$\begin{aligned} \frac{\partial u_\Gamma(\vec{x})}{\partial n_{\vec{x}}} &= \lim_{\vec{x} \rightarrow \Gamma^+} \frac{\partial}{\partial n_{\vec{x}}} \int_\Gamma \frac{1}{\|\vec{x} - \vec{x}'\|} \sigma(\vec{x}') d\Gamma' \\ &= \lim_{\vec{x} \rightarrow \Gamma^+} \lim_{\epsilon \rightarrow 0} \left[\frac{\partial}{\partial n_{\vec{x}}} \int_{\Gamma - B(x, \epsilon)} \frac{\sigma(\vec{x}')}{\|\vec{x} - \vec{x}'\|} d\Gamma' + \frac{\partial}{\partial n_{\vec{x}}} \int_{B(x, \epsilon)} \frac{\sigma(\vec{x}')}{\|\vec{x} - \vec{x}'\|} d\Gamma' \right] \\ &= \int_\Gamma^{PV} \frac{\partial}{\partial n_{\vec{x}}} \frac{1}{\|\vec{x} - \vec{x}'\|} \sigma(\vec{x}') d\Gamma' - 2\pi \sigma(\vec{x}) \end{aligned}$$

5.2.7 Dipole Potentials (DP)

If Γ is a flat surface

$$\int_\Gamma^{PV} \frac{\partial}{\partial n_{\vec{x}'}} \frac{1}{\|\vec{x} - \vec{x}'\|} d\Gamma' = 0 \quad \vec{x} \in \Gamma.$$

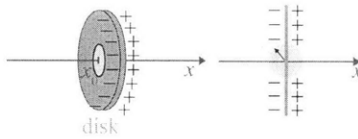
Why? Rewrite using explicit form of integrand

$$\int_\Gamma^{PV} \frac{(\vec{x} - \vec{x}')^T n_{\vec{x}'}}{(\|\vec{x} - \vec{x}'\|)^3} d\Gamma'$$

Integrand is zero when $\vec{x} - \vec{x}'$ orthogonal to surface normal

5.2.8 DP Flat Surface

Flat result applies locally on smooth surfaces.



5.2.9 DP General Surface

If Γ is a general surface

$$u_{\Gamma}(\vec{x}) = 2\pi\mu(\vec{x}) + \int_{\Gamma}^{PV} \frac{(\vec{x} - \vec{x}')^T n_{\vec{x}'}}{(\|\vec{x} - \vec{x}'\|)^3} \mu(\vec{x}') d\Gamma'$$

when \vec{x} approaches Γ from outside Ω , and

$$u_{\Gamma}(\vec{x}) = -2\pi\mu(\vec{x}) + \int_{\Gamma}^{PV} \frac{(\vec{x} - \vec{x}')^T n_{\vec{x}'}}{(\|\vec{x} - \vec{x}'\|)^3} \mu(\vec{x}') d\Gamma'$$

when \vec{x} approaches Γ from inside.

Note 1 *Derivation of the Dipole Surface Potential*

The following derivation goes through the step-by-step process of deriving the dipole surface potential when \vec{x} approaches Γ from inside.

$$\begin{aligned} u(\vec{x}) &= \int_{\Gamma} \frac{\partial}{\partial n_{\vec{x}'}} \frac{1}{\|\vec{x} - \vec{x}'\|} \mu(\vec{x}') d\Gamma' \quad \vec{x} \in \Gamma \\ &= \lim_{\epsilon \rightarrow 0} \int_{\Gamma - B(x_0, \epsilon)} \frac{\partial}{\partial n_{\vec{x}'}} \frac{1}{\|\vec{x} - \vec{x}'\|} \mu(\vec{x}') d\Gamma' + \lim_{\epsilon \rightarrow 0} \int_{B(x_0, \epsilon)} \frac{\partial}{\partial n_{\vec{x}'}} \frac{1}{\|\vec{x} - \vec{x}'\|} \mu(\vec{x}') d\Gamma' \\ &= \int_{\Gamma}^{PV} \frac{\partial}{\partial n_{\vec{x}'}} \frac{1}{\|\vec{x} - \vec{x}'\|} \mu(\vec{x}') d\Gamma' + \lim_{\epsilon \rightarrow 0} \int_{B(x_0, \epsilon)} \hat{n}_{x'} \cdot \nabla \frac{1}{\|\vec{x} - \vec{x}'\|} \mu(\vec{x}') d\Gamma' \\ &= \int_{\Gamma}^{PV} \frac{(\vec{x} - \vec{x}')^T n_{\vec{x}'}}{\|\vec{x} - \vec{x}'\|^3} \mu(\vec{x}') d\Gamma' + \lim_{r \rightarrow 0} \int_0^{\pi} \int_0^{\pi} \hat{n}_{x'} \cdot \nabla \left(\frac{1}{r} \right) \mu r^2 \sin \phi \, d\phi \, d\theta \\ &= 0 + \lim_{r \rightarrow 0} \int_0^{\pi} \int_0^{\pi} \mu \frac{\partial}{\partial r} \left(\frac{1}{r} \right) r^2 \sin \phi \, d\phi \, d\theta \\ &= \lim_{r \rightarrow 0} \int_0^{\pi} \int_0^{\pi} \mu \left(\frac{-1}{r^2} \right) r^2 \sin \phi \, d\phi \, d\theta \\ &= - \int_0^{\pi} \int_0^{\pi} \mu \sin \phi \, d\phi \, d\theta = \left[\int_0^{\pi} \mu \cos \phi \, d\theta \right]_0^{\pi} = \int_0^{\pi} -2\mu \, d\theta \\ &= \left[-2\mu\theta \right]_0^{\pi} = -2\pi\mu \end{aligned}$$

6 Ansatz Formulations

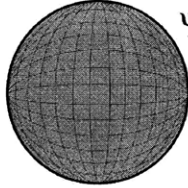
6.1 Dirichlet Problem

6.1.1 Monopole Potential

For an interior or exterior problem:

$$u_{\Gamma}(\vec{x}) = \int_{\Gamma} \frac{1}{\|\vec{x} - \vec{x}'\|} \sigma(\vec{x}') d\Gamma'$$

$$\Psi(x) = 1 \quad \sigma(x) = ?$$

$$x \in \Gamma$$


$$\Psi(x) = \int_{\Gamma} \frac{\sigma(x')}{\|x - x'\|} d\Gamma'$$

$$\Rightarrow \sigma = 4\pi$$

What radiation condition?

The point of the above slide is to show that when using monopole potentials to solve Dirichlet problems on a single domain bounded by γ , the equations are the same for either the interior or the exterior problem.

6.1.2 Dipole Potential

For an exterior problem:

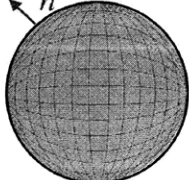
$$u_{\Gamma}(\vec{x}) = 2\pi\mu(\vec{x}) + \int_{\Gamma}^{PV} \frac{(\vec{x} - \vec{x}')^T n_{\vec{x}'}}{(\|\vec{x} - \vec{x}'\|)^3} \mu(\vec{x}') d\Gamma'$$

and for an interior problem:

$$u_{\Gamma}(\vec{x}) = -2\pi\mu(\vec{x}) + \int_{\Gamma}^{PV} \frac{(\vec{x} - \vec{x}')^T n_{\vec{x}'}}{(\|\vec{x} - \vec{x}'\|)^3} \mu(\vec{x}') d\Gamma'$$

Normal points from interior to exterior.

INTERIOR PROBLEM

$$u(x) = u_{\Gamma}(x) \quad x \in \Gamma$$


$$u_{\Gamma}(x) = -2\pi\mu(x) + \int_{\Gamma}^{PV} \frac{(x - x')^T n_x}{\|x - x'\|} \mu(x') d\Gamma'$$

Note that the radiation condition satisfied by the monopole potential is different than the radiation condition satisfied by the dipole potential, even when used to solve the same Dirichlet problem.

6.2 Neumann Problem

6.2.1 Monopole Potential (MP)

Derivative of the monopole potential “jumps” as x passes through Γ , so

$$\frac{\partial}{\partial n_{\vec{x}}} \int_{\Gamma} \frac{1}{\|\vec{x} - \vec{x}'\|} \sigma(\vec{x}') d\Gamma'$$

takes different values just inside and just outside Γ .

6.2.2 MP Int/Ext

For an exterior problem

$$\frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}} = -2\pi\sigma(\vec{x}') - \int_{\Gamma}^{PV} \frac{(\vec{x} - \vec{x}')^T n_{\vec{x}}}{(\|\vec{x} - \vec{x}'\|)^3} \sigma(\vec{x}') d\Gamma'$$

and for an interior problem

$$\frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}} = +2\pi\sigma(\vec{x}') - \int_{\Gamma}^{PV} \frac{(\vec{x} - \vec{x}')^T n_{\vec{x}}}{(\|\vec{x} - \vec{x}'\|)^3} \sigma(\vec{x}') d\Gamma'$$

Normal points from interior to exterior.

Note that the signs for the Neumann monopole potential integral equation are different than the signs for the integral equation in the Dirichlet Dipole case. The sign changes are due to the location of the derivative evaluation. In the Neumann Monopole potential case, the derivative is taken with respect to $n_{\vec{x}}$ whereas for the Dirichlet Dipole potential case, the derivative is taken with respect to $n_{\vec{x}'}$.

6.3 Kinds of Equations

First Kind Equations

- Dirichlet Monopole potential integral equation.

Second Kind Equations

- Dirichlet Dipole potential integral equation.
- Neumann Monopole potential integral equation.

Dipole potential for Neumann?

7 Green's Theorem

Laplace's Equation Green's Function

$$\nabla^2 G(\vec{x}) = 4\pi\delta(\vec{x})$$

$\delta(\vec{x}) \equiv$ impulse in 3-D

Defined by its behavior in an integral

$$\int \delta(\vec{x}') f(\vec{x}') d\Omega' = f(0)$$

Not too hard to show

$$G(\vec{x}) = \frac{1}{\|\vec{x}\|}$$

Note 1

4π

Just as an aside, Green's Function may be defined using a different scaling variable depending upon which source one is using. Sometimes Laplace's equation will be written:

$$\nabla^2 G(\vec{x}) = \delta(\vec{x}).$$

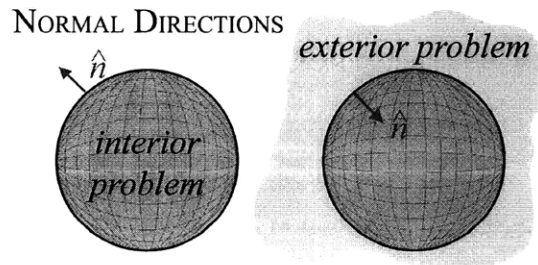
Where you can see that the value of 4π has been left off. This will simply mean that the Green's Function is now written:

$$G(\vec{x}) = \frac{4\pi}{\|\vec{x}\|}.$$

For our purposes, we will be using the notation in the above slide, and not the notation given in this note.

7.1 Normal Directions

A note here about normal directions is essential. In the above section, the "normal points from interior to exterior" whereas, in the image below, the normal points from inside the domain to outside the domain. How would this impact the solution?



When we go through Green's Theorem in the following section, remember that the normal always points "out" of Ω , as it does in the above figure.

7.2 Divergence Thm

The general Divergence Theorem:

For any sufficiently smooth \vec{F}

$$\int_{\Omega} \nabla \cdot \vec{F}(x) dV = \int_{\Gamma} \vec{F} \cdot \vec{n}_x dS$$

where Γ is the surface which encloses Ω .

Green's theorem follows from the **divergence theorem**.

7.3 Volume Theorem

If u satisfies Laplace's equation in Ω , then

$$4\pi u(\vec{x}) = \int_{\Gamma} \left[\frac{\frac{\partial u(\vec{x}')}{\partial n}}{\|\vec{x} - \vec{x}'\|} - \frac{\partial}{\partial n_{\vec{x}'}} \frac{1}{\|\vec{x} - \vec{x}'\|} u(\vec{x}') \right] d\Gamma' \quad \vec{x} \in \Omega$$

where the normal points out of Ω .

7.3.1 Surface

Using the Principle Value Integral:

$$2\pi u(\vec{x}) = \int_{\Gamma} \frac{\frac{\partial u(\vec{x}')}{\partial n}}{\|\vec{x} - \vec{x}'\|} d\Gamma' - \int_{\Gamma}^{PV} \frac{\partial}{\partial n_{\vec{x}'}} \frac{1}{\|\vec{x} - \vec{x}'\|} u(\vec{x}') d\Gamma' \quad \vec{x} \in \Gamma$$

where the normal points out of Ω .

This is one of the cases where it is generally easier to define the normal as pointing out of Ω rather than having the normal point from the interior to the exterior.

7.3.2 Boundary Conditions

The boundary conditions, Dirichlet or Neumann, can be determined by using the surface form of Green's Theorem. For Dirichlet Problems, $u = u_{\Gamma}$ when $\vec{x} \in \Gamma$. So, put the known values for u into Green's Theorem for the surface, and put these known terms on the right hand side, leaving the unknown on the left hand side. Likewise, for Neumann problems, $\frac{\partial u}{\partial n} = \frac{\partial u_{\Gamma}}{\partial n}$ denotes that the derivative of u is known on the boundary when $\vec{x} \in \Gamma$. Again, put the known term on the right hand side, so that the unknown value, u is on the left hand side.

For Dirichlet Problems

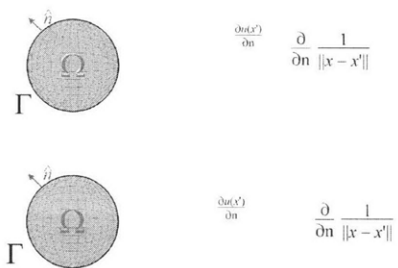
$$\int_{\Gamma} \frac{\frac{\partial u(\vec{x}')}{\partial n}}{\|\vec{x} - \vec{x}'\|} d\Gamma' = 2\pi u_{\Gamma}(\vec{x}) + \int_{\Gamma}^{PV} \frac{\partial}{\partial n_{\vec{x}'}} \frac{u_{\Gamma}(\vec{x}')}{\|\vec{x} - \vec{x}'\|} d\Gamma'$$

For Neuman Problems

$$2\pi u(\vec{x}) + \int_{\Gamma}^{PV} \frac{\partial}{\partial n_{\vec{x}'}} \frac{1}{\|\vec{x} - \vec{x}'\|} u(\vec{x}') d\Gamma' = \int_{\Gamma} \frac{\frac{\partial u_{\Gamma}(\vec{x}')}{\partial n}}{\|\vec{x} - \vec{x}'\|} d\Gamma'$$

where normal points out of Ω (interior or exterior!)

7.4 Overview



8 Summary

Laplace Problems

Exterior Radiation Condition

Potential Representations

Monopole and Dipole potentials

Principle Value

Ansatz and Green's Theorem

Dirichlet and Neumann problems

First and Second Kind Equations.

1.5 First and Second Kind Theory, part 2

Numerical Methods for PDEs

Boundary Element Methods, Lecture 5
First and Second Kind Theory, part 2

Notes by L. Proctor, C. Coelho and J. White

December 8, 2008

1 Outline

Interior Neumann

Use Nystrom to Solve
Look at 2-D Problems

Fredholm Alternative

Connection to Linear Algebra

First Kind Convergence Theory

2 Exterior Formulations

2.1 Dirichlet Problem

2.1.1 Monopole Potential

For an 3D exterior problem:

$$u_{\Gamma}(\vec{x}) = \int_{\Gamma} \frac{1}{\|\vec{x} - \vec{x}'\|} \sigma(\vec{x}') d\Gamma'$$

For an 2D exterior problem:

$$u_{\Gamma}(\vec{x}) = \int_{\Gamma} \log \|\vec{x} - \vec{x}'\| \sigma(\vec{x}') d\Gamma'$$

2.1.2 Dipole Potential

For a 3-D exterior problem:

$$u_{\Gamma}(\vec{x}) = 2\pi\mu(\vec{x}) + \int_{\Gamma}^{PV} \frac{(\vec{x} - \vec{x}')^T n_{\vec{x}'}}{(\|\vec{x} - \vec{x}'\|)^3} \mu(\vec{x}') d\Gamma'$$

For a 2-D exterior problem:

$$u_{\Gamma}(\vec{x}) = \pi\mu(\vec{x}) + \int_{\Gamma}^{PV} \frac{(\vec{x} - \vec{x}')^T n_{\vec{x}'}}{(\|\vec{x} - \vec{x}'\|)^2} \mu(\vec{x}') d\Gamma'$$

Normal points from interior to exterior.

2.2 Neumann Problem

2.2.1 Monopole

For an exterior 3D problem

$$\frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}} = -2\pi\sigma(\vec{x}') - \int_{\Gamma}^{PV} \frac{(\vec{x} - \vec{x}')^T n_{\vec{x}'}}{(\|\vec{x} - \vec{x}'\|)^3} \sigma(\vec{x}') d\Gamma'$$

For an exterior 2D problem

$$\frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}} = -\pi\sigma(\vec{x}') - \int_{\Gamma}^{PV} \frac{(\vec{x} - \vec{x}')^T n_{\vec{x}}}{(\|\vec{x} - \vec{x}'\|)^2} \sigma(\vec{x}') d\Gamma'$$

Normal points from interior to exterior.

3 Interior Example

3.1 3D Case

3.1.1 Monopole Potential

Surface Potential

$$u_{\Gamma}(\vec{x}) = \int_{\Gamma} \frac{1}{\|\vec{x} - \vec{x}'\|} \sigma(\vec{x}') d\Gamma' \quad \vec{x} \in \Gamma$$

Surface Normal Derivative

$$\frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}} = \frac{\partial}{\partial n_{\vec{x}}} \int_{\Gamma} \frac{1}{\|\vec{x} - \vec{x}'\|} \sigma(\vec{x}') d\Gamma' \quad \vec{x} \in \Gamma$$

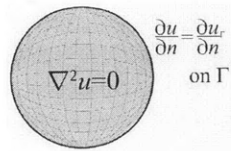
Normal points to exterior.

3.1.2 Interior Neumann

Monopole Potential Using the P.V. Integral

$$\frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}} = 2\pi\sigma(\vec{x}) + \int_{\Gamma}^{PV} \frac{\partial}{\partial n_{\vec{x}}} \frac{1}{\|\vec{x} - \vec{x}'\|} \sigma(\vec{x}') d\Gamma'$$

looks 2nd Kind Equation, Try Nystrom.



3.1.3 Nystrom Method

Set quadrature points = collocation points

$$\frac{\partial u_{\Gamma}(\vec{x}_1)}{\partial n_{\vec{x}}} = 2\pi\sigma_{n1} + \sum_{j=1}^n w_{1,j} \frac{\partial}{\partial n_{\vec{x}}} \frac{1}{\|\vec{x}_1 - \vec{x}_j\|} \sigma_{nj}$$

$$\frac{\partial u_{\Gamma}(\vec{x}_n)}{\partial n_{\vec{x}}} = 2\pi\sigma_{nn} + \sum_{j=1}^n w_{n,j} \frac{\partial}{\partial n_{\vec{x}}} \frac{1}{\|\vec{x}_n - \vec{x}_j\|} \sigma_{nj}$$

n equations in n unknowns

$j = i$ case (self-term)?

3.1.4 $i = j$

For the monopole 3-D Neumann Formulation,

$$G(\vec{x}, \vec{x}') = \frac{\partial}{\partial n_{\vec{x}}} \frac{1}{\|\vec{x} - \vec{x}'\|} = -\frac{(\vec{x} - \vec{x}')^T n_{\vec{x}}}{(\|\vec{x} - \vec{x}'\|)^3}$$

PROBLEM: $G(\vec{x}, \vec{x}')$ blows up as $\vec{x} \rightarrow \vec{x}'$.

3.2 2-D Case

Monopole Neumann Formulation

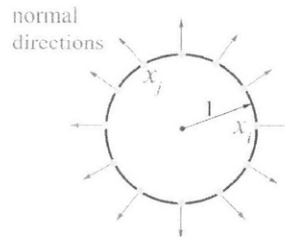
$$\frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}} = \pi\sigma(\vec{x}) + \int_{\Gamma}^{PV} \frac{\partial}{\partial n_{\vec{x}}} \log \|\vec{x} - \vec{x}'\| \sigma(\vec{x}') d\Gamma'$$

Simplifying the Green's function,

$$G(\vec{x}, \vec{x}') = \frac{\partial}{\partial n_{\vec{x}}} \log \|\vec{x} - \vec{x}'\| = -\frac{(\vec{x} - \vec{x}')^T n_{\vec{x}}}{(\|\vec{x} - \vec{x}'\|)^2}$$

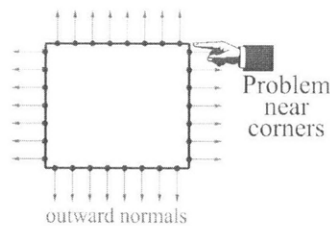
3.2.1 Smooth Γ

$G(\vec{x}, \vec{x}')$ finite as $\vec{x} \rightarrow \vec{x}'$ if Γ smooth.

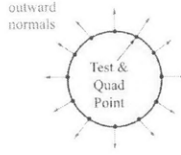


3.2.2 Nonsmooth Γ

$G(\vec{x}, \vec{x}')$ not finite as $\vec{x} \rightarrow \vec{x}'$ if x' is on a corner.



3.2.3 Disk Example



$$\frac{\partial u_{\Gamma}(x_i)}{\partial n_{x_i}} = \pi \sigma(x_i) - \frac{2\pi}{N} \sum \frac{\vec{n}_{x_i} \cdot (x_i - x_j)}{\|x_j - x_i\|^2} \sigma(x_j)$$

Note uniform quadrature weights on the circle.
Resulting matrix is singular! Why?

4 Second Kind Theorem

4.1 Theorem

Given

$$(I + K)\sigma = \Psi \quad (\text{Integral Eqn.})$$

$$(I + K_n)\sigma_n = \Psi_n \quad (\text{Discretized Eqn.})$$

AND

$$\|(I + K)^{-1}\| < C \quad \text{Unique solvability}$$

If

$$\lim_{n \rightarrow \infty} \|(K - K_n)\| \rightarrow 0 \quad \text{and} \quad \|\Psi - \Psi_n\| \rightarrow 0$$

Then

$$\lim_{n \rightarrow \infty} \|\sigma - \sigma_n\| \rightarrow 0$$

4.1.1 Scaled Example

Define Scaled Variables

$$\Psi \equiv \frac{1}{\pi} \frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}}$$

$$K \equiv \frac{1}{\pi} \int_{\Gamma}^{PV} \frac{\partial}{\partial n_{\vec{x}}} \log \|\vec{x} - \vec{x}'\| \sigma(\vec{x}') d\Gamma'$$

The 2-D Neumann problem becomes

$$(I + K)\sigma = \Psi$$

4.1.2 Key Property

Main assumption of second kind theory:

$$(I + K)^{-1} \text{ is bounded.}$$

Is $(I + K)^{-1}$ bounded for Interior Neumann Problem?

4.2 Linear Algebra

Given $Ax = b$, $A \in \mathfrak{R}^{n \times n}$, $x, b \in \mathfrak{R}^n$

A^{-1} exists and is bounded iff

$$Ay = 0 \text{ implies } y = 0 \text{ (no null space)}$$

If $Ay = 0$ for $y \neq 0$ then either

$Ax = b$ has an infinite # of solutions

$$Ax = b \text{ then } A(x + \alpha y) = b$$

OR

$Ax = b$ does not have a solution

b is not in the column space of A

4.3 3-D Null Space

Consider $\tilde{\sigma}$ defined by

$$u_{\Gamma}(\vec{x}) = 1 = \int_{\Gamma} \frac{1}{\|\vec{x} - \vec{x}'\|} \tilde{\sigma}(\vec{x}') d\Gamma' \quad \vec{x} \in \Gamma$$

Then

$$\frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}} = 0 = 2\pi \tilde{\sigma}(\vec{x}) + \int_{\Gamma} \frac{\partial}{\partial n_{\vec{x}}} \frac{1}{\|\vec{x} - \vec{x}'\|} \tilde{\sigma}(\vec{x}') d\Gamma'$$

$\tilde{\sigma}$ is in the Null space of $I + K$

$$(I + K)^{-1} \text{ is not bounded!!}$$

4.4 Fredholm Alternative

General Theorem

For $I + K$ either

$$(I + K)\sigma = \Psi \text{ has an infinite \# of solutions}$$

OR

$$(I + K)\sigma = \Psi \text{ has no solution}$$

4.4.1 2D Example

Scaled Equations:

$$\frac{1}{\pi} \frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}} = \sigma(\vec{x}) + \frac{1}{\pi} \int_{\Gamma}^{PV} \frac{\partial}{\partial n_{\vec{x}}} \log \|\vec{x} - \vec{x}'\| \sigma(\vec{x}') d\Gamma'$$

For a solution to exist

$$\int_{\Gamma} \frac{\partial u_{\Gamma}(\vec{x})}{\partial n_{\vec{x}}} d\Gamma = 0$$

2D Neumann Second Kind Integral equation
has a one-dimensional Null space.

4.4.2 Fixes

Add a point constraint

Fix u at some point

Force σ orthogonal to null space

Need the null space

May need to solve 1st kind equation

Use SVD to solve singular system

Can be computationally expensive

5 1st Kind Convergence

Three-dimensional Laplace's equation

- Unknowns might be physically meaningful.

$$u(x) = \int \frac{1}{\|\vec{x} - \vec{x}'\|} \underbrace{\sigma(\vec{x}')}_{\substack{\text{charge} \\ \text{density}}} dS'$$

- Might match boundary conditions
 - Dirichlet and $\frac{1}{R}$ radiation condition

5.1 Nonsingular Green's Function

Denote the integral operator as K

$$K\sigma \equiv \int_{-1}^1 |x - x'| \sigma(x') dS' \Rightarrow K\sigma = \Psi$$

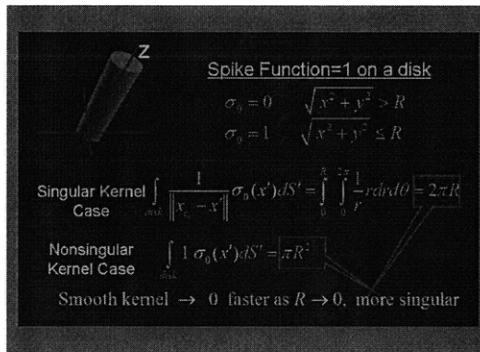
The integral operator is **singular** : K has a null space

$$\sigma_0(x) = 0, x \neq 0, \sigma_0(0) = 1$$

$$K\sigma_0 = \int_{-1}^1 |x - x'| \sigma_0(x') dS' = 0 \tag{2}$$

If $K\sigma^a = \Psi$ then $K(\sigma^a + \sigma_0) = \Psi$

5.2 The Singular Kernel



5.3 Convergence Analysis

Partial Differential Equation form:

$$\begin{aligned} \nabla^2 u &= f & \text{in } \Omega & \quad \Omega \text{ is the volume domain} \\ u &= 0 & \text{on } \Gamma & \quad \Gamma \text{ is the problem surface} \end{aligned}$$

“Nearly” Equivalent weak form

$$\underbrace{\int_{\Omega} \nabla u \nabla v dx}_{a(u,v)} = \underbrace{\int_{\Omega} f v dx}_{l(v)} \quad \forall v \in H^1(\Omega)$$

Introduced an abstract notation for the equation, u must satisfy:

$$a(u, v) = l(v) \quad \forall v \in H^1(\Omega)$$

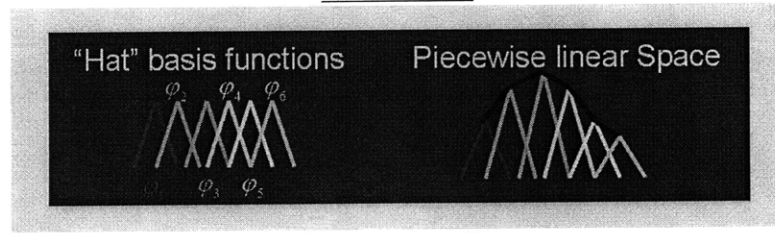
Introduce an approximate solution $u^n = \sum_{i=1}^n \alpha_i \varphi_i$

$\Rightarrow u^n$ is a weighted sum of basis functions

The basis functions define a space

$$X_n = \left\{ v \in X_n \mid v = \sum_{i=1}^n \beta_i \varphi_i \text{ for some } \beta_i \text{'s} \right\}$$

Example



Key Idea

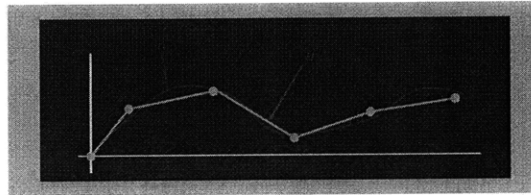
$a(u, u)$ defines a norm on $H_0^1(\Omega)$ $a(u, u) \equiv \|u\|$
 u is restricted to be 0 at 0 & 1!

Using the norm properties, it is possible to show

$$\text{If } a(u^n, \varphi_i) = l(\varphi_i) \quad \forall \varphi_i \in \{\varphi_1, \varphi_2, \dots, \varphi_n\}$$

$$\text{Then } \underbrace{\|u - u^n\|}_{\text{Solution Error}} = \min_{w_n \in X_n} \underbrace{\|u - w^n\|}_{\text{Projection Error}}$$

5.3.1 Optimality Result



How well can you fit the exact solution with a member of X_n ?
 You must measure the error in the $\| \cdot \|$ norm

5.3.2 Sobolov Space

“Weak” Form for the integral equation

$$\underbrace{\iint_{\Gamma} v(x) \frac{1}{\|x - x'\|} \sigma(x') dS' dS}_{a(\sigma, v)} = \underbrace{\int_{\Gamma} v(x) \Psi(x) dS}_{l(v)} \quad \forall v \in \bar{H}(\Gamma)$$

The difficulty is defining $\bar{H}(\Gamma)$ with right properties

Must exclude $\sigma(x)$'s where $\int \frac{1}{\|x - x'\|} \sigma(x') dS' = 0$

$\bar{H}(\Gamma)$ is a fractional Sobolev Space

We won't say more about this!

5.3.3 Use FEM Key Idea

$$\begin{array}{l}
 a(\sigma, \sigma) \text{ defines a norm on } \bar{H}(\Gamma) \\
 \sigma^n = \sum_{i=1}^n \alpha_i \underbrace{\varphi_i(x)}_{\substack{\text{Basis} \\ \text{Functions}}} \\
 a(\sigma, \sigma) = \|\sigma\| \\
 X_n = \left\{ v \in X_n \mid v = \sum_{i=1}^n \beta_i \varphi_i \text{ for some } \beta_i \text{'s} \right\}
 \end{array}$$

5.3.4 FEM Idea Cont.

Using the norm properties, it is possible to show

If

$$a(u^n, \varphi_i) = l(\varphi_i) \quad \forall \varphi_i \in \{\varphi_1, \varphi_2, \dots, \varphi_n\}$$

$$\begin{array}{l}
 \text{Then } \underbrace{\|u - u^n\|}_{\substack{\text{Solution} \\ \text{Error}}} = \min_{w_n \in X_n} \underbrace{\|u - w^n\|}_{\substack{\text{Projection} \\ \text{Error}}}
 \end{array}$$

6 Summary

Interior Neumann

Use Nystrom to Solve

Look at 2-D Problems

Fredholm Alternative

Connection to Linear Algebra

First Kind Convergence Theory

Mostly Waved hands.

1.6 Fast Algorithms for Integral Equation Methods

Numerical Methods for PDEs

Boundary Element Methods, Lecture 6
Fast Algorithms for Integral Equations

L. Proctor, S. De, K. Nabors, J. Phillips, B. Buchmann, & J. White

December 10, 2008

1 Outline

Reasons for Fast Solvers

Collocation System Reminder

Fast Solver General Approach

Using Iterative methods

Fast matrix-vector products

Fast Multipole Algorithms

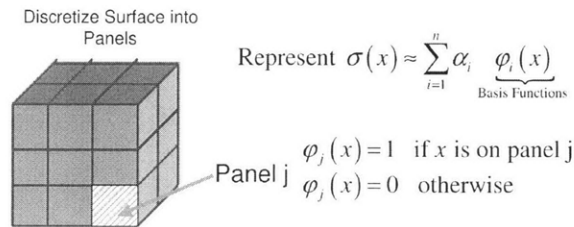
Precorrected-FFT Algorithms

2 Background

2.1 Discretize Surface Into Panels

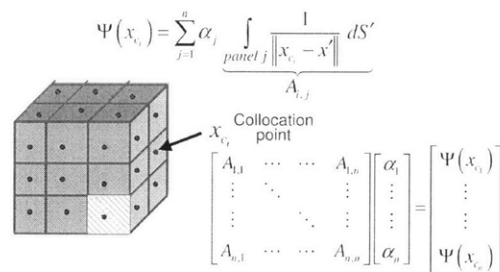
2.1.1 Piecewise Constant Basis

Integral Equation : $\Psi(x) = \int_{surface} \frac{\sigma(x')}{\|x-x'\|} dS'$



2.1.2 Centroid Collocation

Put collocation points at panel centroids



2.2 Dense Matrix

2.2.1 Resultant Dense Matrix

Matrix Entries Are Never Zero

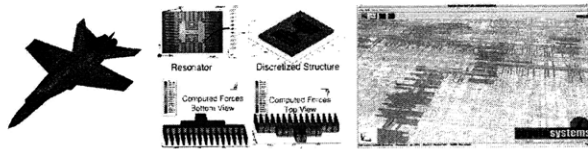
$$A_{i,j} = \int_{panel_j} \frac{1}{\|x - x'\|} dS'$$

Distant Elements Decay Slowly

$$\propto \frac{1}{\|x - x'\|}$$

Too Slow To Ignore.

2.2.2 Complicated Examples



Need More than 100,000 unknowns!!

Need 100 Gigabytes to Store Matrix.

2.2.3 Gaussian Elimination

```

For i = 1 to n-1 {      "For each Row"
  For j = i+1 to n {   "For each Row below pivot"
    For k = i+1 to n { "For each element beyond Pivot"
       $A_{jk} \leftarrow A_{jk} - \frac{A_{ji}}{A_{jj}} A_{jk}$ 
      Multiplier
    }
  }
}

```

A_{ji} Pivot
 A_{jk} Form n-1 reciprocals (pivots)
 Form $\sum_{i=1}^{n-1} (n-i) = \frac{n^2}{2}$ multipliers
 Perform $\sum_{i=1}^{n-1} (n-i)^2 = \frac{2}{3}n^3$ Multiply-adds

n^3 – Too Expensive!

3 Iterative Methods

3.1 Electrostatics Application

General Iterative “Algorithm”

0 : Guess at panel charges $\vec{\alpha}$

1 : Compute the centroid potentials from the charges
 $A\vec{\alpha}$

2 : Compare the computed to known potentials
 $\mathbf{R} = \Psi - \mathbf{A}\vec{\alpha}$

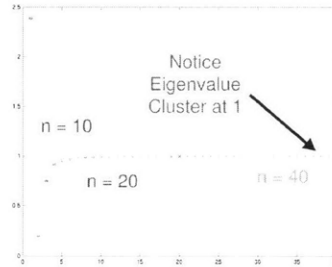
3 : Fix the panel charges, go to Step 1.
 $\vec{\alpha}$

3.2 Conjugate Gradient (CG)

Conjugate Gradient (CG) Methods are iterative methods useful for solving systems of equations involving symmetric matrices $\mathbf{A} = \mathbf{A}^T$. The rate of iteration convergence for CG can be related to the ratio of the maximum to the minimum eigenvalue of A .

3.2.1 CG for 2nd Kind

Eigenvalues for 2nd Kind Integral Equation



3.2.2 CG for 2nd Kind Cont.

Conjugate-Gradient convergence rate

$$\|r^k\| \leq 2 \left(\sqrt{\frac{\lambda_{\max}}{\lambda_{\min}}} - 1 \middle/ \sqrt{\frac{\lambda_{\max}}{\lambda_{\min}}} + 1 \right)^k \|r^0\|$$

For discretized Second Kind equations

$$\frac{\lambda_{\max}}{\lambda_{\min}} \text{ is bounded independent of } n$$

Number of CG iterations independent of n!!

3.2.3 Steps of CG

The k^{th} step of the Conjugate Gradient Algorithm

$\text{compute } Ap_k$ $\alpha_k = \frac{(r^k)^T (Ap_k)}{(Ap_k)^T (Ap_k)}$ $x^{k+1} = x^k + \alpha_k p_k$ $r^{k+1} = r^k - \alpha_k Ap_k$ $p_{k+1} = r^{k+1} - \frac{(Ar^{k+1})^T (Ap_k)}{(Ap_k)^T (Ap_k)} p_k$	<div style="border: 1px solid black; padding: 2px; width: fit-content;">For discretized Integral equations, A is dense</div> <div style="border: 1px solid black; padding: 2px; width: fit-content;">Determine optimal step size in kth search direction</div> <div style="border: 1px solid black; padding: 2px; width: fit-content;">Update the solution and the residual</div> <div style="border: 1px solid black; padding: 2px; width: fit-content;">Compute the new orthogonalized search direction</div>
---	---

3.2.4 Cost of CG

Complexity of the Conjugate Gradient Method

$\text{compute } Ap_k$ $\alpha_k = \frac{(r^k)^T (Ap_k)}{(Ap_k)^T (Ap_k)}$ $x^{k+1} = x^k + \alpha_k p_k$ $r^{k+1} = r^k - \alpha_k Ap_k$ $p_{k+1} = r^{k+1} - \frac{(Ar^{k+1})^T (Ap_k)}{(Ap_k)^T (Ap_k)} p_k$	<div style="border: 1px solid black; padding: 2px; width: fit-content;">Dense Matrix-vector product costs $O(n^2)$</div> <div style="border: 1px solid black; padding: 2px; width: fit-content;">Vector inner products, $O(n)$</div> <div style="border: 1px solid black; padding: 2px; width: fit-content;">Vector Adds, $O(n)$</div> <div style="border: 1px solid black; padding: 2px; width: fit-content;">Inner products, total cost $O(n)$</div>
---	--

Algorithm is $O(n^2)$ for integral equations even though # of iterations, k , is small!

3.2.5 Accelerate CG?

Accelerate the Conjugate Gradient Method

Exactly compute Ap_k

Dense matrix-vector (M-V) product costs $O(n^2)$

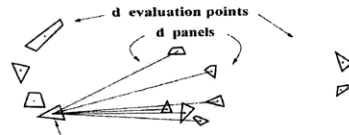
Approximately compute Ap_k

Reduces M-V product costs to $O(n)$ or $O(n \log n)$

Need a fast approximation for matrix-vector products

4 Fast Solvers

4.1 Direct Computation



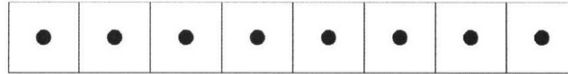
- Physical interpretation:
 $Ap = N$ "potentials" due to N charges.

- $O(N^2)$ if done naively

4.2 1D Strip of Charge in 3D Space

4.2.1 Simplification of the A Matrix

1-D Strip of Charge in 3-D Space



$$\begin{bmatrix} A_{11} & A_{12} & \cdots & A_{18} \\ A_{21} & A_{22} & \cdots & A_{28} \\ \vdots & \vdots & \ddots & \vdots \\ A_{81} & A_{82} & \cdots & A_{88} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_8 \end{bmatrix} = \begin{bmatrix} \Psi_1 \\ \Psi_2 \\ \vdots \\ \Psi_8 \end{bmatrix}$$

What can one say about the **A** matrix?

4.2.2 Properties of A.

The **A** matrix is:

- Symmetric
Panel i exerts the exact same charge on j that j exerts on i
- All the Diagonal Values are the Same

$$A_{ii} = \int_{\text{panel}_i} \frac{1}{\|\vec{x}' - \vec{x}_{c_i}\|}$$

- Each Superdiagonal & Subdiagonal Element is Equal along Its Own Diagonal as Well

4.2.3 More Properties of A.

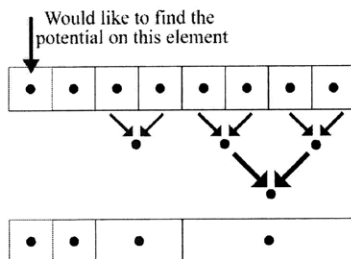
*How many unique entry values are there in **A**?*

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} & \cdots & A_{1n} \\ A_{21} & A_{22} & A_{23} & \cdots & A_{2n} \\ A_{31} & A_{32} & A_{33} & \cdots & A_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ A_{i1} & A_{i2} & A_{i3} & \cdots & A_{in} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ A_{n1} & A_{n2} & A_{n3} & \cdots & A_{nn} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \\ \vdots \\ \alpha_j \\ \vdots \\ \alpha_n \end{bmatrix} = \begin{bmatrix} \Psi_1 \\ \Psi_2 \\ \Psi_3 \\ \vdots \\ \Psi_i \\ \vdots \\ \Psi_n \end{bmatrix}$$

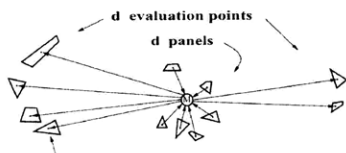
The above matrix is one of a class of matrices called Toeplitz Matrices. In a Toeplitz matrix, the matrix entries along any diagonal have the same value, but the different diagonals can have different values. For this reason, an $n \times n$ Toeplitz matrix has only $2n - 1$ distinct values. A particular important special case of Toeplitz matrices are Circulant Matrices. Circulant matrices are “periodized” Toeplitz matrices in that the first super-diagonal has the same values as the $n - 1$ sub-diagonal, the second super-diagonal has the same values as the $n - 2$ sub-diagonal, etc. Circulant matrices are diagonalized by the discrete Fourier transform, a property we will use in the section covering the Precorrected-FFT methods.

4.2.4 Geometric Simplification

Approximate (by grouping) the elements that are a “reasonable distance” away from the element which you are evaluating

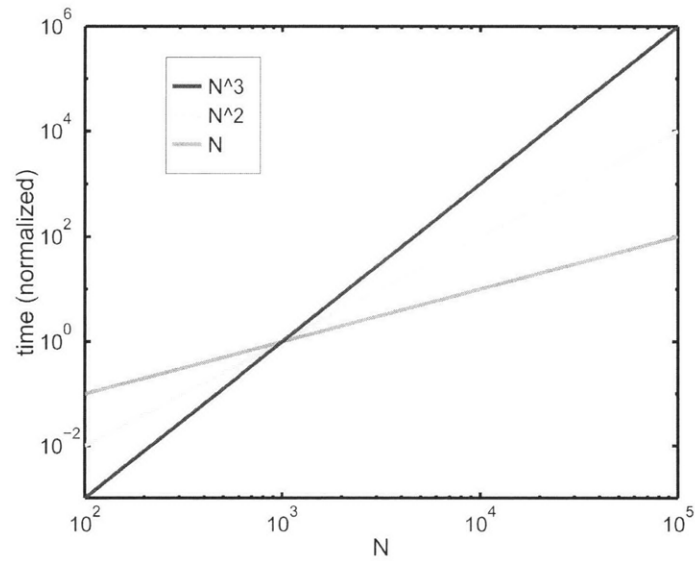


4.3 Fast Potential Concept



- Decompose potential into short- and long- range.
- Approximate long-range part of potential.
- Sum short-range part in normal manner
- Multilevel decomposition for “ $O(N)$ ” algorithm

4.4 Computational Costs

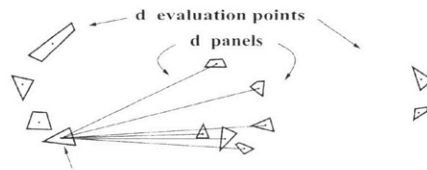


DEC 21164-333		
N	Gaussian Elim	"Fast" $O(N)$
	300 MFLOPS	30 MFLOPS
$5e4$	3 days, 20GB	80sec, 130M
$1e5$	25 days, 80GB	2.5min, 300M
$5e5$	8.8yrs, 2TB	15min, 1.5GB

- Gaussian Elimination: $O(n^3)$ time, $O(n^2)$ memory
- Iterative with direct M-V: $O(n^2)$ time, $O(n^2)$ memory
- Fast Methods: $O(n)$ time, $O(n)$ memory

5 Multipole Algorithms

5.1 Direct Potential Evaluation

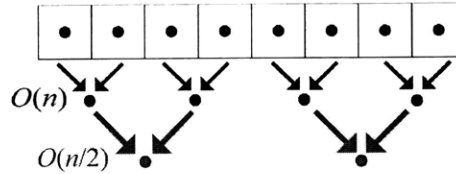


- Potential at point i :
$$v_i(r_i, \phi_i, \theta_i) = \sum_{j=1}^d q_j P_{ij}.$$
- Complete evaluation at d points costs d^2 operations.

5.2 Multipole Representation

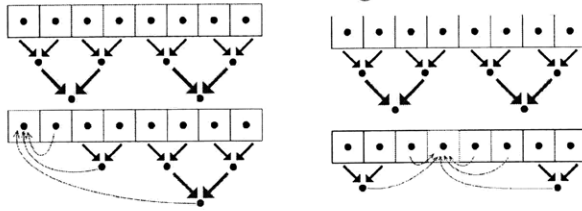
5.2.1 1D Strip in 3D Space

How many operations are needed to form the clusters?



The cost of forming clusters is, in general, $O(n + \frac{n}{2} + \frac{n}{4} + \frac{n}{8} + \dots) \approx O(n)$

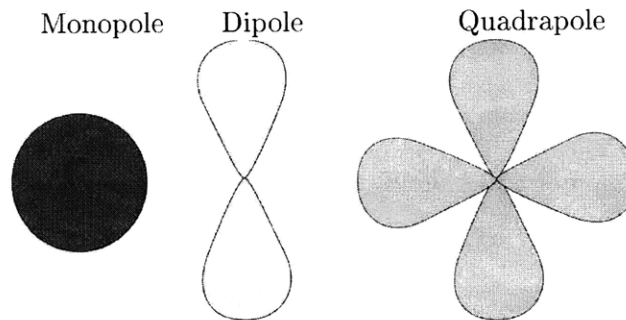
What is the cost of estimating the evaluation point potential?



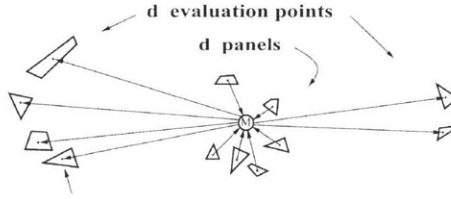
The cost of gathering clusters is $O(n \log n)$

5.2.2 Computational Example

A few multipoles (monopoles, dipoles, quadrupoles, etc) can accurately represent the potential due to a cluster of charges, with the accuracy improving with increasing distance from the cluster. For example, if one is very far from the cluster, the potential due to the cluster will be nearly identical to the potential of a point charge whose location is at the center of the cluster and whose value is the sum of the cluster charges. The accuracy of such a monopole representation can be improved by adding dipole, quadrupole and higher order multipoles. Note, however, that higher order multipoles generate potentials that depend on the multipole's orientation, and that must be considered.



5.2.3 General Case



- Approximate potential at point i :

$$v_i(r_i, \phi_i, \theta_i) \approx \sum_{j=0}^{\text{order}} \sum_{k=-j}^j \frac{M_j^k}{r_i^{j+1}} Y_j^k(\phi_i, \theta_i)$$

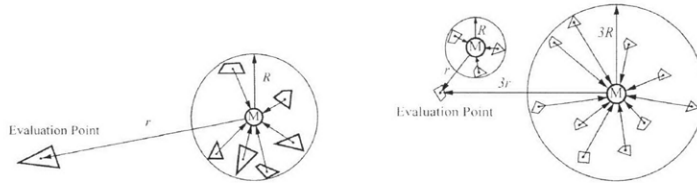
- Multipole coefficients function of panel charges:

$$M_j^k \triangleq \sum_{i=1}^d \frac{q_i}{A_i} \int_{\text{panel } i} \rho^j Y_j^{-k}(\alpha, \beta) dA.$$

- Computing Multipole expansions costs order d operations.
- Each approximate potential evaluation costs order 1 operations.

d potential evaluation due to d panels in order d operations

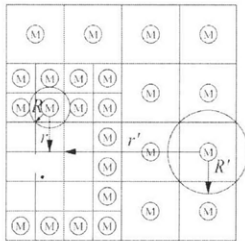
5.3 Error Scale Invariance



$$\text{Error} \leq K \left(\frac{R}{r} \right)^{\text{order}+1}$$

$$\text{Error} \leq K \left(\frac{3R}{3r} \right)^{\text{order}+1}$$

5.4 Multipole Algorithm Hierarchy



Hierarchy guarantees:

- Bounded error:

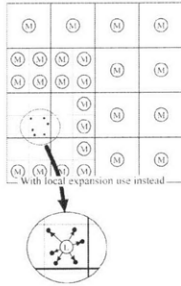
$$\text{Error} \leq K \left(\frac{R}{r} \right)^{\text{order}+1}$$

$$\leq K \left(\frac{1}{2} \right)^{\text{order}+1}$$

$$\text{order} = 2 \text{ yields one percent accuracy.}$$
- Order n ops for n potentials.

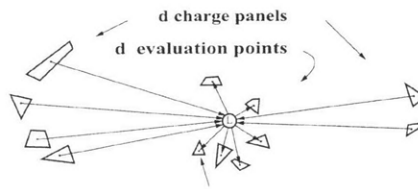
5.5 Local Representations

5.5.1 Cost Reduction



- Construct a local expansion to represent distant charge potentials.
- Evaluate a single local expansion, rather than many multipole expansions, at each evaluation point.

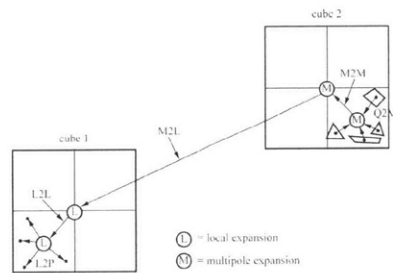
5.5.2 Clustered Evaluations



- Local expansion summarizes the influence of distant charge for clusters of evaluation points.
- Gives $O(n)$ potential evaluation when combined with coalescing of charge done by multipole expansions.

- Approximate potential at point i :
$$v_i(r_i, \phi_i, \theta_i) \approx \sum_{j=0}^{order} \sum_{k=-j}^j L_j^k Y_j^k(\phi_i, \theta_i) r_i^j.$$

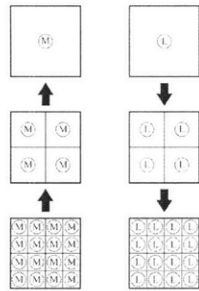
5.5.3 Summary of Operations



- Multiple and local expansions are built using complementary hierarchies.

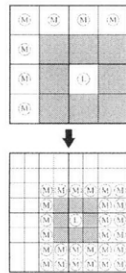
- Complete calculation consists of:
 1. Build multipoles (Upward Pass).
 2. Build locals (Downward Pass).
 3. Evaluate local expansions and nearby charge potential (Evaluation Pass).

5.5.4 Hierarchy Construction



- First build the multipole expansions moving upward from child to parent.
- Then build the local expansions by moving downward from parent to child.
- Computation has a tree structure.

5.5.5 Construction Details

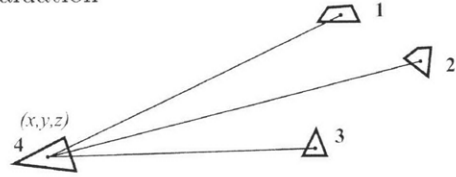


- Conversion of multipole expansions to local expansions.
- A child's local expansion is its parents local expansion plus conversions of multipole expansions in child's interaction range.

5.6 Adaptive Algorithm

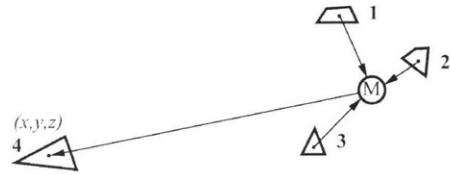
5.6.1 Multipole Inefficiency

Direct Evaluation



$$v_4(x, y, z) = q_1 P_{41} + q_2 P_{42} + q_3 P_{43}$$

Multipole Evaluation



$$v_4(x, y, z) \approx \bar{M}_0^0 \frac{1}{r} + \bar{M}_1^0 \frac{z}{r^3} - \bar{M}_1^1 \frac{x}{2r^3} - \bar{M}_1^1 \frac{y}{2r^3}$$

Using Multipole MORE expensive than Direct.

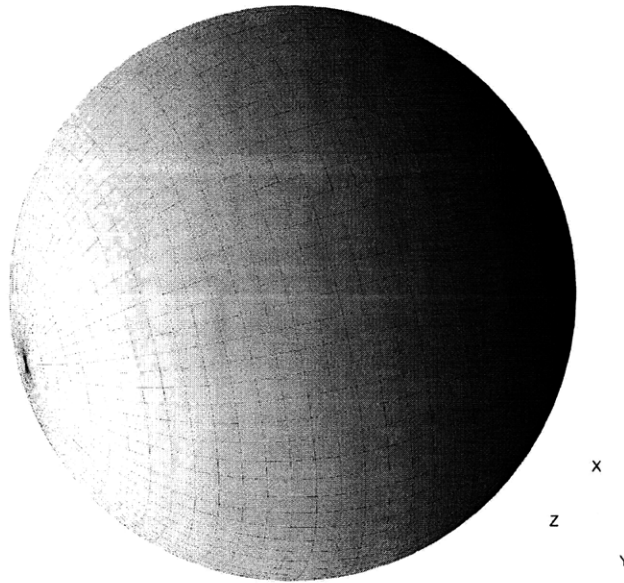
5.6.2 Simple Adaptive Scheme

If there are fewer panels than multipole coefficients, calculate the panels' influence directly.

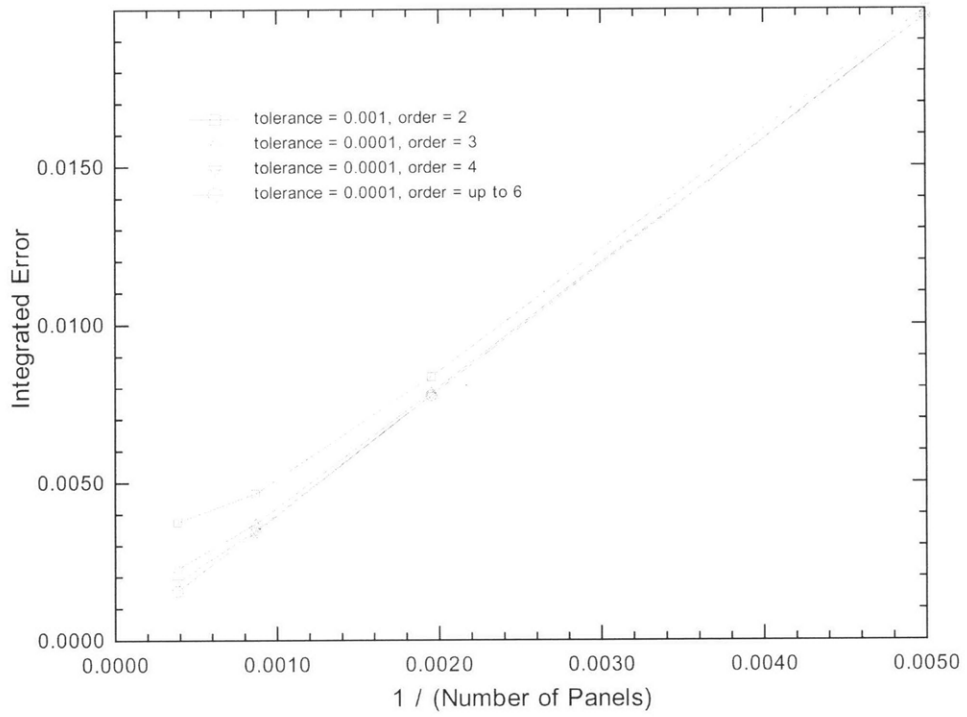
- Similarly, local expansions are not used if there are fewer evaluation points than local expansion coefficients.
- Retains $O(mn)$ complexity for nonuniform panel distributions.

5.7 Computational Examples

5.7.1 Sphere Potential Distribution

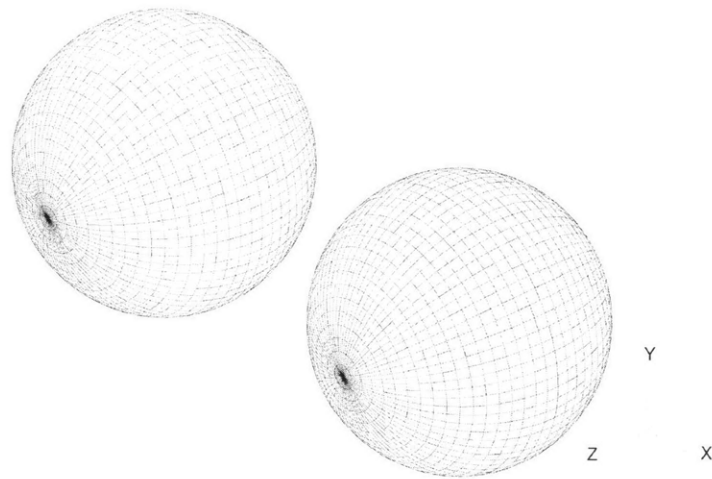


- Potential given by $\psi(x) = -\frac{x_3}{2\|x\|^3}$.
- Charge given by $\sigma(x) = \frac{-3}{8\pi}x_3$.

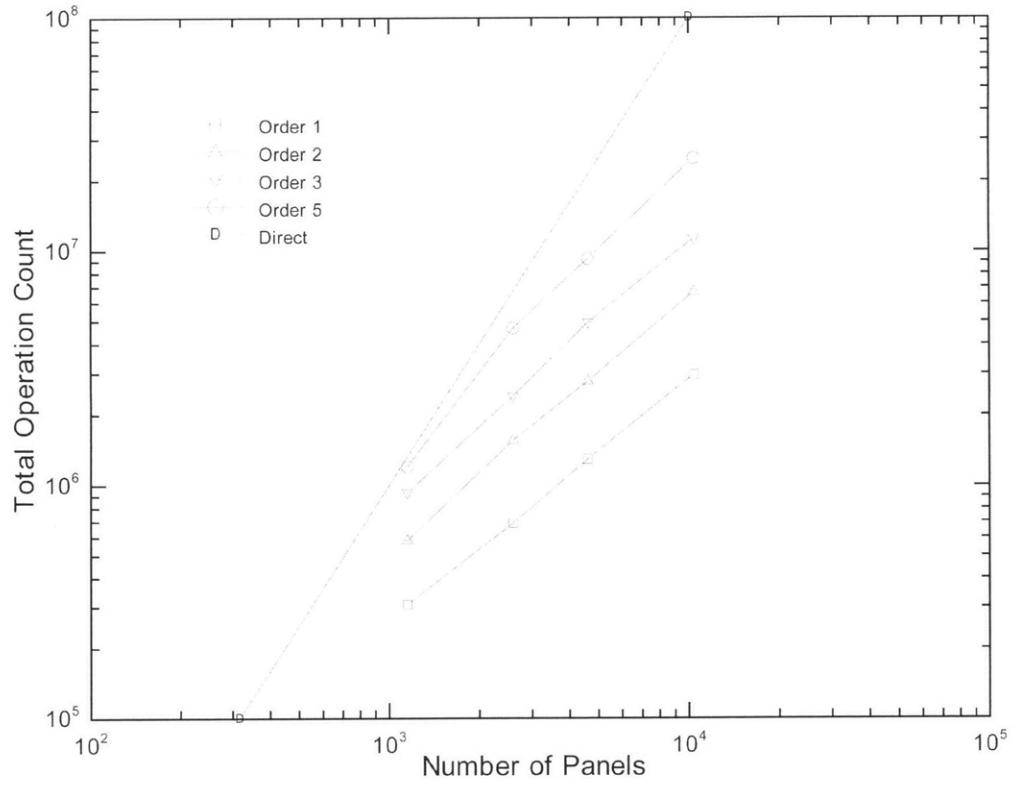


- Error should decay like $\frac{1}{n}$.
- Multipole approximations eventually interfere.
- Higher-order multipole expansions needed for higher accuracy.

5.7.2 Two Sphere Example



- Potential on each sphere: $\psi(x) = -\frac{x_3}{2\|x\|^3}$.
- Does not correspond to a simple physical problem.



- Direct matrix-vector product cost increases like n^2 .
- Multipole matrix-vector product cost increases like n .
- The slope for the multipole algorithm depends on accuracy.
- For order 2 expansions, breakpoint is about $n = 400$.

5.8 Complexity Summary

For an integral equation discretized with n panels:

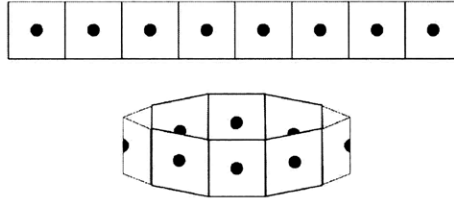
- Gaussian elimination: $O(n^3)$.
- Iterative Matrix Solution, direct M-V $O(n^2)$.
- Multipole accelerated Iterative method $O(n)$.

6 Precorrected-FFT

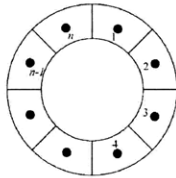
6.1 Introduction

Strip of Charge in Space

Bring the ends of the strip of charge together to form a ring.



Flattening the ring leads to the figure shown below on the left. Forming a ring from the strip of charges results in a system of equations with even more structure than the Toeplitz matrix system described above. The matrix in the ring case will be *circulant*.

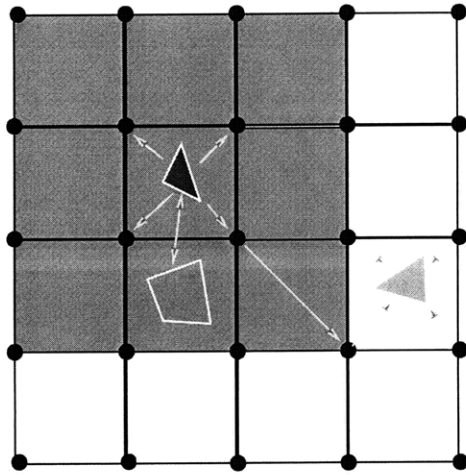


Produces a “**Circulant Matrix**”

$$\begin{bmatrix}
 A_{11} & A_{12} & A_{13} & \dots & A_{1n} \\
 A_{21} & A_{22} & A_{23} & \dots & A_{2n} \\
 A_{31} & A_{32} & A_{33} & \dots & A_{3n} \\
 \vdots & \vdots & \vdots & \ddots & \vdots \\
 A_{n1} & A_{n2} & A_{n3} & \dots & A_{nn}
 \end{bmatrix}
 \begin{bmatrix}
 \alpha_1 \\
 \alpha_2 \\
 \alpha_3 \\
 \vdots \\
 \alpha_j \\
 \vdots \\
 \alpha_n
 \end{bmatrix}
 =
 \begin{bmatrix}
 \Psi_1 \\
 \Psi_2 \\
 \Psi_3 \\
 \vdots \\
 \Psi_j \\
 \vdots \\
 \Psi_n
 \end{bmatrix}$$

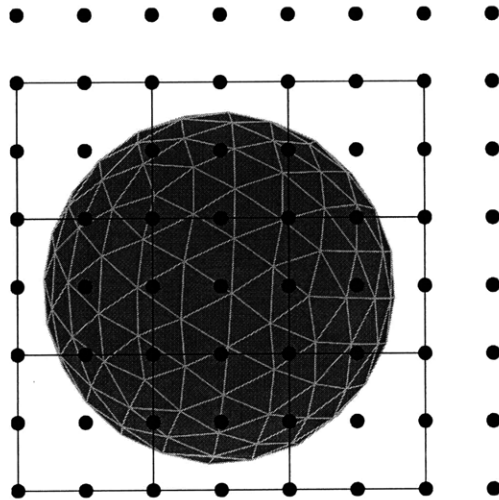
The above circulant matrix is the matrix representation of periodic convolution. This convolutional structure is partly due to the homogeneity of the geometry, and partly because the Green’s function is translation invariant. The Green’s function for Laplace’s equation, $G(x, x')$, is translationally invariant because $G(x, x')$ only depends on the difference, $x - x'$. As mentioned above, the discrete Fourier transform diagonalizes circulant matrices, and therefore circulant matrices can be inverted in $n \log n$ time using the fast Fourier transform.

6.2 Algorithm Outline



1. Project panel charges on grid
2. Calculate grid-charge potentials ϕ
3. Interpolate grid potentials onto p
4. Local corrections
[compute nearby interactions dire

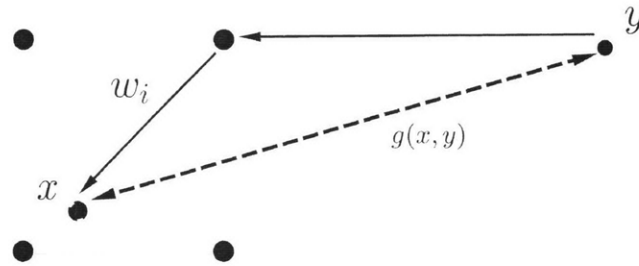
6.2.1 PFFT Grid Balances Costs



- Grid Selected So Direct Cost equals FFT Cost.
- Finer Discretizations Usually Yield Finer Grids.

6.3 Algorithm Analysis

6.3.1 Interpolation and Projection



Approximate potential Ψ at x due to charge at y by *interpolating* potential using points and weights x_i, w_i

Interpolate: potential at x due to unit charge at y

$$\Psi(x|y) \simeq \hat{\Psi}(x|y) = \sum w_i g(x_i, y)$$

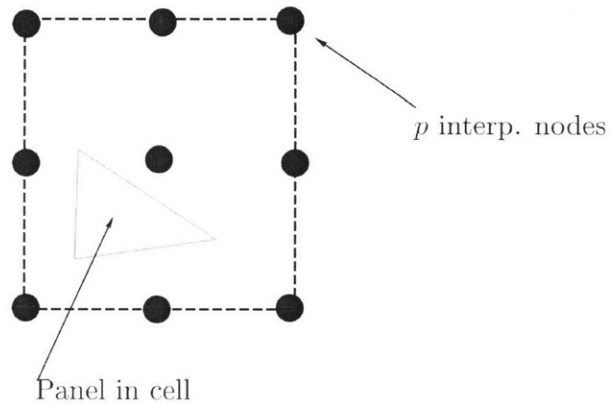
Anterpolate: potential at y due to unit charge at x

$$\Psi(y|x) \simeq \hat{\Psi}(y|x) = \sum w_i g(y, x_i)$$

So

$$\hat{\Psi}(y|x) = \hat{\Psi}(x|y)$$

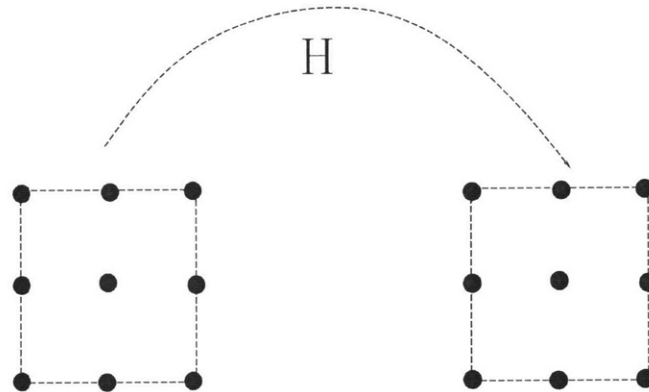
Same as representing charge at x with w_i and evaluating at y



Equivalent conditions:

- Approx Potential in cell due to charge at large R .
- Approx Potential at large R due to charge in cell.
- Cost is $O(N)$

6.3.2 Grid Potentials



- Let H be grid charge-potential mapping

$$H : q_g \rightarrow \Psi_g$$

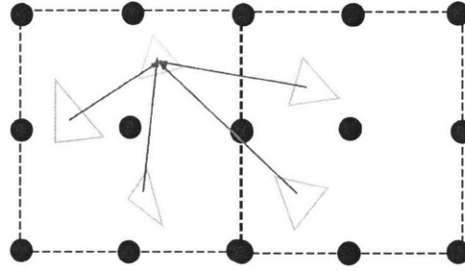
- H is Toeplitz
- Embed H in circulant matrix

$$\begin{bmatrix} \psi_g \\ x \end{bmatrix} = \begin{bmatrix} H & X \\ X & X \end{bmatrix} \begin{bmatrix} q_g \\ 0 \end{bmatrix}$$

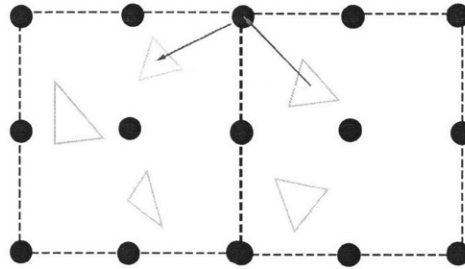
- Use FFT for matrix multiply
Must Have Translation Invariance
- Cost $O(M \log_2 M)$, $M = \#$ cells

6.3.3 Nearby Interactions

Direct interactions
 Cost $O(N[n_c])$
 $[n_c] = \max \# \text{ panels / cell}$



Local corrections
 Cost $O(1) - O(N) - O(Nn_c^2)$



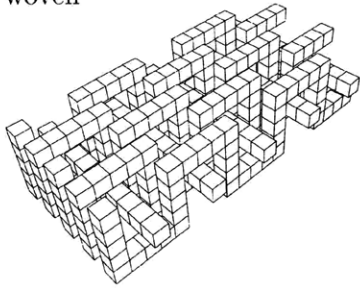
6.3.4 Inhomogeneity Problem



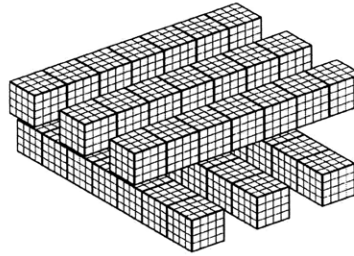
- Empty Grid due to FFT - Inefficient

6.4 Examples

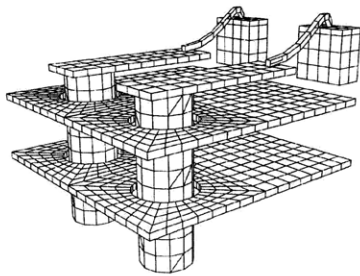
woven



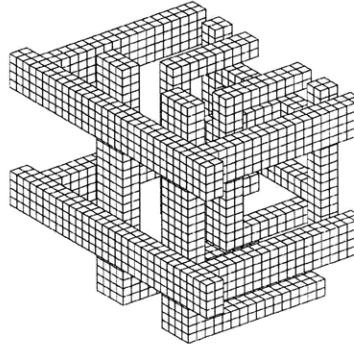
bus3x8



via



sram



6.5 PFFT vs. Multipole

SLIDE 49

Example	CPU	Memory	Product	Error
via	0.61	0.37	0.23	0.18
woven5x5	0.45	0.48	0.22	0.09
cube	0.38	0.32	0.12	0.12
bus3x8	0.27	0.27	0.07	0.01
SRAM	0.39	0.43	0.17	0.07
mean	0.42	0.37	0.16	0.09

- Comparisons: PFFT $p = 3$ to Multi $l = 2$

- Faster with $10\times$ better accuracy !

6.6 PFFT vs. direct

6.6.1 Memory

Example		Memory Usage	
Name	Panels[conductors]	P/FFT	Direct
via	6120[4]	21 Mb	(286 Mb)
woven5x5	9360[10]	50 Mb	(668 Mb)
woven15	82080[30]	246 Mb	(50.2 Gb)
cube	126150[1]	225 Mb	(119 Gb)

6.6.2 Time

Example	CPU Usage		
Name	P/FFT	Dir. Iter.	Gauss. Elim.
via	1.1 min	(5.6 min)	(1.9 hrs)
woven5x5	5.2 min	(42 min)	(6.9 hrs)
woven15	1.7 hrs	(11.5 days)	(194 days)
cube	3.3 min	(8.4hrs)	(2.7 yrs)

7 Summary

Reasons for Fast Solvers

Collocation System Reminder

Fast Solver General Approach

Using Iterative methods

Fast matrix-vector products

Two Fast Methods

Fast Multipole - Multiresolution

Precorrected-FFT - Translation Invariance

Chapter 2

Foundations of Algorithms and Computational Techniques in Systems Biology

Engineering has always played a role in biology, specifically in the past couple of decades the field of computational biology has emerged and contributed greatly. Foundations of Algorithms and Computational Techniques in Systems Biology is a course that gives an overview of topics of interest to a computational biologist. The course covers protein modeling, modeling networks, and image processing. These are the top three areas in computational biology, and this course shows how one may use computational techniques to solve various problems with a biological application. This is very interesting to both the biologist and the computer scientist.

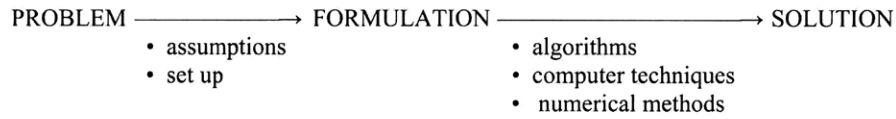
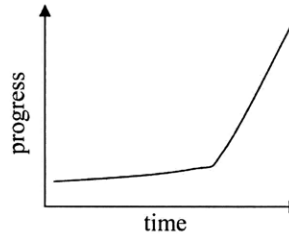
2.1 Motivation/Overview

MIT 6.581/BE.482
 FOUNDATIONS OF ALGORITHMS AND COMPUTATIONAL
 TECHNIQUES IN SYSTEMS BIOLOGY
 Spring 2006

7 February 2006
 Tuesday

MOTIVATION/OVERVIEW

There is a disconnect between biology and computer science.
 The biologist will pose the problem statement, but it may not be amenable for the computer scientist to solve it.
 There is a need for scientists who possess the breadth of knowledge to marry the two realms.



ecologies

populations

individuals

⇕

organ systems

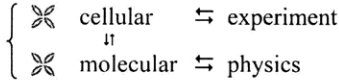
⇕

organs

⇕

tissues

⇕



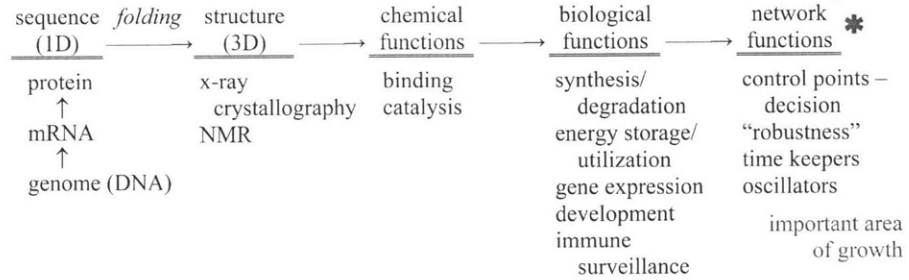
focus of this course

MOLECULAR LEVEL (atoms)	CELLULAR LEVEL (concentration of biomolecules)	IMAGING	
✓		✓	Fast Fourier Transform
✓	✓		Combinatorial Search
	✓		Model Reduction
✓	✓	✓	Singular Value Decomposition
✓			Multipole Algorithm
✓	✓	✓	Numerical Differentiation
✓	✓	✓	Optimization
✓	✓		Newton Methods

PHYSICAL, CHEMICAL, & BIOLOGICAL MODELING OF PROTEINS

Proteins:

- biological polymers of about 20 amino acids
polymers are any kind of large molecules made of repeating identical or similar subunits called monomers
- “perfect” homogeneous, pure synthesis
- around 10k copies in a cell
- linear, unbranched chains of a unique sequence
- generally fold to characteristic structure with no additional information



Why Model?

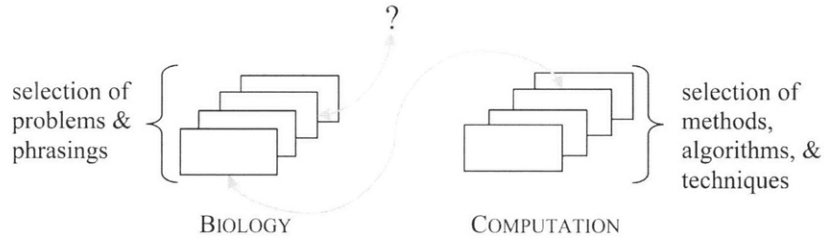
- Understanding : model facilitates development of understanding reason for properties
 - mechanistic basis for function
 - disease
- Prediction
 - experiment planning
 - validate a model or select among models
- Design
 - perturbation : improve properties
 - intervention : repair

2.2 Models of Proteins

MIT 6.581/BE.482
 FOUNDATIONS OF ALGORITHMS AND COMPUTATIONAL
 TECHNIQUES IN SYSTEMS BIOLOGY
 Spring 2006

9 February 2006
 Thursday

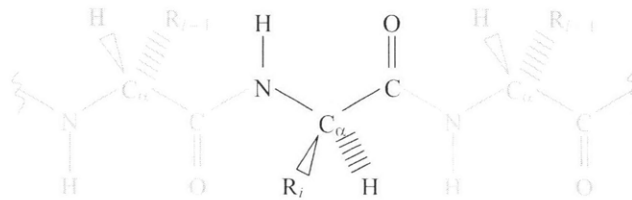
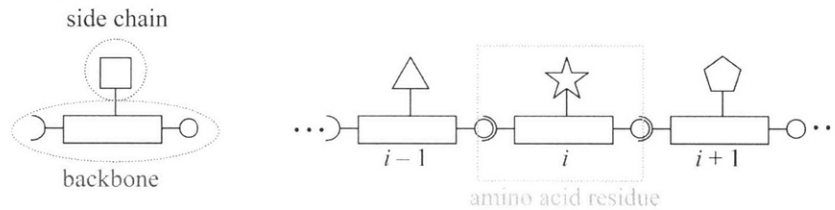
LECTURE 2 : MODELS OF PROTEINS



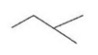
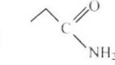
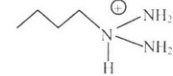
Fundamental role of models:

Understanding
 Prediction
 Design

DNA - mRNA - Protein
 (genome)



The " R_i " groups are chosen from the common 20 amino acid side chains – chemical diversity

- (1) size: small – large
 $R_{Gly} : -H \rightarrow R_{Trp}$
- (2) polarity: hydrophobic – polar – charged
 $R_{Leu} :$  – $R_{Asn} :$  – $R_{Arg} :$ 
- (3) uniformity of character
- (4) local backbone flexibility
 Gly (flexible) – Pro (rigidity)

Coordinate systems:

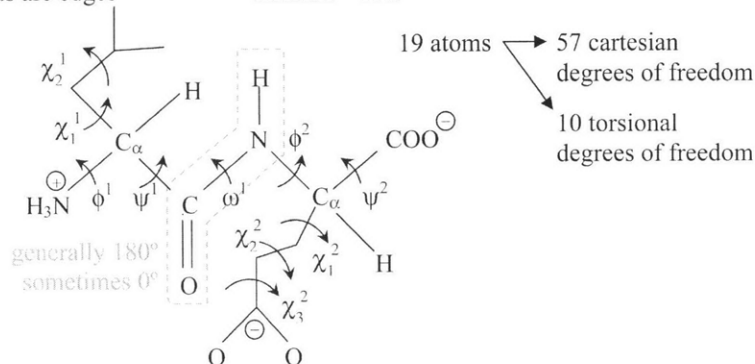
1) Absolute Cartesian Coordinates

$$\vec{X} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ \vdots \\ x_{3N-1} \\ x_{3N} \end{bmatrix} \left. \begin{array}{l} \text{cartesian coordinates of 1}^{\text{st}} \text{ atom} \\ \\ \\ \\ \\ \\ \\ \text{N}^{\text{th}} \text{ atom} \end{array} \right\}$$

2) Relative Coordinates – Internal

Think of the molecules as graphs where

- atoms are vertices bond lengths & bond angles – rigid
- bonds are edges torsions – soft



Desire : Mapping $\vec{X}^{3N} \rightarrow E(\vec{X}^{3N})$
 "energy"
 scalar value

⇒ Bias toward mechanistic basis for model

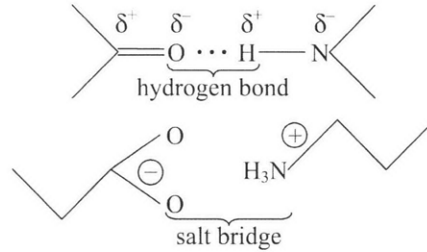
Chemistry – Physics (Quantum Mechanics)

Schrödinger Equation: $i\hbar \frac{\partial \Psi}{\partial t} = -\frac{\hbar^2}{2m} \nabla^2 \Psi + V(x) \Psi(x,t) \equiv \hat{H} \Psi(x,t)$
 nuclear & electrons

Linus Pauling

Observations:

- bond lengths, angles – fixed
- torsions – “soft” & sinusoidal
- atoms appear to have a fixed spherical size & approach to contact neighbors
- complementary electrostatics



Molecular Mechanics Potential:

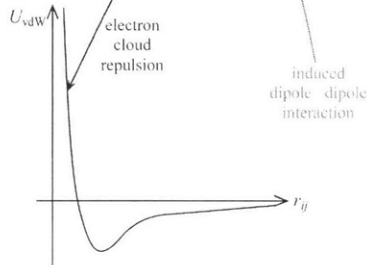
$$E(\vec{X}^{3N}) = U_{\text{COVALENT}} + U_{\text{NON-COVALENT}}$$

\uparrow bonded \uparrow through space

$$U_{\text{COVALENT}} = \sum_{i:\text{bonds}} \frac{1}{2} k_{b,i} (b_i - b_{o,i})^2 + \sum_{i:\text{angles}} \frac{1}{2} k_{\theta,i} (\theta_i - \theta_{o,i})^2 + \sum_{i:\text{impropers}} \frac{1}{2} k_{\phi,i} (\Phi_i - \Phi_{o,i})^2 + \sum_{i:\text{torsions}} \frac{1}{2} k_{\phi,i} [1 + \cos(n_i \phi_i - \delta_i)]$$

$$U_{\text{NON-COVALENT}} = \sum_{i>j} \left(\frac{B_{ij}}{r_{ij}^{12}} - \frac{C_{ij}}{r_{ij}^6} \right) + \sum_{i>j} \frac{q_i q_j}{\epsilon r_{ij}}$$

van der Waals → Lennard-Jones Electrostatics → Coulombic



2.3 Discrete Conformational Search

6.001 / DE 90Z
LECTURE 3: DISCRETE CONFORMATIONAL SEARCH

THURSDAY
14 FEB. 2006

REVIEW:

A) COORDINATE REPRESENTATION - PROTEIN STRUCTURE

- 1) CARTESIAN - ABSOLUTE
- 2) INTERNAL - RELATIVE

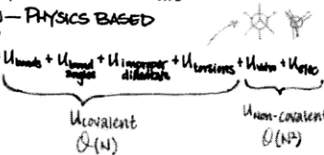
- bond lengths } "hard"
- bond angles } "hard"
- torsional } "soft"

$$X^{atom} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_{3N-2} \\ x_{3N-1} \\ x_{3N} \end{bmatrix}$$

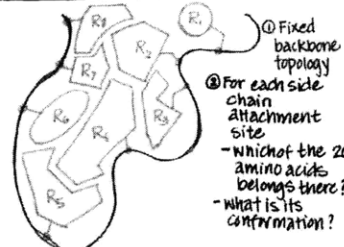
1st atom
Nth atom

B) ENERGY FUNCTION - PHYSICS BASED

- level of "atom"
- $X^{atom} \rightarrow E(X^{atom}) = U_{bonds} + U_{bond\ angles} + U_{improper\ dihedral} + U_{torsions} + U_{vdw} + U_{elec}$



PROBLEM STATEMENT:



- ① Fixed backbone topology
- ② For each side chain attachment site - Which of the 20 amino acids belongs there? - what is its conformation?

Our energy function

- adopts relaxed local (covalent) conformations
 - vdw potential: efficient packing of space
 - elec: complementary charge patterns
- 2 simplifications**
1. fix all bond lengths & bond angles
 2. discretize side chain torsions

WHAT CAN WE DO WITH AN ENERGY FUNCTION?

- * Fold a protein
 - Find the final, folded shape
 - Understand how a protein finds its final shape
- Average folding time for a protein: ~1ms - is typical
- Most detailed simulations: 1cpu day @ 1ms
- 1ms $\Rightarrow 10^6$ CPU days $\Rightarrow 2,740$ years

Two Studies of Historical Note:

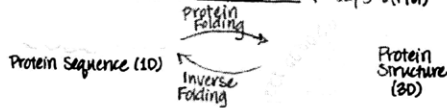
- 1) Levitt & Warshel *Nature* 253: 694-698 (1975).
 - Simplified model of a small protein
 - BPTI 58 residues
 - Every center represents an amino acid
 - ran detailed atom-level simulations on pairs of amino acids to produce higher level model
 - Accelerated Simulation Methodology
 - RESULT**
 - Authors claim close to "native" structure
 - Subsequent analysis question this
- 2) Duan & Kollman *Science* 282: 740-744 (1998).
 - 1 μ s simulation of 36-residue protein (villin headpiece subdomain)
 - model: all-atoms plus ~3,000 water molecules
 - took ~4 months on a 256-processor Cray T3E
 - RESULT**
 - Marginally stable state "close" to native state adopted at ~150ns, persisted briefly, then dissolved

FUNDAMENTAL PROBLEM

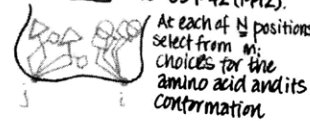
1. What result does the model produce?
 - Interplay: Simulation Methodology \leftrightarrow Model Pathway \leftrightarrow Final Structure
2. Can't relate problems to individual terms in the energy function

Two Visionary Statements:

- c.o. Pabo *Nature* 301: 200 (1983)
- E. Drexler *Proc Natl Acad Sci USA* 78: 5275-8 (1981)



The Dead-End Elimination Algorithm (DEE)
Desmet et al. *Nature* 356: 639-42 (1992).



Solve for $\vec{m} = \begin{bmatrix} m_1 \\ m_2 \\ m_3 \\ \vdots \\ m_N \end{bmatrix}$ that minimizes $E(\vec{m})$ over all possible \vec{m}

Imagine $m_i = m \Rightarrow$ space is m^N

NOTE: OUR ENERGY FUNCTION HAS SIMPLE FORM

$$E_{\vec{m}} = \sum_{i=1}^N E_{m_i}^{self} + \sum_{i=1}^N \sum_{j=i+1}^N E_{m_i, m_j}^{pair} + const.$$

terms between SC & BB SC-SC

Looking at just one position:

$$E_{m_i}^{self} + \sum_{j=1}^N E_{m_i, m_j}^{pair}$$

is the contribution of I choose m_i at position i

But I can bound the contribution

Lower Bound: $E_{m_i}^{self} + \sum_{j=1}^N \min_k E_{m_i, m_j}^{pair}$

Upper Bound: $E_{m_i}^{self} + \sum_{j=1}^N \max_k E_{m_i, m_j}^{pair}$

DEE Singles Criterion

$$\text{If } E_{m_i}^{self} + \sum_{j=1}^N \min_k (E_{m_i, m_j}^{pair}) > E_{m_i}^{self} + \sum_{j=1}^N \max_k (E_{m_i, m_j}^{pair})$$

For some tps, then m_i can not be in the global optimum

this method is probably correct, but can't prove remaining time

Approach:

1. Eliminate iteratively singles until no more possible
2. Develop higher-order eliminations (pairs) $\rightarrow \uparrow$
3. When all eliminations done, enumerate remaining space

Application:

Insulin: 76 positions: 2.7×10^{76} conformations
iterative approach (9 iterations)

$2.7 \times 10^{76} \rightarrow 7200 \rightarrow$ search by enumeration

\downarrow
93% of buried positions "correct"

2.4 Binding and Docking

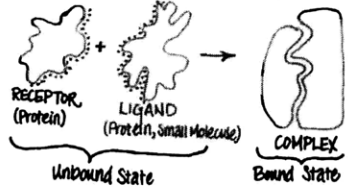
16 FEB 06

LECTURE 4: BINDING & DOCKING

END OF PREVIOUS LECTURE

FORWARD
model - simple
model - ???
REVERSE
model - ???

- Doherty & Mayo Science 278: 82-87 (1997)
Known protein ZIF-268
At each position allowed the actual & other compatible
Amino acids each in multiple rotomers
Sequence complexity (S)ⁿ ~ 1.9 x 10¹⁷
Structural complexity (M)ⁿ ~ 1.1 x 10¹⁸
DEE → 90 CPU hours → Global Optimum
Synthesize corresponding protein & characterized
- adopted the same fold as the natural sequence
- had cooperative transition

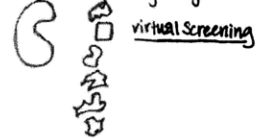


Prediction of Bound State from Unbound

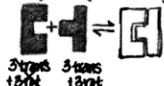
- 1) Understanding
- 2) Predict binding specificity & partnering
phosphorylated protein phospho-binding proteins



- 3) Drug discovery & drug design



Simplify to Rigid Binding

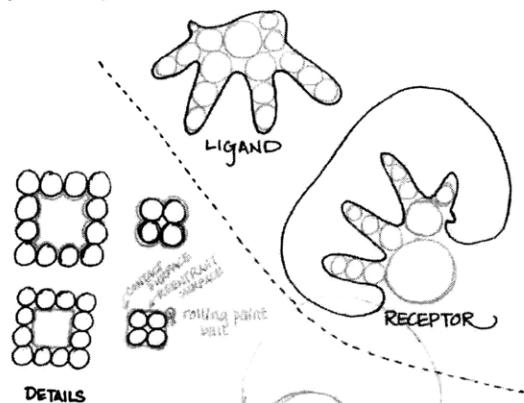


6 dof size of a molecule
0.1 Å grid across 20 Å → 200 grid pts per dimension
→ 1.9 x 10⁶ poses

- IO Kuntz et al J. Mol Biol 161: 269-288 (1992)

- 3 steps:
- 1) Representation - use few abstraction
Ligand - (positive) - spheres
Receptor - (negative) - spheres
 - 2) Matching
- recognize similar features in ligand & receptor

3) Fitting

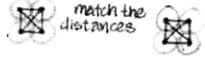


DETAILS

- Construct Sphere sets
1. Sphere touches surface at 1 pt
 2. sphere center lies along the surface normal from 1
 3. Receptor spheres are drawn "outside" surface; Ligand spheres "inside"
- Then produce reduced representation
- at each point, there are (n-1) spheres
 - retain the smallest - removing spheres that cross surface
 - preference for $\theta < 90^\circ$
 - keep only one sphere per atom
 - largest from contact points on receptor (convex)
 - largest from reentrant pts on ligand (concave)

First Result: overlapping receptor spheres, tend to represent binding sites (known)

Sphere pattern in ligand ↔ Sphere pattern in receptor



2.5 Binding and Docking - Molecular Dynamics Simulation

6.561/BE 402
LECTURE 5: BINDING & DOCKING
MOLECULAR DYNAMICS SIMULATION

THURSDAY
23 FEB 2006

BINDING & DOCKING

R + L = C

FRIZZ FACTOR FOR VOXELS → WANT TO LOOK AT ALL NEIGHBORS

GEOMETRIC HASHING Fischer et al J Mol Biol 246:469-477 (1995)

BASIC STRATEGY

- 1) Molecular Surface: Ligand & Receptor
- 2) Identify "critical points" on each surface
 - Centers of convex patches on ligand
 - Centers of concave patches on receptor
 } & surface normals
- 3) Create total orientations
 - In each orientation, line up 2 critical points & their avg. surface normal from ligand with 2 more & avg. surface normal from receptor
 - In each orientation, count number of other critical point coincidences
- 4) Collect high scorers (based on geometry alone) & evaluate with energy function

TIMING: $2 \binom{m}{2} \cdot \binom{n}{2} \sim m^2 n^2$ trial orientations
 $m \cdot n \rightarrow$ # of distances to measure in each orientation
 $\sim m^2 n^2$ total \rightarrow quite bad

EFFICIENT IMPLEMENTATION

- 1) For each pair of critical points on ligand, transform into "standard orientation"
 - mark positions of other critical points in space & store in hash table
 - index by voxel of critical point \rightarrow store the current ligand transformation
- 2) For each pair of critical points on the receptor, carry out the complementary transformation
 - look up each receptor critical points voxel in hash table. Add a vote for each ligand transformation with a critical point in that voxel.
 - tally votes for each ligand transformation
 - pass high scorers on to energy function

TIMING: 1) $\binom{m}{2}$ pairs of ligand points, $O(m)$ operations on each $\rightarrow m^3$
 2) $2 \binom{n}{2}$ pairs of receptor points & $O(n)$ operations on each $\rightarrow \sim n^3$
 $\sim m^3 + n^3$ Total

ASK Questions like:

- Extract Diffusion Constants
- or flow rates, density
- Temperature, Pressure
- Ensemble Averages

Extraction of Continuum Properties

Important Forces - Non-bonded

Lennard-Jones: $\sum_{i>j} \frac{B_{ij}}{r_{ij}^{12}} - \frac{C_{ij}}{r_{ij}^6}$ ← Potential functions

Electrostatics: $\sum_{i>j} \frac{D_{ij}}{r_{ij}}$ ← Potential functions

Look at VMD simulator

one or two molecules interacting (with water or salt)

QUESTIONS:

- Conformational Changes
- Different foldings
- How hard are they to pull apart

More Important Forces:

(bonded/covalent)

Force between A_1 & A_2 : $k_b \|\vec{r}_1 - \vec{r}_2\|$
 Spring potential energy: $\frac{1}{2} k_b (\|\vec{r}_1 - \vec{r}_2\| - d_{min})^2$ (don't worry)

Molecular Dynamics

Force on atom $i = -\frac{\partial}{\partial \vec{r}_i} E(\vec{r})$ (all atom positions)

$F(\vec{r}) = -\nabla_{\vec{r}} E(\vec{r})$

HAMILTONIAN POINT OF VIEW

$H(q, p) = \sum_i \frac{1}{2} m_i \dot{p}_i^2 + E(q)$

position (m, p) kinetic energy

Potential Energy $E(q)$ bonds, angles, vdw, etc.

$\frac{d}{dt} \vec{q} = M^{-1} \vec{p} = \frac{\partial}{\partial \vec{p}} H(q, p)$

$\frac{d}{dt} \vec{p} = -\frac{\partial}{\partial \vec{q}} H(q, p) = -\frac{\partial}{\partial \vec{q}} E(q)$

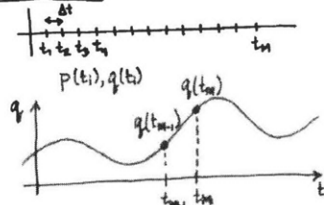
Invariant: $H(p, q) = \text{const.}$

Material Science approach

track thousands of molecules

THURSDAY
25 FEB 06

DISCRETE TIME



$$\frac{d}{dt} q \approx \frac{q(t_n) - q(t_{n-1})}{\Delta t}$$

FORWARD EULER

$$\begin{cases} q(t_{n+1}) = q(t_n) + \Delta t \frac{\partial}{\partial q} H(q(t_n), p(t_n)) \\ p(t_{n+1}) = p(t_n) + \Delta t \frac{\partial}{\partial p} H(q(t_n), p(t_n)) \end{cases}$$

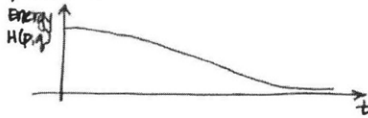
known

$$\begin{cases} q(t_n) = q(t_n) + \Delta t H(q(t_n), p(t_n)) \\ p(t_n) = \dots \end{cases}$$



BACKWARD EULER

$$\begin{cases} q(t_{n+1}) = q(t_n) + \Delta t \frac{\partial}{\partial q} H(q(t_{n+1}), p(t_{n+1})) \\ p(t_{n+1}) = \dots \end{cases}$$



2.6 Molecular Dynamics and Electrostatics

@581/BE482
LECTURE 6: MOLECULAR DYNAMICS & ELECTROSTATICS

TUESDAY
28 FEB 2006

OVERVIEW:

- 1) Remind about E(x)
- 2) Time Discretization
- 3) Electrostatics Evaluation
- COSTS LOW

$$E(x) = \sum_{\text{bonds}} + \sum_{\text{bond angles}} + \sum_{\text{torsions}} + \sum_{\text{all pairs}} \frac{B_{ij}}{r_{ij}^2} - \frac{C_{ij}}{r_{ij}^6} + \sum_{\text{all pairs}} \frac{q_i q_j}{r_{ij}}$$

problematic term

$$m_i \frac{d}{dt} v_{i,x,y,z} = - \frac{\partial}{\partial x_i, y_i, z_i} E(\vec{x})$$

↑
ith atomic mass

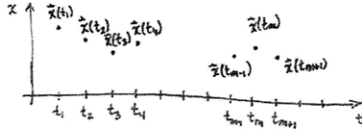
$$M \frac{d}{dt} \vec{V} = - \nabla_x E(\vec{x})$$

$$\frac{d}{dt} \vec{x} = \vec{V}$$

atomic positions atomic velocities

System of 2nd order Equations

$$M \frac{d^2}{dt^2} \vec{x} = - \nabla_x E(\vec{x})$$



$$\frac{d}{dt} \vec{x}(t_m) \approx \frac{\vec{x}(t_{m+1}) - \vec{x}(t_m)}{t_{m+1} - t_m} = \frac{\vec{x}(t_m) - \vec{x}(t_{m-1})}{t_m - t_{m-1}}$$

$$\Delta t = t_{m+1} - t_m = t_m - t_{m-1}$$

$$\frac{d^2}{dt^2} \vec{x}(t_m) \approx \frac{\vec{x}(t_{m+1}) - \vec{x}(t_m)}{\Delta t} - \frac{\vec{x}(t_m) - \vec{x}(t_{m-1})}{\Delta t}$$

$$\frac{d}{dt} \vec{v}(t_m) \approx \frac{\vec{x}(t_{m+1}) - 2\vec{x}(t_m) + \vec{x}(t_{m-1}))}{\Delta t^2}$$

verlet

$$M \frac{\vec{x}(t_{m+1}) - 2\vec{x}(t_m) + \vec{x}(t_{m-1}))}{\Delta t^2} = - \nabla_x E(\vec{x})$$

ALGORITHM:

$$\vec{x}(t_{m+1}) = 2\vec{x}(t_m) - \vec{x}(t_{m-1}) - \Delta t^2 M^{-1} E(\vec{x})$$

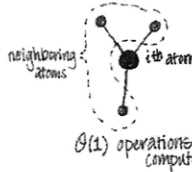
explicit integration scheme

on each time step:

Evaluate $E_{(x)}(\vec{x}(t_m))$ + other stuff

BOND CONTRIBUTIONS:

independent ~ O(N)



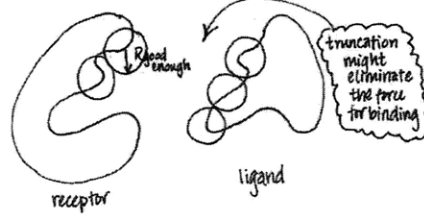
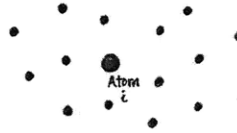
$$E(x) = \sum_{\text{bonds}} + \sum_{\text{bond angles}} + \sum_{\text{torsions}} + \sum_{\text{all pairs}} \frac{B_{ij}}{r_{ij}^2} - \frac{C_{ij}}{r_{ij}^6} + \sum_{\text{all pairs}} \frac{q_i q_j}{r_{ij}}$$

van der Waals:

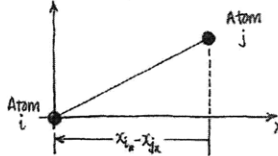


all pairwise interactions
O(N) operations per atom
but truncate
⇒ O(1)

Electrostatics:

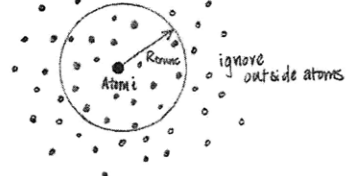


electrostatic force

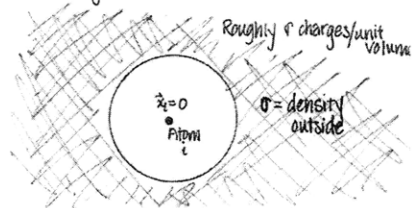


$$F = E_x = - \frac{\partial}{\partial x_i} \frac{q_i q_j}{|\vec{x}_i - \vec{x}_j|} = \frac{q_i q_j (x_i - x_j)}{|\vec{x}_i - \vec{x}_j|^3}$$

try truncation:



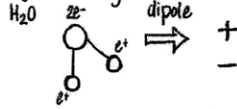
idea of ignored force:



bound on ignored force: ignoring direction cancellation

$$\begin{aligned} \text{Force Bound} &\approx \int_{\text{Volume Outside}} \frac{q(\vec{r})}{|\vec{r}|^2} dV \\ &\approx \int_0^{2\pi} \int_0^\pi \int_{R_{\text{min}}}^{R_{\text{max}}} \frac{1}{R^2} R^2 \sin\theta dR d\theta d\phi \\ &\approx \int_0^{2\pi} \int_0^\pi [R]_{R_{\text{min}}}^{R_{\text{max}}} \sin\theta d\theta d\phi \rightarrow \infty \end{aligned}$$

charge neutrality:



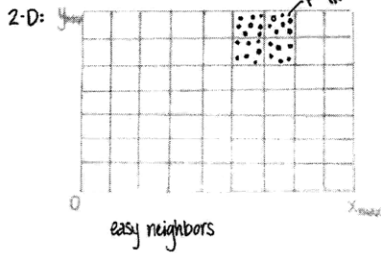
Potential due to dipole: $\frac{\vec{p} \cdot (\vec{x}_i - \vec{x}_j)}{|\vec{x}_i - \vec{x}_j|^3}$
dipole moment

force due to dipole $\propto \frac{1}{R^3}$

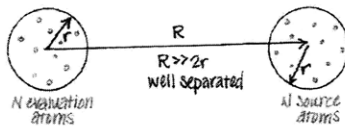
bound on ignored force (dipole assumption):

$$\begin{aligned} \text{Force bound} &\approx \int_{\Omega} \frac{\rho(\vec{r})}{R^3} dV \quad \leftarrow \text{dipole density} \\ &\approx \int_0^{2\pi} \int_0^\pi \int_{R_{\text{min}}}^{R_{\text{max}}} \frac{\rho}{R^3} R^2 \sin\theta dR d\theta d\phi \\ &\approx \int_0^{2\pi} \int_0^\pi [\rho \ln R]_{R_{\text{min}}}^{R_{\text{max}}} \sin\theta d\theta d\phi \rightarrow \infty \end{aligned}$$

STILL INFINITY!

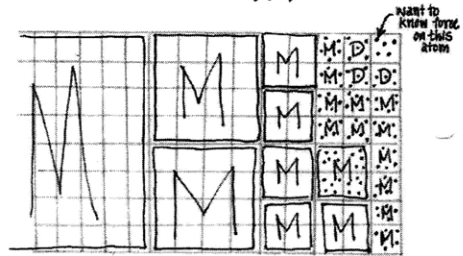
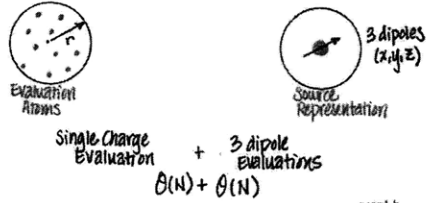
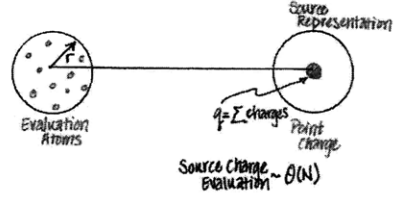


MULTIPOLE IDEA:



direct force calculation is $\Theta(N^2)$

for each of N evaluation atoms: $F_i = \sum_{j \in \text{source atoms}} F_{ij}$

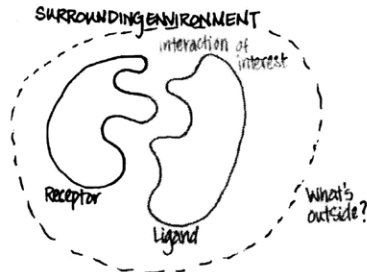


Multiresolution $\Theta(N \log N)$

2.7 Continuum Electrostatic Modeling I

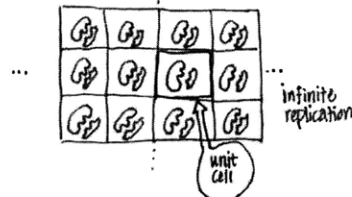
4.001 / DS 907
LECTURE 7: CONTINUUM ELECTROSTATIC MODELING I

11 MAR 2006
2 MARCH 2006

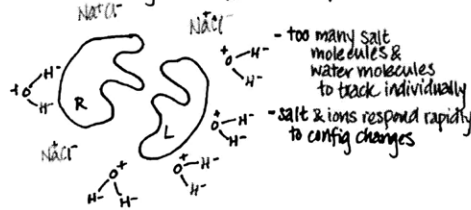


TWO APPROACHES:

1) PERIODIC ARRAY - "CRYSTALS"



2) SURROUNDING BY WATER OR SALT



TODAY: Look at Approach 1 - infinitely periodic unit cell

Electrostatics

Integral Representation

Point charge: $\psi(\vec{r}) = \frac{q}{\epsilon |\vec{r} - \vec{r}_c|}$

Volume charge: $\psi(\vec{r}) = \int_V \frac{\rho(\vec{r}')}{\epsilon |\vec{r} - \vec{r}'|} dV$

evaluation point \vec{r}

point charge \vec{r}_c

point charge \vec{r}'

Superposition integral

$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{\rho(\vec{r}')}{|\vec{r} - \vec{r}'|} d\vec{r}' = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{1}{|\vec{r} - \vec{r}'|} d\vec{r}' = 1$

point charge is a "delta" function in charge density
 $\rho(\vec{r}) = q \delta(\vec{r} - \vec{r}_c)$

KEY δ PROPERTY:

$\int G(\vec{r}) \delta(\vec{r} - \vec{r}_c) dV = G(\vec{r}_c)$

$\int \frac{q \delta(\vec{r} - \vec{r}_c)}{\epsilon |\vec{r} - \vec{r}_c|} dV = \frac{q}{\epsilon |\vec{r} - \vec{r}_c|}$

Differential Representation:

Volume Charge Potential: $\nabla^2 \psi(\vec{r}) = \frac{\partial^2 \psi(\vec{r})}{\partial x^2} + \frac{\partial^2 \psi(\vec{r})}{\partial y^2} + \frac{\partial^2 \psi(\vec{r})}{\partial z^2} = -\frac{\rho(\vec{r})}{\epsilon}$ (Poisson Equation)

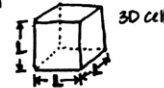
Laplacian: $\nabla^2 \psi(\vec{r}) = -\Delta \psi(\vec{r}) = 0$

Point Charge Potential: $\nabla^2 \psi(\vec{r}) = -\frac{q \delta(\vec{r} - \vec{r}_c)}{\epsilon} = 0$ except at $\vec{r} = \vec{r}_c$

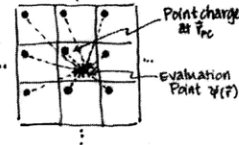
- 1) Ewald Sums Problems
- 2) Algorithms for the periodic case - mixture of the differential & integral forms

Periodic Case:

Cubic unit cell



2D picture:



consider 1 charge in a periodic case -

$\psi(\vec{r}) = \sum_{n_x, n_y, n_z} \frac{q}{\epsilon |\vec{r} - (\vec{r}_c + n_x \hat{x} + n_y \hat{y} + n_z \hat{z})|}$

at \vec{r}_c

same cell

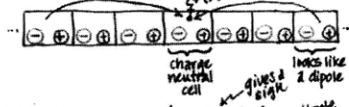
consider 1D (along position x) -

$\frac{1}{\epsilon |r - r_0|} + \frac{1}{\epsilon |r - r_0 + L|} + \frac{1}{\epsilon |r - r_0 + 2L|} + \dots$

terms $\rightarrow \frac{1}{\epsilon L} \frac{1}{n}$ for distant cells

$\sum_{n=k}^{\infty} \frac{1}{n} \rightarrow \infty$ divergent sum

Charge Neutrality Saves You... "sort of"



DIPOLE POTENTIAL:
$$\frac{(\vec{r} - \vec{r}_{dipole}) \cdot \vec{p}}{E |\vec{r} - \vec{r}_{dipole}|^3}$$

decays like $\frac{1}{r^2} \Rightarrow$ not fast enough in 3D

Symmetry ensures cancellation
left & right dipoles nearly cancel
difference dies like $\frac{1}{r^3} \Rightarrow$ fast enough

Computed sum depends on ORDER
Symmetric

3	2	3	3	3
3	2	2	2	3
3	2	1	2	3
3	2	2	2	3
3	3	3	3	3

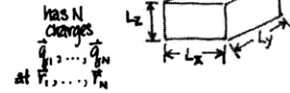
"shell ordering"

vs. asymmetric (diverges)

3	3	3	3	3
3	2	2	2	3
3	2	1	2	3
3	2	2	2	3
3	3	3	3	3

EWALD SUM ALGORITHMS

consider a 3-D cell:



only finite potential if $\sum q_i = 0$
charge neutral

finite potential \leftarrow must be periodic

Periodic Series are described by
Fourier Series:

potential
$$U(\vec{r}) = \sum_{m_x, m_y, m_z} \bar{U}[m_x, m_y, m_z] e^{j2\pi m_x \frac{x}{L_x}} e^{j2\pi m_y \frac{y}{L_y}} e^{j2\pi m_z \frac{z}{L_z}}$$

charge density is also periodic:
determine

charge density
$$\rho(\vec{r}) = \sum \bar{\rho}[m_x, m_y, m_z] e^{j2\pi (\frac{m_x}{L_x} x + \frac{m_y}{L_y} y + \frac{m_z}{L_z} z)}$$

charge density $\vec{m} \cdot \vec{r} = \left[\frac{m_x}{L_x} \quad \frac{m_y}{L_y} \quad \frac{m_z}{L_z} \right] \cdot \begin{bmatrix} x \\ y \\ z \end{bmatrix}$

Poisson Eqn: $\nabla^2 u = -4\pi\rho$

plug in the Fourier Representation on both sides
match Fourier terms...

$$\bar{U}[m_x, m_y, m_z] = \frac{-4\pi}{(2\pi)^3 \bar{m}^2} \bar{\rho}[m_x, m_y, m_z]$$

(scalar)
(same Fourier components)

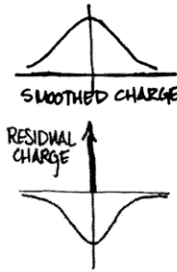
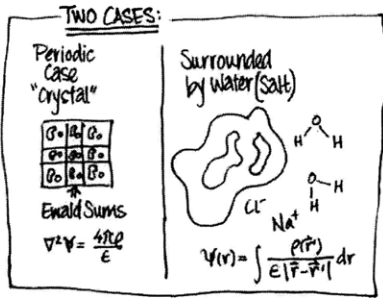
to solve periodic problem:

- take charge density
 - compute Fourier Series
 - easily compute Fourier Series for periodic potential
- hard due to point charges

2.8 Continuum Electrostatic Modeling II

U2B11/DE 482
LECTURE 8: CONTINUUM ELECTROSTATICS MODELING II

THU-FRI
7 MARCH 2006

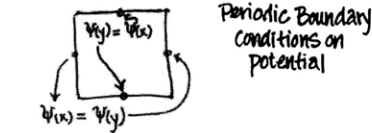


Fourier series of potential due to smoothed σ decays exponentially fast in Fourier space
few Fourier Coeffs can be truncated

Potential dies exponentially fast in space can be truncated

NOTE: Not charge neutral

PERIODIC CASE:



Fourier Series Coeffs for potential indices \rightarrow periodic function

$$\sum_{m_x, m_y, m_z} \bar{U}(m_x, m_y, m_z) e^{j2\pi(\frac{m_x}{L_x}x + \frac{m_y}{L_y}y + \frac{m_z}{L_z}z)}$$

Fourier Series Coefficients for Potential

$$\bar{U}(m_x, m_y, m_z) = \frac{1}{\Omega} \int \rho(\mathbf{r}) e^{-j2\pi(\frac{m_x}{L_x}x + \frac{m_y}{L_y}y + \frac{m_z}{L_z}z)} d\mathbf{r}$$

(σ^2 doesn't mix frequency) FS. for charge

$$\sigma(\mathbf{r}) = \sum_{\mathbf{m}} \bar{\sigma}(\mathbf{m}) e^{j2\pi \mathbf{m} \cdot \mathbf{r}}$$

\rightarrow multiply by $e^{-j2\pi \mathbf{m} \cdot \mathbf{r}}$ & integrate

$$\bar{\sigma}(\mathbf{m}) = \frac{1}{L_x L_y L_z} \int \sigma(\mathbf{r}) e^{-j2\pi(\frac{m_x}{L_x}x + \frac{m_y}{L_y}y + \frac{m_z}{L_z}z)} d\mathbf{r}$$

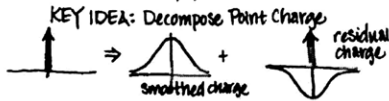
$$\bar{\sigma}(\mathbf{r}) = \sum q_i \delta(\mathbf{r} - \mathbf{r}_i)$$

$$\bar{\sigma}(\mathbf{m}) = \frac{1}{V} \sum_{i=1}^N q_i e^{-j2\pi \mathbf{m} \cdot \mathbf{r}_i}$$

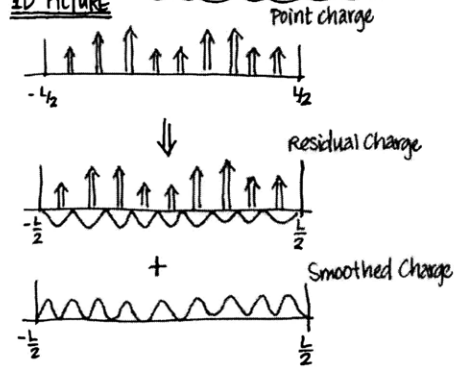
FROM MATCHING:

$$\bar{U}(\mathbf{m}) \approx \frac{1}{\epsilon V} \sum q_i \frac{e^{-j2\pi \mathbf{m} \cdot \mathbf{r}_i}}{|\mathbf{m}|^2}$$

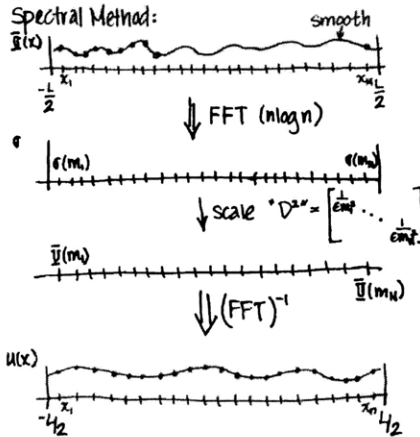
$\bar{U}(\mathbf{m}) \propto \frac{1}{|\mathbf{m}|^2} \Rightarrow$ SLOW DECAY



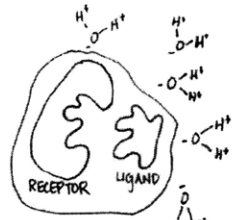
4D PICTURE



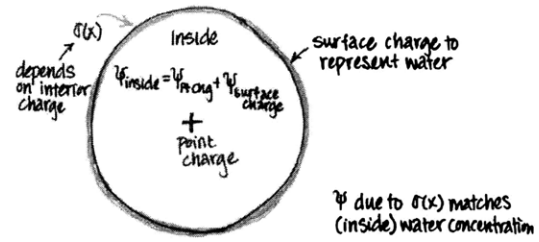
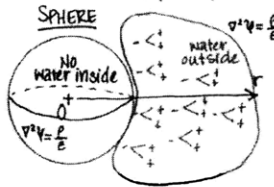
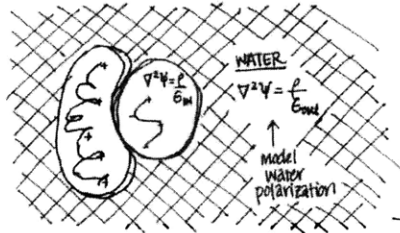
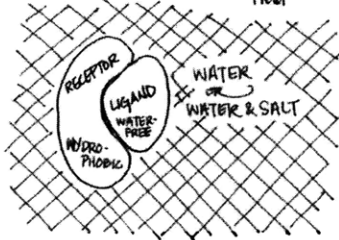
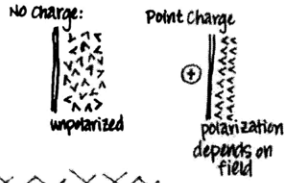
- 1) Potential due to Residual affects nearby points only ($O(N)$ ops constant # pts per residual charge)
- 2) Use Spectral Method to compute: $\nabla^2 U_{smooth} = \frac{\rho_{smooth}}{\epsilon}$
- 3) Add results



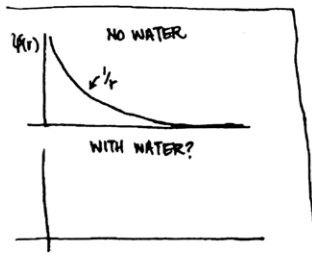
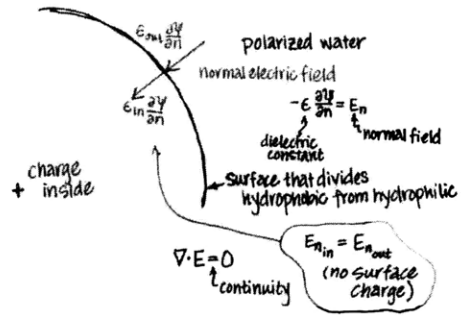
SURROUNDED BY WATER (SALT)



Polar molecules will react to the fields



Original Problem:



2.9 Electrostatic Contributions to Binding and Design

THURSDAY
9 MARCH 2006

6.581/BE 422
 LECTURE 9: ELECTROSTATIC CONTRIBUTIONS
 TO BINDING AND DESIGN

NOTE: Prof. White will finish his previous lecture in the next lecture.

REVIEW:

Molecular Models: $U(\vec{r}^{3N}) = U_{\text{COVALENT}} + U_{\text{NON-COVALENT}}$

Protein Folding:

Protein Binding:

Simulation & Energy Evaluation:

- Integrating Eqs of Motion
 - movie
 - response to force
 - statistical mechanical ensemble
- Computing long-range electrostatic energies
 - explicit treatment - numerical approximation to arbitrary accuracy
 - implicit methods - continuum ideas

Linearized PB Eqn:

$$\vec{\nabla} \cdot [\epsilon(\vec{r}) \vec{\nabla} \psi(\vec{r})] - \kappa^2 \psi(\vec{r}) = -4\pi \rho(\vec{r})$$

Linear \Leftrightarrow Superposition

$$f(x+y) = f(x) + f(y)$$

$$f(ax) = a \cdot f(x)$$

Free Energy:

$$G = \frac{1}{2} \int \rho(\vec{r}) \psi(\vec{r}) d\vec{r} = \frac{1}{2} \sum_{i=1}^N q_i \psi_i$$

$$= \frac{1}{2} [q_1, q_2, \dots, q_N] \cdot \begin{bmatrix} \psi_1 \\ \psi_2 \\ \vdots \\ \psi_N \end{bmatrix}$$

$\psi =$ fictitious surface charge

$$G = \frac{1}{2} \vec{q}^T M \vec{q} = \vec{q}^T M \vec{q}$$

where $M = \frac{1}{2} M'$

$$G = \vec{q}^T M \vec{q}$$

symmetric due to reciprocity

Assume Rigid binding

$$\Delta G_{\text{binding}}^{\text{elec}} = G_{\text{bound}} - G_{\text{unbound}} = \vec{q}^T M_{\text{bound}} \vec{q} - \vec{q}^T M_{\text{unbound}} \vec{q} = \vec{q}^T (\Delta M) \vec{q}$$

SOLVED BY:

- Finite Difference
- Finite Element
- Boundary Element*

Explicit Solvent

- ⊕ physically more accurate
- ⊕ cheaper to compute a single energy
- ⊖ thermodynamic properties are not expressed in a single frame

Continuum Solvent

- ⊕ mean field properties represented well in single frame
- ⊕ cheaper to compute free energies
- ⊖ less accurate physically

Poisson-Boltzmann Equation

$$\vec{\nabla} \cdot [\epsilon(\vec{r}) \vec{\nabla} \psi(\vec{r})] - \kappa^2 \sinh(\psi(\vec{r})) = -4\pi \rho(\vec{r})$$

non-linear

Linearized Poisson-Boltzmann Eqn:

$$\vec{\nabla} \cdot [\epsilon(\vec{r}) \vec{\nabla} \psi(\vec{r})] - \kappa^2 \psi(\vec{r}) = -4\pi \rho(\vec{r})$$

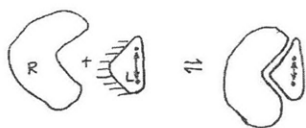
ionic strength $\kappa=0$ for no salt

$\kappa^2 = \kappa_D^{-2}$ (linearized)

$\Rightarrow \Delta G_{\text{binding}}^{\text{elec}} = \begin{bmatrix} q_L \\ q_R \end{bmatrix}^T \begin{bmatrix} L_{LL} & C_{LR} \\ C_{RL} & R_{RR} \end{bmatrix} \begin{bmatrix} q_L \\ q_R \end{bmatrix}$

$= q_L^T L_{LL} q_L + 2q_L^T C_{LR} q_R + q_R^T R_{RR} q_R$

ligand-ligand interactions, ligand-receptor interactions, receptor-receptor interactions



$$\Delta G_{\text{binding}}^{\text{elec}} = \underbrace{\vec{q}_L^T \mathbf{L} \vec{q}_L}_{\text{ligand desolvation}} + \underbrace{2\vec{q}_L^T \mathbf{C} \vec{q}_R}_{\text{interactions recovered}} + \underbrace{\vec{q}_R^T \mathbf{R} \vec{q}_R}_{\text{receptor desolvation}}$$

$$\begin{aligned} \vec{\nabla}_{\vec{q}_L} (\Delta G_{\text{binding}}^{\text{elec}}) &= 2\mathbf{L} \vec{q}_L + 2\mathbf{C} \vec{q}_R = \vec{0} \\ &\Rightarrow \mathbf{L} \vec{q}_L = -\mathbf{C} \vec{q}_R \\ &\Rightarrow \boxed{\vec{q}_L^{\text{opt}} = -\mathbf{L}^{-1} \mathbf{C} \vec{q}_R} \end{aligned}$$

JPhys B 105: 380-388 (2001).
chorismate mutase



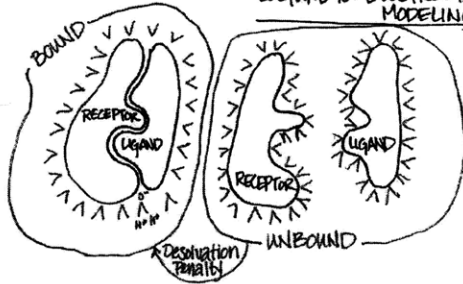
JACS 125: 5598-9 (2003).
10-fold improvement in binding

2.10 Electrostatics Modeling

6.501/BE 482

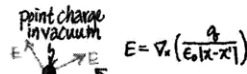
LECTURE 10: ELECTROSTATICS MODELING

TUESDAY
14 MARCH 2006

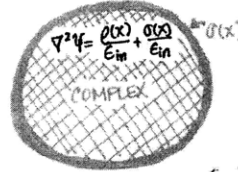
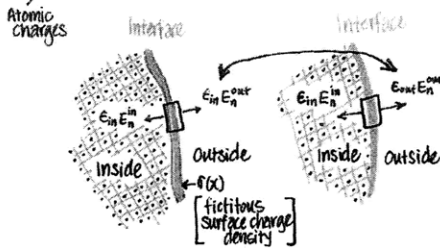
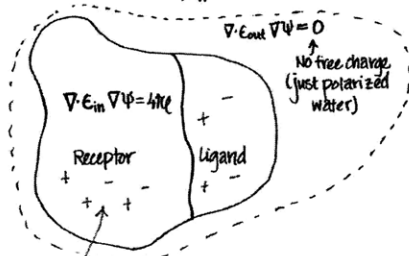
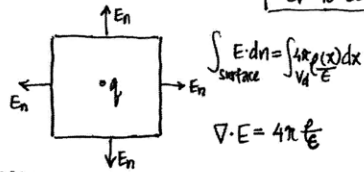
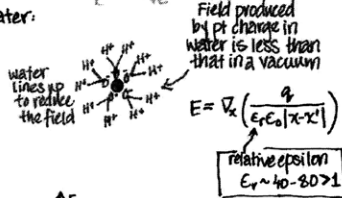


CONTINUUM MODEL

outside water:



in water:



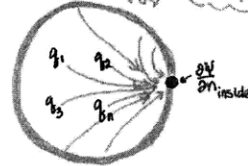
$\nabla^2 \psi = 0$

Choose $f(x)$ so that:

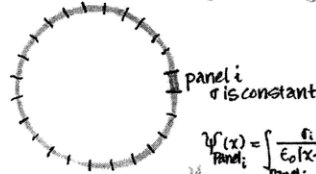
$\frac{\partial \psi}{\partial n} = \epsilon_r \frac{\partial \psi}{\partial n}$ inside	original
$\frac{\partial \psi}{\partial n} - \frac{\partial \psi}{\partial n} = -4\pi \frac{\sigma(x)}{\epsilon_{in}}$ outside	fictitious problem

After: $\epsilon_{in} (1 - \epsilon_r) \frac{\partial \psi(\mathbf{x}_i)}{\partial n} = \sigma(\mathbf{x}_i)$

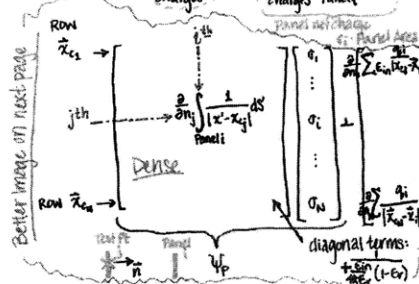
center of the i th panel



Atomic charges interaction between atomic charges is changed by surrounding water



$\psi(x) = \sum_{atomic\ charges} \frac{q_i}{\epsilon_{in} |\mathbf{x} - \mathbf{x}_i|} + \sum_{panel\ charges} \int \frac{\sigma_i}{\epsilon_{in} |\mathbf{x} - \mathbf{x}'|} dS'$



$$\underbrace{\frac{\partial \psi_p(\vec{x}_i)}{\partial n}}_{\text{Panel Charge Potential}} - \frac{q_i}{4\pi\epsilon_0(1-\epsilon_r)} = \underbrace{-\frac{\partial \psi_A(\vec{x}_i)}{\partial n}}_{\text{Atomic Charge Potential}} \quad \text{known}$$

Dense)
 $P\vec{q}_p = \vec{\psi}_A$
 Iterate: Guess at \vec{q}_p^0
 Calculate $\vec{\psi}_A - P\vec{q}_p^0 = \text{Residual}$
 $\vec{q}_p^1 = f(\text{Residual})$
 Matrix Vector Product
 N potential derivatives from N charges using Multipole $\mathcal{O}(N)$

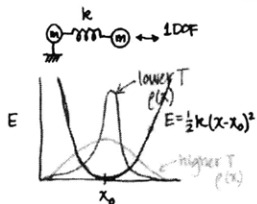
$$\begin{array}{c} \text{Row } \vec{x}_1 \rightarrow \\ \vdots \\ \text{jth} \rightarrow \\ \vdots \\ \text{Row } \vec{x}_N \rightarrow \end{array}
 \left[\begin{array}{c} \text{DENSE MATRIX} \\ \text{(Every element is nonzero)} \end{array} \right]
 \begin{array}{c} \vdots \\ q_1 \\ \vdots \\ q_i \\ \vdots \\ q_N \end{array}
 = -
 \begin{array}{c} \psi_A \\ \vdots \\ \frac{\partial}{\partial n_j} \sum \frac{q_i}{\epsilon_{in} |\vec{x}_{ej} - \vec{x}_i|} \\ \vdots \\ \frac{\partial}{\partial n_N} \sum \frac{q_i}{\epsilon_{in} |\vec{x}_{en} - \vec{x}_i|} \end{array}$$

Diagonal terms:
 $+$ $\frac{1}{4\pi\epsilon_0(1-\epsilon_r)}$

2.11 Statistical Mechanics

6.581 / BE.482
LECTURE 11: STATISTICAL MECHANICS

MIT COURSE 16 MARCH 2006



Imagine a dense system of these packed together → see probability distributions above
 $p(x) = \text{probability of finding a copy of the system with configuration } x = \frac{e^{-\frac{E(x)}{k_B T}}}{\int e^{-\frac{E(x')}{k_B T}} dx'}$

$k_B = \text{Boltzmann Constant} = 1.38 \cdot 10^{-23} \frac{\text{kcal}}{\text{mol} \cdot \text{K}}$
 $T = \text{Absolute Temperature} = 273.15 + t \text{ (}^\circ\text{C)}$

Note: $\int_{-\infty}^{\infty} p(x) dx = \int_{-\infty}^{\infty} \frac{e^{-\frac{E(x)}{k_B T}}}{Q} dx = \frac{1}{Q} Q = 1$

What is average position?

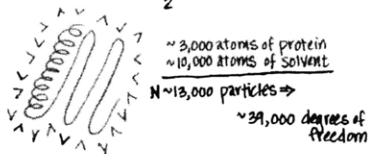
$\bar{x} = \sum_{\text{All possible configurational states}} (\text{probability of } x) \cdot (\text{value at that position, } x)$
 $= \int_{-\infty}^{\infty} p(x) x dx = \int_{-\infty}^{\infty} \frac{x e^{-\frac{E(x)}{k_B T}}}{Q} dx = \int_{-\infty}^{\infty} \frac{x e^{-\frac{1}{2} k(x-x_0)^2}}{Q} dx$
 change of variables: $y = x - x_0$, $dy = dx$
 $= \frac{1}{Q} \int_{-\infty}^{\infty} y e^{-\frac{1}{2} k y^2} dy + \frac{x_0}{Q} \int_{-\infty}^{\infty} e^{-\frac{1}{2} k y^2} dy = x_0$

Computing Average Fluctuation:

$\sqrt{\langle (x-\bar{x})^2 \rangle} = \sqrt{\frac{1}{Q} \int_{-\infty}^{\infty} (x-x_0)^2 e^{-\frac{k(x-x_0)^2}{2k_B T}} dx}$
 $= \sqrt{\frac{k_B T}{k}}$
 → higher T ⇒ greater fluctuation
 → higher k ⇒ lower fluctuation

Computing Average Potential Energy

$\bar{E} = \langle E \rangle = \frac{1}{Q} \int_{-\infty}^{\infty} E(x) e^{-\frac{E(x)}{k_B T}} dx = \frac{1}{Q} \int_{-\infty}^{\infty} \frac{1}{2} k(x-x_0)^2 e^{-\frac{k(x-x_0)^2}{2k_B T}} dx = \frac{k_B T}{2}$

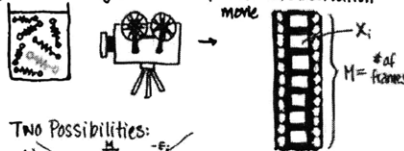


Compute Average Str. of Protein

$\langle \bar{X}^{3N} \rangle = \frac{\int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \bar{X}^{3N} e^{-\frac{E(\bar{X})}{k_B T}} d\bar{X}^{3N}}{\int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} e^{-\frac{E(\bar{X})}{k_B T}} d\bar{X}^{3N}}$

first problem: $E(\bar{X}^{3N})$ is not analytically integrable → must do it numerically
 second problem: 0.1 Å on a 100 Å grid
 1000 points per dimension
 (1000)^{3,000} points total

Imagine Constructing $e(x)$ from Experimental Observation



Two Possibilities:

1) $\langle E \rangle = \frac{1}{M} \sum_{i=1}^M E_i e^{-\frac{E_i}{k_B T}}$
 ~~$\langle E \rangle = \frac{1}{M} \sum_{i=1}^M E_i$~~

If we use a molecular dynamics simulation as the movie:

Geometric Quantity: $\langle |r_i - r_j| \rangle = \frac{1}{M} \sum_{i=1}^M |r_i - r_j(\bar{r}_k)|$
 Interaction Energy: $\langle U_{i-j} \rangle = \frac{1}{M} \sum_{i=1}^M U_{i-j}(\bar{r}_k)$
 Total Potential Energy: $\langle U \rangle = \frac{1}{M} \sum_{i=1}^M U(\bar{r}_k)$

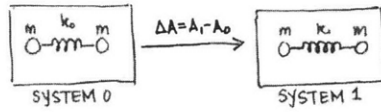
Metropolis Monte Carlo, like MD, also produces a trajectory that converges to a statistical mechanical ensemble
 Metropolis et al, J Chem Phys 21: 1087-1092 (1953).

Problem trying to get the free energy -

I attempted to write: $\langle A \rangle = \frac{1}{M} \sum_{i=1}^M A_i$ (this doesn't exist)

Statistical Mechanics gives us a definition:
 $A = -k_B T \ln Q = -k_B T \ln \left[\int_{-\infty}^{\infty} e^{-\frac{E(x)}{k_B T}} dx \right]$

Doesn't help, because Q is 1000^{3,000} dimensional integral



Thermodynamic Integration (Kirkwood, 1934)

1) Construct a hybrid potential that smoothly connects 0 \rightarrow 1.

$$U(\lambda) = (1-\lambda)U_0 + \lambda U_1 \quad \lambda: (0 \rightarrow 1)$$

$U: (U_0 \rightarrow U_1)$

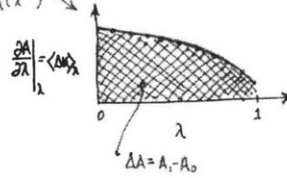
2) Fundamental Theorem of Integral Calculus

$$\Delta A = A_1 - A_0 = \int_0^1 \frac{\partial A}{\partial \lambda} d\lambda$$

$$A(\lambda) = -k_B T \ln \left[\int e^{-\frac{U(\vec{x}^{3N}, \lambda)}{k_B T}} d\vec{x}^{3N} \right]$$

$$\frac{\partial A(\lambda)}{\partial \lambda} = -k_B T \frac{\int e^{-\frac{U(\vec{x}^{3N}, \lambda)}{k_B T}} \left[-\frac{1}{k_B T} \frac{\partial U(\vec{x}^{3N}, \lambda)}{\partial \lambda} \right] d\vec{x}^{3N}}{\int e^{-\frac{U(\vec{x}^{3N}, \lambda)}{k_B T}} d\vec{x}^{3N}} = \left\langle \frac{\partial U}{\partial \lambda} \right\rangle_\lambda = \langle U_1 - U_0 \rangle_\lambda = \langle \Delta U \rangle_\lambda$$

$$\Delta A = A_1 - A_0 = \int_0^1 \langle \Delta U \rangle_\lambda d\lambda$$



2.12 Statistical Mechanics

PSET 3
DUE ON TUES
4 APRIL

6.581/EE.482
LECTURE 12: STATISTICAL MECHANICS

TUESDAY
21 MARCH 2002

SUMMARY OF LAST TIME

- Uniform Sampling: $f(\vec{x}) = \langle f(\vec{x}) \rangle = \int d\vec{x} p(\vec{x}) f(\vec{x})$
- Boltzmann Sampling: $\langle f(\vec{x}) \rangle = \frac{1}{N} \sum_{i=1}^N f(\vec{x}_i)$

Probability density function: $p(\vec{x}) = \frac{e^{-\beta E(\vec{x})}}{\int d\vec{x} e^{-\beta E(\vec{x})}}$

Configuration partition function: $Q = \int d\vec{x} e^{-\beta E(\vec{x})}$

Free Energy: $A = -k_B T \ln Q$

Problem: Free Energy doesn't exist in single frame $f(\vec{x})$ so can't compute this way

Solution: Can compute free energy changes from simulation (ensemble) averages

Then $\Delta A = A_1 - A_0 = \int d\lambda \langle \Delta U \rangle_\lambda$

Science, 244:1069-1072 (1989)

Hemoglobin

ALLOSTERIC PROTEIN WHICH BINDS OXYGEN

Independent binding events

cooperative binding events

Oxygen binding graph: 100% vs 0% O₂. Shows sigmoidal curves for 'muscle' and 'lungs'.

wild type: $DEOXY \rightleftharpoons OXY$ (low O₂ affinity) \rightleftharpoons (high O₂ affinity)

mutant: $DEOXY \rightleftharpoons OXY$ (Hb Radcliffe mutation) $Asp\ G1(99)\beta \rightarrow Ala$

Al for "allosteric"

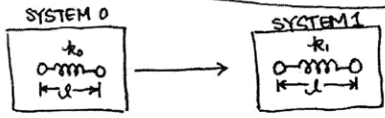
Experimentally:

$\Delta \Delta A = \Delta A' - \Delta A = -3.4$ kcal/mol/interface

- Reduced cooperativity & increased oxygen affinity
- Shifts equilibrium toward Oxy'
- Less ability to release O₂ at muscles

DEOXY (wild type) $Asp\ 99\beta$ $Asn\ G4(97)\alpha$ $Tyr\ C7(42)\alpha$

OXY $Asp\ 99\beta$ $Asp\ 94\alpha$ $Tyr\ 42\alpha$ $Asn\ 97\alpha$



Run a series of simulations with

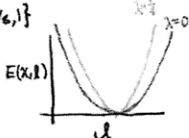
$$k(\lambda) = (1-\lambda)k_0 + \lambda k_1$$

$$= \frac{1}{2}(1-\lambda)k_0(x-1)^2 + \frac{1}{2}\lambda k_1(x-1)^2$$

$$= \frac{1}{2}[k_0 + \lambda(k_1 - k_0)](x-1)^2$$

$\lambda = \{0, 1/8, 1/4, 3/8, 1/2, 5/8, 3/4, 7/8, 1\}$

$$\Delta U = \frac{1}{2}(k_1 - k_0)(x-1)^2$$



$$\Delta U = U_1 - U_0 = \sum_i \Delta U_i$$

$$\Delta A = \int_0^1 d\lambda \langle \Delta U \rangle = \int_0^1 d\lambda \langle \sum_i \Delta U_i \rangle = \sum_i \int_0^1 d\lambda \langle \Delta U_i \rangle = \sum_i \Delta A_i$$

thermodynamic cycle $\rightarrow 0 = \Delta A_{deoxy} + \Delta A'_{ox} - \Delta A_{oxy} - \Delta A_{al}$

$$\Delta \Delta A = \Delta A'_{al} - \Delta A_{al} = \Delta A_{oxy} - \Delta A_{deoxy}$$

$$\Delta \Delta A = -5.5$$
 kcal/mol/interface

$\Delta \Delta A$ [kcal/mol/interface]

Total -5.5

Solvent +22.5 \leftarrow strong solvation in oxy state of Asp

Tyr42a -12.7 \leftarrow consistent w/ loss of h-bond

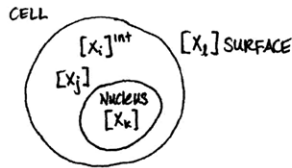
Asp94a -22.4 \leftarrow strong repulsion in oxy form

Yal96a 5.5

Asn97a 3.3 \leftarrow inconsistent w/ std. explanation

2.13 Formulating Models

DIFFERENTIAL EQUATION MODEL



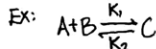
VERY GENERAL SETTING

$$\begin{bmatrix} [X_1] \\ \vdots \\ [X_n] \end{bmatrix} = \vec{X} \quad \vec{U} = \begin{bmatrix} [U_1] \\ \vdots \\ [U_m] \end{bmatrix}$$

time-variation of concentrations

$$\frac{d}{dt} [X_i] = V^*(\vec{X}, \vec{U}) - V^*(\vec{X}, \vec{U})$$

rate of production rate of consumption



$$\frac{d}{dt} [C] = k_1 [A][B] - k_2 [C]$$

$$\frac{d}{dt} [A] = k_2 [C] - k_1 [A][B]$$

$$\frac{d}{dt} [B] = k_2 [C] - k_1 [A][B]$$

+ Initial conditions
 $A(0), B(0), C(0)$

EX: Dimerization



$$\frac{d}{dt} [A_2] = k_1 [A][A] - k_2 [A_2]$$

$$\frac{d}{dt} [A] = 2k_2 [A_2] - 2k_1 [A][A]$$

STEADY-STATE (Equilibrium)

Reactions are in balance

$$\frac{d}{dt} [X_i] = 0$$

continuum approximation
bc molecules are reacting, but overall the net rate is balanced
↓
Lots of copies of X

Dimerization case:
 $k_1 [A][A] = k_2 [A_2]$
 $[A]^2 = \frac{k_2}{k_1} [A_2]$
 $[A] = \sqrt{\frac{k_2}{k_1}} [A_2]^{1/2}$
Also given:
 $A(0) = A_0, A_2(0) = 0$
→ initial conditions
Steady-State:
 $[A] + 2[A_2] = A_0$

CHEMICAL KINETICS FORM

$$\frac{d}{dt} [X_i] = \sum_{j=1}^N \gamma_j^i [X_j] \quad \text{destination}$$

(Typical Case $\gamma_i^i < 0$)
degradation process

$$\frac{d}{dt} [X_i] = \sum_{j=1}^N \gamma_j^i [X_j] + \sum_{j=1}^N \sum_{k=1}^N \gamma_{jk}^i [X_j][X_k]$$

(Typical Consumption Process $\gamma_{ik}^i < 0$)
(Typical Production Process $\gamma_{jk}^i > 0$)

$$\frac{d}{dt} [X_i] = \sum_{j=1}^N \gamma_j^i [X_j] + \sum_{j=1}^N \sum_{k=1}^N \gamma_{jk}^i [X_j][X_k] + \sum_{j=1}^M \alpha_j^i [U_j] + \sum_{j=1}^M \sum_{k=1}^N \alpha_{jk}^i [U_j][X_k] + \sum_{j=1}^M \sum_{k=1}^M \beta_{jk}^i [U_j][U_k]$$

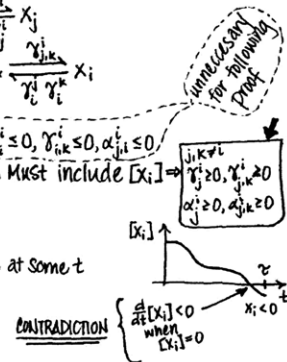
MULTIPLE INPUTS

$\gamma_i^i < 0 \Rightarrow$ models degradation
 $\gamma_j^i \neq 0 \Rightarrow$ models conversion $X_j \xrightarrow{\gamma_j^i} X_i$
 γ_{jk}^i
 $X_j + X_k \xrightarrow{\gamma_{jk}^i} X_i$

If $[X_i]$ is involved only in reactions that consume $[X_i] \Rightarrow \gamma_i^i \leq 0, \gamma_{ik}^i \leq 0, \alpha_{ji}^i \leq 0$
and reactions that produce $[X_i]$ must include $[X_i] \Rightarrow \gamma_j^i > 0, \gamma_{jk}^i > 0, \alpha_j^i > 0, \alpha_{jk}^i > 0$

Then All $[X_i] \geq 0$ for all t

Proof: Suppose $[X_i] < 0$ for some i at some t



$$\frac{d}{dt} \vec{X} = A^{(1)} \vec{X} + A^{(2)} \begin{bmatrix} [x_1][x_1] \\ [x_1][x_2] \\ \vdots \\ [x_1][x_n] \\ [x_2][x_1] \\ [x_2][x_2] \\ \vdots \\ [x_2][x_n] \\ [x_3][x_1] \\ \vdots \\ [x_n][x_n] \end{bmatrix}$$

terms $\sum \gamma_j^i [x_j]$
 $\Rightarrow A_{ij}^{(1)} = \gamma_j^i$

terms appear twice... not a big deal
 N^2 terms

$$\begin{matrix} \uparrow N \\ A^{(2)} \\ \downarrow N^2 \end{matrix}$$

$$\Rightarrow A_{ij}^{(2)} = \gamma_{j,k}^i$$

$[(j-1)n+k+1]$

KROENEKER PRODUCT:

$$\vec{X} \otimes \vec{X} = \begin{bmatrix} x_1 \vec{x} \\ x_2 \vec{x} \\ x_3 \vec{x} \\ \vdots \\ x_n \vec{x} \end{bmatrix}$$

$$\frac{d}{dt} \vec{X} = A^{(1)} \vec{X} + A^{(2)} \vec{X} \otimes \vec{X} + B \vec{u} + C \vec{u} \otimes \vec{X} + D \vec{u} \otimes \vec{u}$$

$$\begin{bmatrix} u_1 \vec{x} \\ u_2 \vec{x} \\ \vdots \\ u_m \vec{x} \end{bmatrix} = \begin{bmatrix} [u_1][x_1] \\ [u_1][x_2] \\ \vdots \\ [u_1][x_n] \\ [u_2][x_1] \\ \vdots \\ [u_m][x_n] \end{bmatrix}$$

2.14 Nonlinear Dynamics and Stability

6.581/BE.482

LECTURE 15: NONLINEAR DYNAMICS & STABILITY

THURSDAY
6 APRIL 2006

COMPUTING STEADY-STATES

$$\frac{d\vec{x}}{dt} = A^{(0)}\vec{x} + A^{(1)}\vec{x} \otimes \vec{x} + B^{(0)}\vec{u} + B^{(1)}\vec{u} \otimes \vec{x} + B^{(2)}\vec{u} \otimes \vec{x} \otimes \vec{x}$$

known ignore for simplicity



KROENECKER PRODUCT:

$$A \otimes B = \begin{bmatrix} A_{11}B & A_{12}B & \dots & A_{1n}B \\ A_{21}B & A_{22}B & \dots & A_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ A_{m1}B & A_{m2}B & \dots & A_{mn}B \end{bmatrix}$$

$$\vec{x} \otimes \vec{x} = \begin{bmatrix} x_1x_1 \\ x_1x_2 \\ \vdots \\ x_1x_n \\ x_2x_1 \\ \vdots \\ x_nx_n \end{bmatrix}$$



In steady-state: $\frac{d}{dt}[X_i] = 0$

$\frac{d}{dt}\vec{x} = 0$ (vector case)

FIND \vec{x}^* 's $\rightarrow A^{(0)}\vec{x}^* + A^{(1)}\vec{x}^* \otimes \vec{x}^* + B^{(0)}\vec{u} = 0$

IS THE STEADY-STATE STABLE?

Suppose $\vec{x}^* : A^{(0)}\vec{x}^* + A^{(1)}\vec{x}^* \otimes \vec{x}^* + B^{(0)}\vec{u} = 0$

If $\vec{x}(t) = \vec{x}^* + \vec{\epsilon}$ for small $\vec{\epsilon}$, does soln $\rightarrow \vec{x}^*$?

$\frac{d}{dt}\vec{x} = A^{(0)}\vec{x} + A^{(1)}\vec{x} \otimes \vec{x} + B^{(0)}\vec{u}$

only STABLE steady-states are observable

stable means \vec{x}^* is approached when starting from $\vec{x}^* + \vec{\epsilon}$

CASE: $A^{(0)} = 0 \leftarrow$ NOT Biological

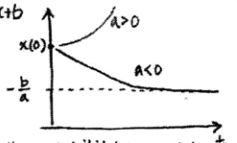
$\frac{d}{dt}\vec{x}(t) = A\vec{x}(t) + B\vec{u} \leftarrow$ constant in time

$A\vec{x}^* + B\vec{u} = 0 \Rightarrow \vec{x}^* = -A^{-1}B\vec{u}$

FOR SCALAR CASE:

$\frac{d}{dt}x = ax + b$

$x^* = -\frac{b}{a}$



observe that stability can only occur if $a < 0$

FOR VECTOR CASE:

$\frac{d}{dt}\vec{x} = A\vec{x} + B\vec{u} \Rightarrow \vec{x} = -A^{-1}B\vec{u}$

In order for this to be stable, what must be said about A?

Re {eigenvalues of A} < 0

BUT, WHY EIGENVALUES?

$A\vec{s}_i = \lambda_i \vec{s}_i$ eigenvector

$A[\vec{s}_1, \vec{s}_2, \dots, \vec{s}_n] = \begin{bmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \lambda_n \end{bmatrix} [S]$

EIGEN (SPECTRAL) DECOMPOSITION:

$\vec{x}(t) = d_1(t)\vec{s}_1 + d_2(t)\vec{s}_2 + \dots + d_n(t)\vec{s}_n \Rightarrow \vec{x}(t) = [S]\vec{\alpha}(t)$

$\frac{d}{dt}\vec{x} = A\vec{x} + B\vec{u}$

$\Rightarrow \frac{d}{dt}\vec{\alpha}(t) = \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_n \end{bmatrix} \vec{\alpha}(t) + [S]^{-1}B\vec{u}$

$\Rightarrow \begin{cases} \frac{d}{dt}d_1 = \lambda_1 d_1 + (S^{-1}B\vec{u})_1 \\ \frac{d}{dt}d_2 = \lambda_2 d_2 + (S^{-1}B\vec{u})_2 \\ \vdots \\ \frac{d}{dt}d_n = \lambda_n d_n + (S^{-1}B\vec{u})_n \end{cases}$

Re $(\lambda_i) < 0 \forall i \Rightarrow d_i$'s don't blow up

CONSIDER A GENERAL DYNAMIC, NONLINEAR SYSTEM OF EQUATIONS

$\frac{d}{dt}\vec{x} = F(\vec{x}) = \begin{bmatrix} f_1(x_1, x_2, \dots, x_n) \\ f_2(x_1, x_2, \dots, x_n) \\ \vdots \\ f_n(x_1, x_2, \dots, x_n) \end{bmatrix}$

$F(\vec{x}) = A^{(0)}\vec{x} + A^{(1)}\vec{x} \otimes \vec{x} + B\vec{u}$

STEADY-STATE EQN: Find \vec{x}^* s.t. $F(\vec{x}^*) = 0$

$F(\vec{x}^* + \vec{\epsilon}) = F(\vec{x}^*) + J_F(\vec{x}^*)\vec{\epsilon}$

(for scalar case: $f(x+\epsilon) \approx f(x) + \frac{df(x)}{dx}\epsilon + H.O.T.$)

SYSTEM ABOUT $\vec{x}^* \Rightarrow$

$\frac{d}{dt}\vec{x} = F(\vec{x}) \Rightarrow \frac{d}{dt}(\vec{x}^* + \vec{\epsilon}) = F(\vec{x}^* + \vec{\epsilon})$

$\Rightarrow \frac{d}{dt}\vec{\epsilon} \approx J_F(\vec{x}^*)\vec{\epsilon}$

$J_F(\vec{x}^*)\vec{\epsilon} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \dots & \frac{\partial f_n}{\partial x_n} \end{bmatrix} \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$

\vec{x}^* is a stable, steady-state IF Re {eigenvalues of $J_F(\vec{x}^*)$ } < 0

CASE: $A^{(2)} \neq 0$ (Biological Case)

$$F(\vec{x}) = A^{(1)}\vec{x} + A^{(2)}\vec{x} \otimes \vec{x} + B\vec{u}$$

$$J_F(\vec{x}) = A^{(1)} + A^{(2)}(\mathbf{I} \otimes \vec{x} + \vec{x} \otimes \mathbf{I})$$

STEADY-STATE PROBLEM:

① FIND \vec{x}^* s.t. $A^{(1)}\vec{x}^* + A^{(2)}\vec{x}^* \otimes \vec{x}^* + B\vec{u} = 0$

② VERIFY $\lambda(A^{(1)} + A^{(2)}(\mathbf{I} \otimes \vec{x}^* + \vec{x}^* \otimes \mathbf{I}))$ have negative real parts

use NEWTON'S METHOD to solve

Problem: Find \vec{x}^* s.t. $F(\vec{x}^*) = 0$

FOR THE SCALAR CASE:

$$f(x^*) = 0$$

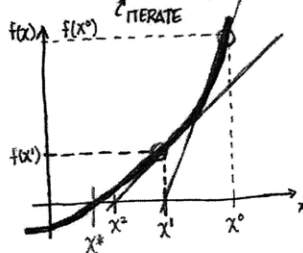
use Taylor's expansion for nonlinear f :

$$f(x^*) = f(x) + \frac{df(x)}{dx}(x - x^*) + \text{H.O.T.}$$

$$\Rightarrow \frac{df}{dx}(x - x^*) = -f(x)$$

GUESS x^0

$$\text{COMPUTE } x^1: x^1 - x^0 = -\left(\frac{df}{dx}\right)^{-1} f(x^0)$$



\vec{x}^* s.t. $F(\vec{x}^*) = 0$

guess at \vec{x}^0

$$\begin{cases} J_F(\vec{x}^0)(\vec{x}^1 - \vec{x}^0) = -F(\vec{x}^0) \Rightarrow \vec{x}^1 \\ J_F(\vec{x}^1)(\vec{x}^2 - \vec{x}^1) = -F(\vec{x}^1) \Rightarrow \vec{x}^2 \end{cases}$$

where

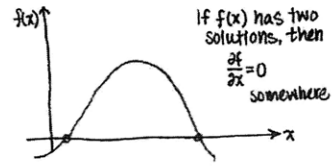
$$F(\vec{x}^0) = A^{(1)}\vec{x}^0 + A^{(2)}\vec{x}^0 \otimes \vec{x}^0 + B\vec{u}$$

$$J_F(\vec{x}^0) = A^{(1)} + A^{(2)}(\mathbf{I} \otimes \vec{x}^0 + \vec{x}^0 \otimes \mathbf{I})$$

thus

$$\begin{cases} [A^{(1)} + A^{(2)}(\mathbf{I} \otimes \vec{x}^0 + \vec{x}^0 \otimes \mathbf{I})](\vec{x}^1 - \vec{x}^0) = -F(\vec{x}^0) \\ [A^{(1)} + A^{(2)}(\mathbf{I} \otimes \vec{x}^1 + \vec{x}^1 \otimes \mathbf{I})](\vec{x}^2 - \vec{x}^1) = -F(\vec{x}^1) \end{cases}$$

ONCE AGAIN, FOR THE SCALAR CASE:



SO, FOR $F(\vec{x}) = 0$ TO HAVE MULTIPLE SOLNS

$$\Rightarrow J_F(\vec{x})\vec{v} = 0 \quad \text{For some } \vec{x}, \text{ there is a flat direction}$$

IF $J_F(\vec{x})$ IS NONSINGULAR FOR ALL \vec{x} THEN $F(\vec{x}) = 0$ HAS ONE SOLUTION

IF $A^{(1)} + A^{(2)}(\mathbf{I} \otimes \vec{x} + \vec{x} \otimes \mathbf{I})$ IS ALWAYS NONSINGULAR

Then STEADY STATE SOLUTION IS UNIQUE

2.15 Steady-State Problems

6.581/BE.482
LECTURE 16: STEADY-STATE PROBLEMS

TUESDAY
11 APRIL 2006

SIMPLIFIED SYSTEM

$$\frac{d\vec{x}}{dt} = A^{(1)}\vec{x} + A^{(2)}(\vec{x} \otimes \vec{x}) + B\vec{u}$$

STEADY STATE

Solve $A^{(1)}\vec{x} + A^{(2)}(\vec{x} \otimes \vec{x}) + B\vec{u} = 0$
 $F(\vec{x}) = 0$

NEWTON'S METHOD

Guess at \vec{x}

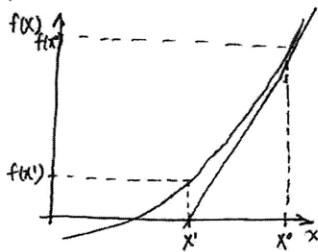
Evaluate $F(\vec{x}^0) \neq 0$

$$\approx F(\vec{x}^0) + J_F(\vec{x}^0)(\vec{x}^1 - \vec{x}^0)$$

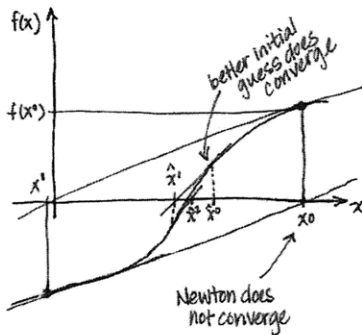
Use Approximation

$$J_F(\vec{x}^0)\Delta x = -F(\vec{x}^0)$$

update $\vec{x}^1 = \vec{x}^0 + \Delta x$



$$J_F(x) = A^{(1)} + A^{(2)}(\vec{x} \otimes I + I \otimes \vec{x})$$



NEWTON PROPERTIES:

1. Converges if you are close enough
 2. Converges very fast
- $J_F(\vec{x})$ is nonsingular

Original F:

$$F(\vec{x}) = A^{(1)}\vec{x} + A^{(2)}\vec{x} \otimes \vec{x} + B\vec{u}$$

Simpler F: $\vec{r} = 0$

$$F(\vec{x}, 0) = A^{(1)}\vec{x} + A^{(2)}(\vec{x} \otimes \vec{x})$$

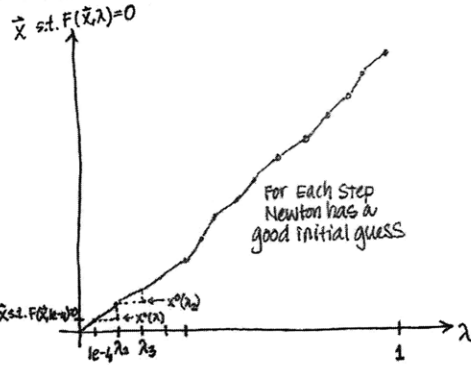
$$\vec{x} \text{ s.t. } F(\vec{x}, 0) = 0 \Rightarrow \vec{x} = 0$$

Continuation scheme:

$$F(\vec{x}, \lambda) = A^{(1)}\vec{x} + A^{(2)}\vec{x} \otimes \vec{x} + \lambda B\vec{u}$$

$$F(\vec{x}, 0) = A^{(1)}\vec{x} + A^{(2)}\vec{x} \otimes \vec{x} \Rightarrow \vec{x} = 0 \text{ \& } \lambda = 0$$

(original problem: $F(\vec{x}, 1) = 0$)



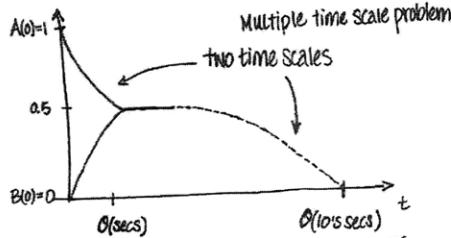
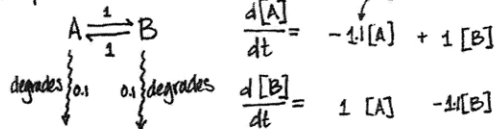
WANTED INFORMATION:

- How many stable steady-states
- What are the stable steady-states

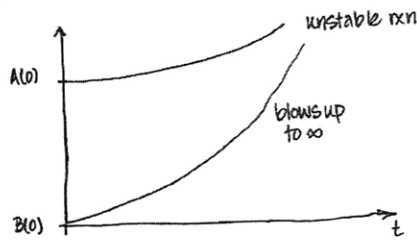
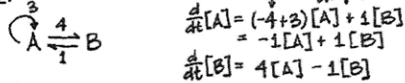
NEWTON GIVES:

- Maybe one steady-state

Example 1:



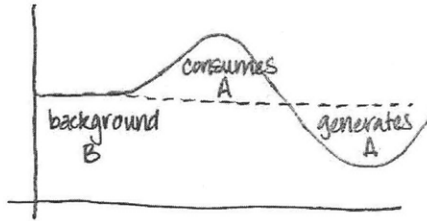
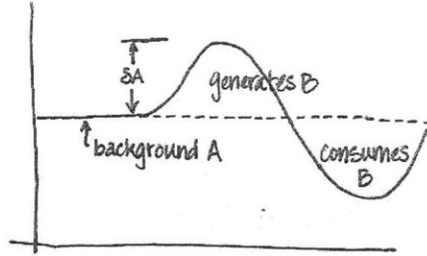
Example 2:



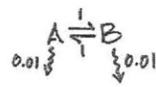
Example 3. Oscillatory System (nonphysical)

$$\frac{d}{dt}[A] = -1[B]$$

$$\frac{d}{dt}[B] = +1[A]$$

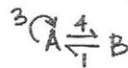


$$\frac{d}{dt} \vec{x} = \begin{bmatrix} -1.01 & 1 \\ 1 & -1.01 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$



$$\lambda_i = \begin{matrix} \begin{bmatrix} -1 \\ -1 \end{bmatrix} & \begin{bmatrix} 1 \\ 1 \end{bmatrix} \\ -2.01, -0.01 \\ \text{Fast Mode} & \text{Slow Mode} \end{matrix}$$

$$\frac{d}{dt} \vec{x} = \begin{bmatrix} -1 & 1 \\ 4 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$



$$\lambda_1 = -3 \quad \lambda_2 = 1$$

$$\frac{d}{dt} \vec{x} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

SA ? SB

$\lambda_1 = j \quad \lambda_2 = -j$
 oscillatory soln w/o decay \Rightarrow purely imaginary eigenvalues
 unstable

\rightarrow go to slides

2.16 Parameter Fitting and Estimation

6.581/BE.482
LECTURE 17: PARAMETER FITTING & ESTIMATION

THURSDAY
13 APRIL 2006

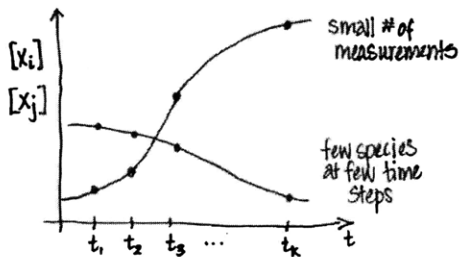
Dynamical Net Cell



known:
[u1]
[u2]

$$\frac{d}{dt} \begin{bmatrix} [X_1] \\ [X_2] \\ [X_3] \end{bmatrix} = A^{(1)}\vec{x} + A^{(2)}(\vec{x} \otimes \vec{x}) + B^{(1)}\vec{u} + B^{(2)}(\vec{x} \otimes \vec{x}) + B^{(3)}(\vec{x} \otimes \vec{u})$$

ignored for simplicity



What was $\vec{x}(0)$? (Initial Condition)

ESTIMATING INITIAL CONDITIONS

1) START by looking at the linear case

$$\dot{\vec{x}} = A^{(1)}\vec{x} + B^{(1)}\vec{u}$$

2) NONLINEAR case

$$\dot{\vec{x}} = A^{(1)}\vec{x} + A^{(2)}(\vec{x} \otimes \vec{x}) + B^{(1)}\vec{u}$$

LINEAR CASE

$$\dot{\vec{x}} = A\vec{x} + B\vec{u}$$

$$\vec{x}(t) = e^{At}x(0) + \int_0^t e^{A(t-\tau)} B u(\tau) d\tau$$

SCALAR CASE:

$$\dot{x} = ax + bu$$

$$x(t) = e^{at}x(0)$$

$$\dot{x} = ax + bu$$

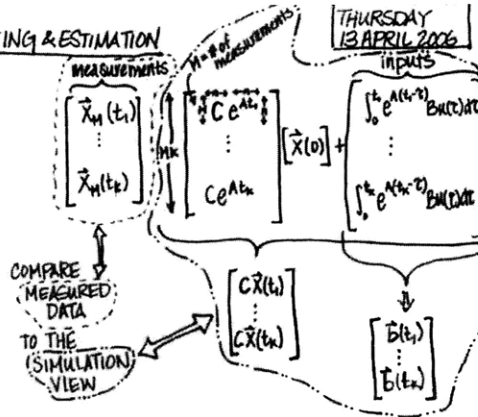
$$x(t) = e^{at}x(0) + \int_0^t e^{a(t-\tau)} bu(\tau) d\tau$$

$$[X_i(t_k)] = \frac{1}{C} \begin{bmatrix} 0 & \dots & 0 & 1 & 0 & \dots & 0 \\ \vdots & & & & & & \\ 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ \vdots & & & & & & \\ 0 & \dots & 0 & 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} X_1(t_k) \\ \vdots \\ X_n(t_k) \end{bmatrix}$$

i-th column

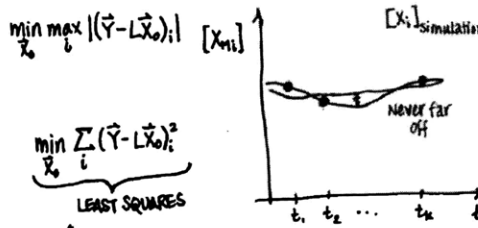
$$\begin{bmatrix} [X_i(t_k)] \\ [X_j(t_k)] \end{bmatrix} = \frac{1}{C} \begin{bmatrix} 0 & \dots & 0 & 1 & 0 & \dots & 0 \\ \vdots & & & & & & \\ 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ \vdots & & & & & & \\ 0 & \dots & 0 & 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} \vec{x}(t_k) \\ \vec{u}(t_k) \end{bmatrix}$$

i-th j-th



$$\min_{\vec{x}_0} \text{measure } \{ \vec{Y} - L\vec{x}_0 \}$$

$$\begin{bmatrix} X_m(t_1) \\ \vdots \\ X_m(t_k) \end{bmatrix} - \begin{bmatrix} \vec{b}(t_1) \\ \vdots \\ \vec{b}(t_k) \end{bmatrix} = \begin{bmatrix} C e^{A t_1} \\ \vdots \\ C e^{A t_k} \end{bmatrix} \vec{x}_0$$



$$\min_{\vec{x}_0} \sum_i (\vec{Y} - L\vec{x}_0)_i^2$$

LEAST SQUARES

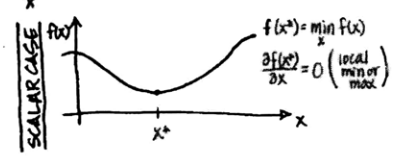
$$\min_{\vec{x}_0} \frac{(\vec{Y} - L\vec{x}_0)^T (\vec{Y} - L\vec{x}_0)}{\vec{e}^T \vec{e} = \sum \vec{e}_i^2}$$

f(x0)

- Dynamical system is linear $\dot{\vec{x}} = A^{(1)}\vec{x} + B^{(1)}\vec{u}$...
- Unknowns are initial conditions
- Measured quantities were linear functions of \vec{x} $C\vec{x}(t)$

ABSTRACT PROBLEM

$f(\vec{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ f : vector \rightarrow scalar
Cost function
 $\min_{\vec{x}} f(\vec{x}) \leftarrow$ optimization problem



VECTOR CASE

vector function of a vector

$$\frac{\partial f}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} = 0 \quad (\text{local min or max})$$

For some functions f
 $\exists \mathbf{x}^*$ s.t. $\frac{\partial f(\mathbf{x}^*)}{\partial \mathbf{x}} = 0$
 means \mathbf{x}^* is the global optimum

min $\vec{y} - L\vec{x}_0$ $f(\vec{x}_0) = (\vec{y} - L\vec{x}_0)^T (\vec{y} - L\vec{x}_0)$

$$\frac{\partial f}{\partial \vec{x}_0} = \frac{\partial}{\partial \vec{x}_0} [\vec{y}^T \vec{y} - 2\vec{y}^T L\vec{x}_0 + \vec{x}_0^T L^T L \vec{x}_0]$$

$$= 2L^T \vec{y} - 2L^T L \vec{x}_0 = 0$$

$$\Rightarrow L^T L \vec{x}_0 = L^T \vec{y}$$

NORMAL EQNS

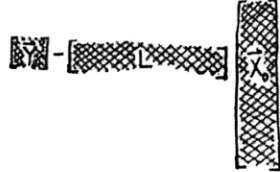
Lots of data $m \times n$

$$\begin{bmatrix} C_1^{Att} \\ \vdots \\ C_m^{Att} \end{bmatrix} \vec{x}_0 = \vec{y}$$

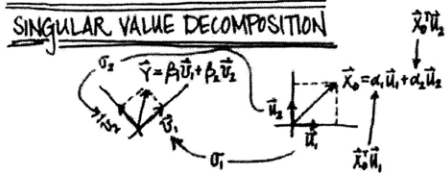
Limited data n \rightarrow typical

$$\begin{bmatrix} C_1^{Att} \\ \vdots \\ C_m^{Att} \end{bmatrix} \vec{x}_0$$

$$\underbrace{\begin{bmatrix} L^T & L \end{bmatrix}}_{n \times n} \underbrace{\begin{bmatrix} \vec{y} \\ \vec{x}_0 \end{bmatrix}}_{n \text{ length vector}} = \begin{bmatrix} L^T \vec{y} \\ \vec{x}_0 \end{bmatrix}$$



SINGULAR VALUE DECOMPOSITION



$$\begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_m \end{bmatrix} = \begin{bmatrix} \leftarrow u_1^T \rightarrow \\ \vdots \\ \leftarrow u_m^T \rightarrow \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \vec{x}_0$$

$$\begin{bmatrix} \uparrow v_1 \\ \vdots \\ \uparrow v_m \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_m \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_m \end{bmatrix}$$

$$\vec{y} = \begin{bmatrix} \uparrow v_1 \\ \vdots \\ \uparrow v_m \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_m \end{bmatrix} \begin{bmatrix} \leftarrow u_1 \rightarrow \\ \vdots \\ \leftarrow u_m \rightarrow \end{bmatrix} \vec{x}_0$$

$\vec{v}_i^T \vec{v}_j = 0$ for $i \neq j$
 (orthogonal coordinate system)

$\vec{u}_i^T \vec{u}_j = 0$ for $i \neq j$

$$\begin{bmatrix} \downarrow u_1 \\ \vdots \\ \downarrow u_m \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_m \end{bmatrix} \begin{bmatrix} \leftarrow v_1^T \rightarrow \\ \vdots \\ \leftarrow v_m^T \rightarrow \end{bmatrix} \vec{y} = \vec{x}_0$$



2.17 Parameter Estimation; Robustness, Fragility, Control

6.581/BE 482

THURSDAY
20 APRIL 2006

LECTURE 18: PARAMETER ESTIMATION; ROBUST, FRAGILITY, CONTROL

BRIEF SYNOPSIS:

- BIOCHEMICAL NETWORKS**
- information processing
 - decision & control
 - effector function

λ PHAGE

- 2 states; long-term stability
- multiple interactions among proteins & binding sites

DIFFERENTIAL EQUATION FORMULATION:

$$\frac{d\vec{x}}{dt} = F(\vec{x}, \vec{u}) = A\vec{x} + A^0(\vec{x}, \vec{u}) + B\vec{u} + C(\vec{u}, \vec{x}) + D(\vec{u}, \vec{x})$$

generally don't write

For steady states:

$$\frac{d\vec{x}}{dt} = 0 \Rightarrow F(\vec{x}, \vec{u}) = 0$$

Solve using **NEWTON'S METHOD**:

$$\frac{dF(\vec{x}, \vec{u})}{d\vec{x}} = \nabla_{\vec{x}} F \approx J_F(\vec{x})$$

$$F(\vec{x}^i; \vec{u}) = F(\vec{x}^0; \vec{u}) + \frac{dF}{d\vec{x}}(\vec{x}^0; \vec{u})(\vec{x}^i - \vec{x}^0) + H.o.T.$$

$J_F(\vec{x}^0)$

$$\text{Set } F(\vec{x}^i; \vec{u}) = 0 \Rightarrow J_F(\vec{x}^0)(\vec{x}^i - \vec{x}^0) = -F(\vec{x}^0; \vec{u})$$

solve for \vec{x}^i

$$\text{iterate: } J_F(\vec{x}^i)(\vec{x}^{i+1} - \vec{x}^i) = -F(\vec{x}^i; \vec{u})$$

RUN SIMULATIONS:

- Forward Euler - explicit
- Backward Euler - implicit (iterative)

PARAMETER SENSITIVITIES:

$$\min_{\Delta A} f(\Delta A) = \min_{\Delta A} \sum_{i=1}^n [C^T \vec{x}^i(\Delta A) - C^T \vec{x}^i_{measured}]^2$$

of experiments [0.0000000]

OPTIMIZE:

- GRADIENT DESCENT

$$\Delta A^k = \Delta A^{k-1} + \beta \nabla_{\Delta A} f$$

pick stepsize β s.t. $\min f(\Delta A^k + \beta \nabla_{\Delta A} f)$

This method does not work well

- ADJOINT OPTIMIZATION

Set gradient to zero

$$\frac{\partial}{\partial \Delta A} [C^T \vec{x}^i(\Delta A) - C^T \vec{x}^i_{measured}]^2 = 2 [C^T \vec{x}^i(\Delta A) - C^T \vec{x}^i_{measured}] \left[\frac{\partial}{\partial \Delta A} C^T \vec{x}^i(\Delta A) \right]$$

straightforward difficult

$$\frac{\partial}{\partial \Delta A} C^T \vec{x}^i(\Delta A) = \frac{\partial}{\partial \Delta A} C^T \left[\frac{\partial \vec{x}^i(\Delta A)}{\partial \Delta A} \right]$$

n = # of concentrations m = # of parameters

$$\frac{dF}{d(\Delta A)} = 0 \Rightarrow \frac{\partial F}{\partial \Delta A} + \frac{\partial F}{\partial \vec{x}} \frac{\partial \vec{x}}{\partial \Delta A} = 0 \Rightarrow \frac{\partial F}{\partial \vec{x}} \frac{\partial \vec{x}}{\partial \Delta A} = - \frac{\partial F}{\partial \Delta A}$$

(*) solve for this term & plug-in HARDER TO SOLVE

Rather than solve the previous problem, solve its **DUAL**:

$$\frac{\partial}{\partial \Delta A} C^T \vec{x}^i(\Delta A) = \vec{z}^T \left(\frac{\partial F}{\partial \Delta A} \right)$$

where \vec{z} is obtained from:

$$n \left[\left(\frac{\partial F}{\partial \vec{x}} \right)^T \vec{z} \right] = - \frac{\partial F}{\partial \Delta A}$$

EASIER TO SOLVE

SHOW THAT THE DUAL IS EQUIVALENT:

$$C^T \frac{\partial \vec{x}^i(\Delta A)}{\partial \Delta A} = C^T \left(- \frac{\partial F}{\partial \vec{x}} \right)^{-1} \left(\frac{\partial F}{\partial \Delta A} \right)$$

$$\Rightarrow \frac{\partial}{\partial \Delta A} C^T \vec{x}^i(\Delta A) = \vec{z}^T \frac{\partial F}{\partial \Delta A}$$

$$\vec{z}^T = C^T \left(- \frac{\partial F}{\partial \vec{x}} \right)^{-1}$$

$$\Rightarrow \vec{z}^T \left(\frac{\partial F}{\partial \vec{x}} \right) = - C^T$$

$$\Rightarrow \left(\frac{\partial F}{\partial \vec{x}} \right)^T \vec{z} = - C \quad \square$$

*Finding steady states

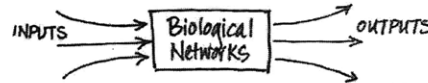
*Running Simulations

- Model Construction

parameter estimation

can also be used in design mode

parameter sensitivities

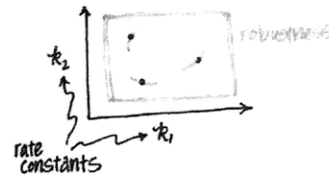


- absence of glucose
- presence of galactose

- synthesize multiple enzymes for galactose utilization
- turns on galactose transporters

PROBLEMS IN BIOLOGY:

- Difficulty in maintaining constant concentration
 - cell volume changes
 - small counts make uniformity difficult
- Highly variable environment
 - Temperature
 - Humidity
- Certain functions persist across all stages of development, cell cycle, tissue types
- Evolve



2.18 2-D and 3-D Light Microscopy; Image Reconstruction

6.581/BE.482

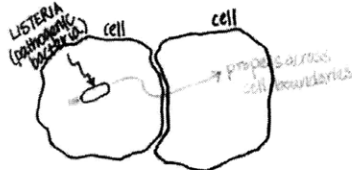
TUESDAY
25 APRIL 2006

LECTURE 19: 2-D & 3-D LIGHT MICROSCOPY; IMAGE RECONSTRUCTION

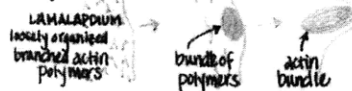
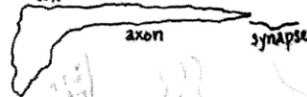
WHAT TYPES OF PHENOMENA MIGHT ONE STUDY?

MOVIES - Time dependent

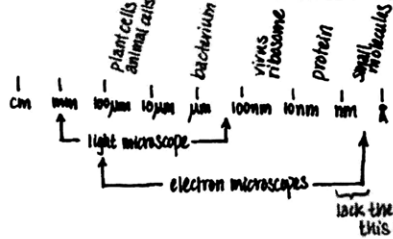
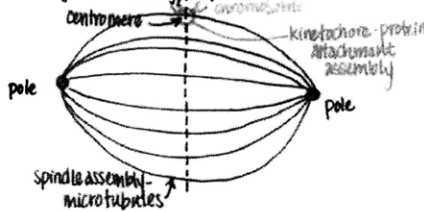
(1) Julie Theriot (Stanford) - Invasive Bacteria



(2) Frank Gertler (MIT Biology) nerve cell

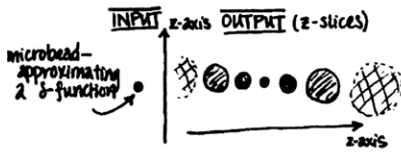
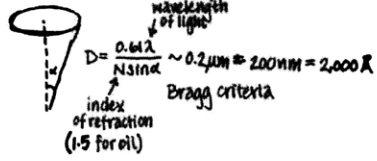


(3) Pete Sorger (MIT Biology/BE)



look the lenses to achieve this range

Resolution limit: How close can two points be & still be resolved?

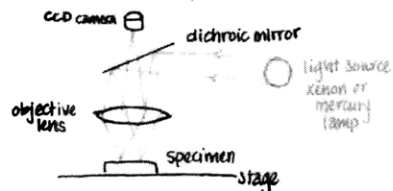


WHY?
1) objective lens & aperture don't collect "all" the light
2) Diffraction off edge of aperture (ring pattern)

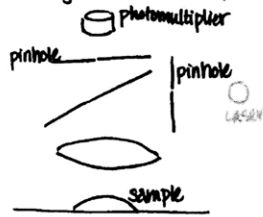


THE PHYSICAL MICROSCOPE - FLUORESCENCE

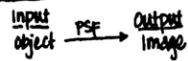
(1) Wide-field microscope



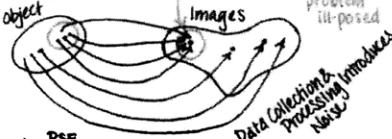
(2) Scanning Confocal Microscope



THE BASIC PROBLEM:



Mathematical procedure for (L → R) is well defined. However, we need to solve the inverse problem which is generally not well-posed



Object \xrightarrow{PSF} Noise-Free Image $\xrightarrow{\text{Data Collection \& Processing (introduces Noise)}}$ Actual Image

Noise-free image is a convolution of the object with a PSF

$$g^*(\vec{x}) = \int K(\vec{x}, \vec{x}') f^*(\vec{x}') d\vec{x}' \quad \text{if } f^*(\vec{x}) \text{ is } \delta(\vec{x} - \vec{x}_0) \Rightarrow g^*(\vec{x}) = K(\vec{x}, \vec{x}_0)$$

Actual image $g(\vec{x}) = g^*(\vec{x}) + n(\vec{x}) = \int K(\vec{x}, \vec{x}') f(\vec{x}') d\vec{x}' + n(\vec{x})$

Often we treat the PSF as $K(\vec{x}, \vec{x}') = K(\vec{x} - \vec{x}')$

2.19 Deconvolution

6.581/BE.402

LECTURE 20: DECONVOLUTION

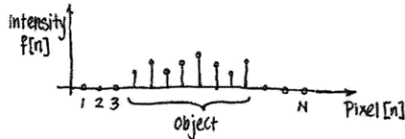
THURSDAY
27 APRIL 2006

DECONVOLUTION

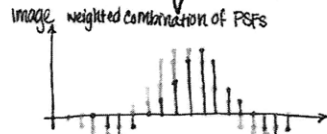
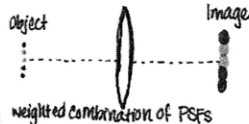
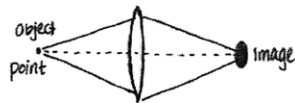
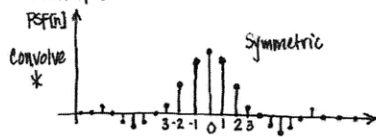
- 1.) Point spread function representation
- 2.) Straightforward deconvolution
- 3.) Handling noise (SVD)

1-D

Object (black & white)



Point Spread Function (PSF)

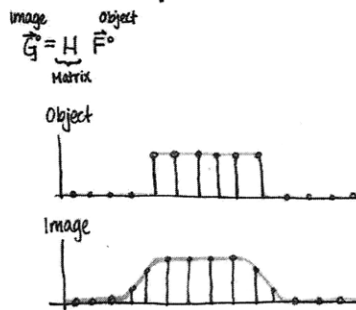
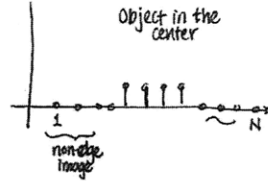


$$\begin{bmatrix} G[0] \\ \vdots \\ G[k] \\ \vdots \\ G[N] \end{bmatrix} = \begin{bmatrix} \text{PSF}[0] & \text{PSF}[1] & \dots & \text{PSF}[N-1] \\ \vdots & \vdots & \ddots & \vdots \\ \text{PSF}[0] & \text{PSF}[1] & \dots & \text{PSF}[N-1] \\ \vdots & \vdots & \ddots & \vdots \\ \text{PSF}[0] & \text{PSF}[1] & \dots & \text{PSF}[N-1] \end{bmatrix} \begin{bmatrix} F[0] \\ \vdots \\ F[k] \\ \vdots \\ F[N] \end{bmatrix}$$

assume periodicity
pretend object repeats periodically

OBJECT $\begin{bmatrix} k-1 \\ k \\ k+1 \end{bmatrix}$
 ↓
 PIXELS IMAGE $\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$

K-1's contribution at K is the +1 contribution from K-1 → PSF[K]



$$\vec{G} = H \vec{F}$$

Matrix
Object
Image

When the smallest bit of noise is added to the image & then reconstructed, the resulting reconstructed image is garbage ⇒ very fragile

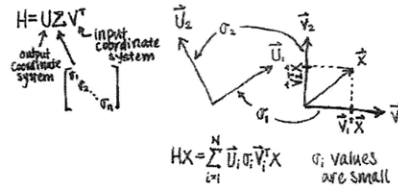
Noise

$$\vec{G} = \vec{G} + \vec{w} = H \vec{F} + \vec{w}$$

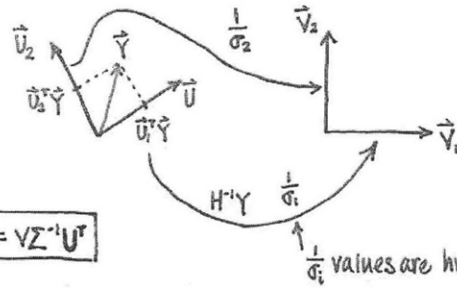
vector of noise

$$\vec{F} = H^{-1} (H \vec{F} + \vec{w}) = \vec{F} + H^{-1} \vec{w}$$

How large?



$$H = U \Sigma V^T$$

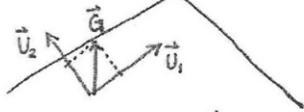


$$H^{-1} = V \Sigma^{-1} U^T$$

$$\vec{G}^o = H \vec{F}^o = \sum_i \hat{u}_i \sigma_i \hat{v}_i^T \vec{F}^o = \sigma_1 \hat{v}_1^T \vec{F}^o \begin{bmatrix} 1 \\ 0 \\ \vdots \end{bmatrix} + \sigma_2 \hat{v}_2^T \vec{F}^o \begin{bmatrix} 0 \\ 1 \\ \vdots \end{bmatrix} + \dots + \sigma_N \hat{v}_N^T \vec{F}^o \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

Perfect Reconstruction (PR)

$$\vec{F}^{PR} = \hat{v}_1^T \vec{F}^o \begin{bmatrix} 1 \\ 0 \\ \vdots \end{bmatrix} + \hat{v}_2^T \vec{F}^o \begin{bmatrix} 0 \\ 1 \\ \vdots \end{bmatrix} + \dots$$



$$\vec{G} = \vec{G}^o + \vec{w} = \sum_i \hat{u}_i \sigma_i \hat{v}_i^T \vec{F}^o + \vec{w}$$

$$\vec{F}^R = \sum_i \hat{v}_i \frac{1}{\sigma_i} \hat{u}_i^T \vec{G}$$

$$\vec{F}^R = \sum_i \hat{v}_i \frac{1}{\sigma_i} (\sigma_i \hat{v}_i^T \vec{F}^o + \hat{u}_i^T \vec{w})$$

$$= \sum_i \left[\underbrace{\hat{v}_i \hat{v}_i^T \vec{F}^o}_{\text{perfect reconstruction}} + \hat{v}_i \frac{1}{\sigma_i} (\hat{u}_i^T \vec{w}) \right]$$

$$\vec{F}^R = (\hat{v}_1 \hat{v}_1^T \vec{F}^o + \hat{v}_1 \frac{1}{\sigma_1} \hat{u}_1^T \vec{w}) + \hat{v}_2 (\hat{v}_2^T \vec{F}^o + \frac{1}{\sigma_2} \hat{u}_2^T \vec{w}) + \dots + \hat{v}_N (\hat{v}_N^T \vec{F}^o + \frac{1}{\sigma_N} \hat{u}_N^T \vec{w})$$

→ Remove terms with small σ

2.20 Deconvolution II

OUTLINE

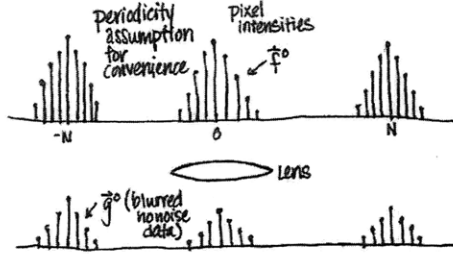
- 1) Reminder about SVD
- 2) Fourier (Freq. Domain) View
- 3) 2D Case

REMINDER object

$$\vec{g} = H\vec{f}^o + \vec{w}$$

↑ data ↑ effects of imaging system ↑ noise

IN 1D



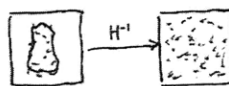
$$\begin{bmatrix} g^o[1] \\ \vdots \\ g^o[N] \end{bmatrix} = \begin{bmatrix} h[0] & h[N-1] & h[N-2] & \dots & h[2] & h[1] \\ h[1] & h[0] & h[N-1] & \dots & h[3] & h[2] \\ h[2] & h[1] & h[0] & \dots & h[4] & h[3] \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ h[N-1] & h[N-2] & h[N-3] & \dots & h[4] & h[3] \end{bmatrix} \begin{bmatrix} f^o[1] \\ \vdots \\ f^o[N] \end{bmatrix}$$

CIRCULANT MATRIX

1st DECONVOLUTION IDEA

$$\vec{f}^{BE} = H^{-1}\vec{g} = \vec{f}^o + H^{-1}\vec{w}$$

(bad estimate)



Using SVD

$$H = U \begin{bmatrix} \sigma_1 & & \\ & \sigma_n & \\ & & 0 \end{bmatrix} V^T$$

orthonormal (define coordinate system)

$$\vec{g} = \sum \tilde{u}_i \sigma_i (\tilde{v}_i^T \vec{f}^o) + \vec{w}$$

$$\vec{f}^{BE} = \sum \frac{1}{\sigma_i} \tilde{v}_i (\tilde{u}_i^T \vec{g})$$

$$\vec{f}^{BE} = \sum \tilde{v}_i [(\tilde{v}_i^T \vec{f}^o) + \frac{1}{\sigma_i} \tilde{u}_i^T \vec{w}]$$

6.581/BE.482
LECTURE 21: DECONVOLUTION II

TUESDAY
2 MAY 2006

ESTIMATOR

$$\vec{f}^{SE} = \alpha_1 (\tilde{v}_1^T \vec{f}^o + \frac{1}{\sigma_1} \tilde{u}_1^T \vec{w}) \tilde{v}_1 + \alpha_2 (\tilde{v}_2^T \vec{f}^o + \frac{1}{\sigma_2} \tilde{u}_2^T \vec{w}) \tilde{v}_2 + \dots + \alpha_n (\tilde{v}_n^T \vec{f}^o + \frac{1}{\sigma_n} \tilde{u}_n^T \vec{w}) \tilde{v}_n$$

(simple estimate)

$\alpha_i = \text{either zero or one}$

image noise image noise

want to keep terms where image is much larger than the noise

$$\text{If } |\tilde{v}_i^T \vec{f}^o| > \text{thresh} \left| \frac{1}{\sigma_i} \tilde{u}_i^T \vec{w} \right| \Rightarrow \alpha_i = 1$$

$$\text{Else } \alpha_i = 0$$

CONNECTION TO FREQUENCY RESPONSE

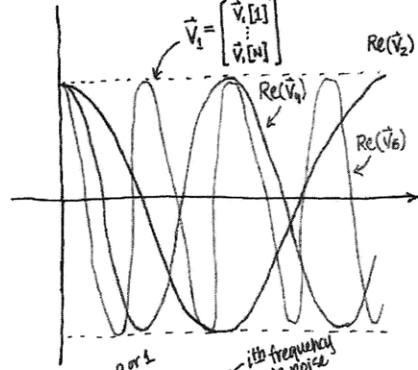
$$H_{\text{CIRCULANT}} = \text{DFT}^{-1} \begin{bmatrix} \times & & & \\ & \times & & \\ & & \times & \\ & & & \times \end{bmatrix} \text{DFT}$$

diagonal discrete Fourier transform

Complex numbers

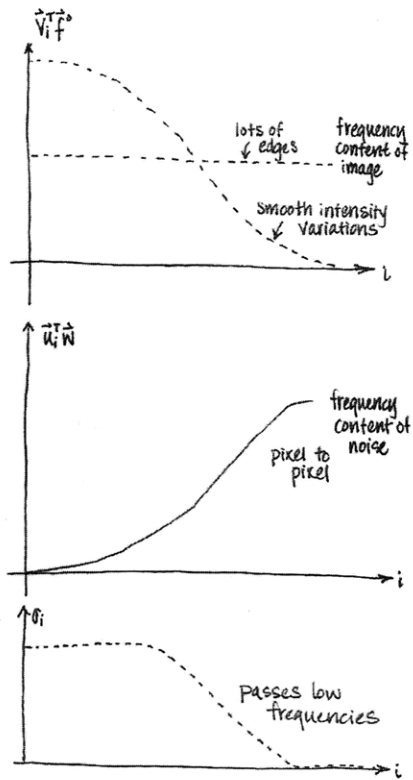
$$\frac{1}{N} \begin{bmatrix} 1 & & & \\ & e^{j2\pi k(1-0)/N} & & \\ & & \ddots & \\ & & & e^{j2\pi k(N-1)/N} \end{bmatrix} \begin{bmatrix} X[1] \\ \vdots \\ X[N] \end{bmatrix} = \begin{bmatrix} x[1] \\ \vdots \\ x[N] \end{bmatrix}$$

$$x[n] = \sum_{k=1}^N e^{j2\pi k(n-1)/N} X[k]$$

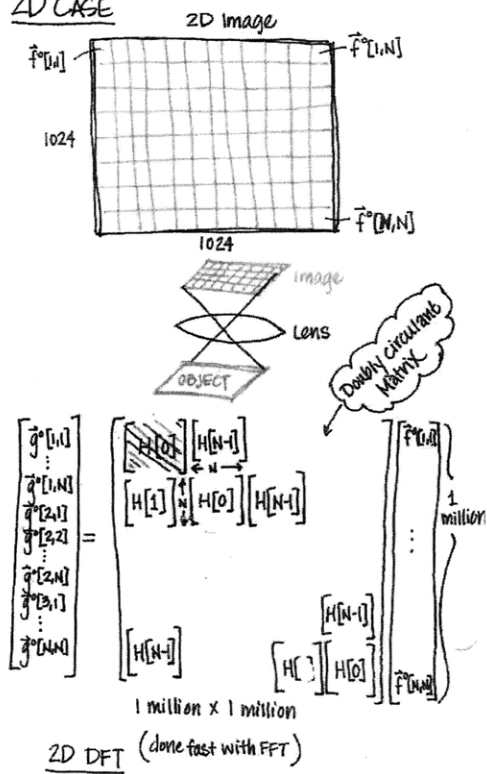


$$\vec{f}^E = \sum \alpha_i (\tilde{v}_i^T \vec{f}^o + \frac{1}{\sigma_i} \tilde{u}_i^T \vec{w}) \tilde{v}_i$$

0 or 1 i-th frequency in noise i-th frequency in object

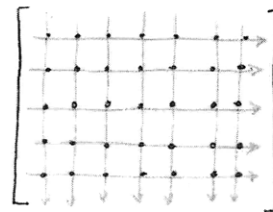


2D CASE



2D DFT (done fast with FFT)

$$X[n_1, n_2] = \frac{1}{N^2} \sum_k e^{-j2\pi \frac{(n_1-1)(k_1-1)}{N}} e^{-j2\pi \frac{(n_2-1)(k_2-1)}{N}} X[k_1, k_2]$$



Which terms to keep

$$\tilde{F}^E = \sum \alpha_i \tilde{V}_i (\tilde{V}_i^T \tilde{F}^0 + \frac{1}{\sigma_i} \tilde{u}_i^T \tilde{w})$$

is $\alpha_i \pm$ or 0 ?

$$\tilde{V}_i^T \tilde{F}^0 \text{ versus } \frac{1}{\sigma_i} \tilde{u}_i^T \tilde{w}$$

Suppose we know the "energy" in the noise: $\tilde{w}^T \tilde{w} = E_{noise}$

& suppose we also know the "energy" in the image: $\tilde{F}^0{}^T \tilde{F}^0 = E_{image}$

\Rightarrow energy in the image is uniform in frequency

$$\tilde{F}^0 = \sum \tilde{V}_i (\tilde{V}_i^T \tilde{F}^0)$$

$$\tilde{F}^0{}^T \tilde{F}^0 = (\sum \tilde{V}_i (\tilde{V}_i^T \tilde{F}^0))^T (\sum \tilde{V}_i (\tilde{V}_i^T \tilde{F}^0))$$

$$|\tilde{V}_i^T \tilde{F}^0| = \frac{\sqrt{E_{image}}}{N}$$

Uniform noise assumption: $|\tilde{u}_i^T \tilde{w}| = \frac{\sqrt{E_{noise}}}{N}$

2.21 Blind Deconvolution

6.581/EE.482
LECTURE 22: BLIND DECONVOLUTION

THURSDAY
4 MAY 2006

$$g(\vec{x}) = \int k(\vec{x}-\vec{x}') f(\vec{x}') d\vec{x}' + w(\vec{x})$$

$(f, k) \rightarrow$ provides infinite number of possible solutions

What types of constraints apply?

- non-negativity
- symmetry on k
- frequency space expected for object
- support for PSF \rightarrow zero beyond some distance
- maximum intensity for object

Noiseless: $g(\vec{x}) = \int k(\vec{x}-\vec{x}') f(\vec{x}') d\vec{x}'$

$$C_g \equiv \{(u,v) : u * v = g\} \rightarrow (f, k) \in C_g$$

convolution

$$C_f \equiv \{(u,v) : u \text{ satisfies constraints on } f\}$$

$$C_k \equiv \{(u,v) : v \text{ satisfies constraints on } k\}$$

Let $C_0 = C_g \cap C_f \cap C_k$

To find a solution in C_0 :

Iterative projections:

Let P_g be an operator that projects into C_g

& P_f operator projects into C_f

& P_k operator projects into C_k

$$(f, k)_{i+1} = P_k P_f P_g (f, k)_i$$

initial guess $(f, k)_0$
 \rightarrow iterations converge to point in C_0

Projection operator: often a form of minimization

minimize $\|(u,v) - (f, k)\|$ subject to constraints
 solution \uparrow operand \uparrow applied with Lagrange multipliers \uparrow $(u,v) \in C_g$ or C_f or C_k

$$\text{Let } J_{(u,v)} \equiv \|g - u * v\|^2 = \int [g(\vec{x}) - (u * v)(\vec{x})]^2 d\vec{x} \geq 0$$

measure of error to measured image
 $(u,v) \in C_g \Rightarrow J_{(u,v)} = 0$
 $C_f \& C_k$

1) initial guess: (f_0, k_0)

$$f_0 = P_f g, [k_0 = 0]$$

2) Solve for $k_i = \arg \min_{k \in C_k} J_{(f_{i-1}, k)}$

3) Solve for $f_i = \arg \min_{f \in C_f} J_{(f, k_i)}$

4) increment $i \rightarrow i+1$

Poorer Alternative

$$E = J_{(f,k)} + E_{\text{restraint}} k(\vec{x}) = 0 \text{ for } \Omega_k$$



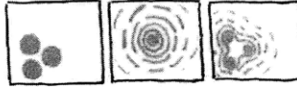
$$E_{\text{restraint}} = \int_{\Omega_f} |f(\vec{x})|^2 d\vec{x} + \int_{\Omega_k} |k(\vec{x})|^2 d\vec{x}$$

$g(\vec{x})$ is invariant to $(f, k) \rightarrow (\alpha f, \frac{k}{\alpha})$

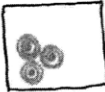
SAMPLE PROBLEM:

$$f(\vec{x}) * k(\vec{x}) \rightarrow g(\vec{x})$$

input:

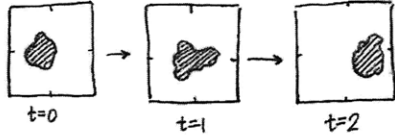


output:

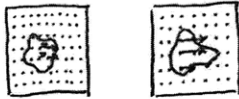


Yang et al. J. Opt. Soc. Am A
11: 2401 (1994)

2.22 Optical Flow



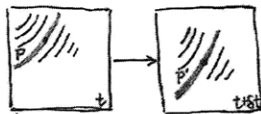
DESCRIPTION:
 • movement of CoM
 • deformation described as movement expansion
 ALTERNATIVE DESCRIPTION:
 • vector field flow



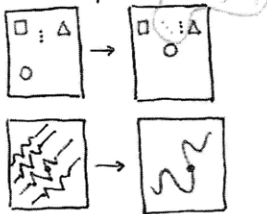
OPTICAL FLOW: measuring "motion" in visual field through changes in image across time
 - usually, but not always, equivalent to motion of objects

2 PATHOLOGICAL CASES (exceptions, rather than rule)

1. uniform Gray Sphere
 - optical flow zero
 - but object is moving
2.
 - moving light
 - optical flow
 - but object is still



- 2 ways of tracking motion:
1. Global triangulation
 2. Local relationships



Imagine a series of 2D images with t as 3rd dimension
 Intensity of each pixel: $E(x,y,t)$
 Our flow field: $(u(x,y), v(x,y))$ at time t

Most naive assumption
 At $t+\delta t$, the location where some index pixel moves to has the same intensity as its index location at time t
 $E(x,y,t) = E(x+\delta x, y+\delta y, t+\delta t)$ $\delta x = u\delta t$
 $= E(x+u\delta t, y+v\delta t, t+\delta t)$ $\delta y = v\delta t$

- Note:
1. At best, only true for small δt
 2. Later will deal with issue that intensity may be different in new location
 3. Need at least one more constraint to solve for (u,v)

Re-express 1st constraint:
 $E(x,y,t) = E(x,y,t) + \frac{\partial E}{\partial x} \delta x + \frac{\partial E}{\partial y} \delta y + \frac{\partial E}{\partial t} \delta t + \dots$
 $\frac{\partial E}{\partial x} \delta x + \frac{\partial E}{\partial y} \delta y + \frac{\partial E}{\partial t} \delta t = 0$
 $\Rightarrow E_x u + E_y v + E_t = 0$

Express a "smoothness" constraint
 Error function for smoothness

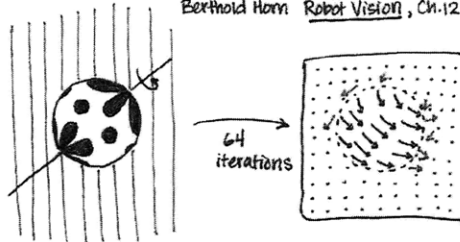
$$e_s = \iint [(u_x^2 + u_y^2) + (v_x^2 + v_y^2)] dx dy \quad u_x = \frac{\partial u}{\partial x}$$

$$e_c = \iint (E_x u + E_y v + E_t) dx dy$$

We will minimize: $e_s + \lambda e_c$

Minimize $\iint F dx dy$
 where $F = [(u_x^2 + u_y^2) + (v_x^2 + v_y^2)] + \lambda (E_x u + E_y v + E_t)$
 if $E_x = E_y = 0$
 LAPLACE'S EQN FORM:
 $\nabla^2 u + \nabla^2 v = 0$

For mathematics & details:
 Berthold Horn Robot Vision, Ch.12



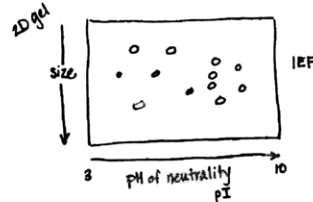
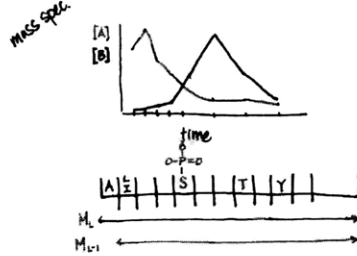
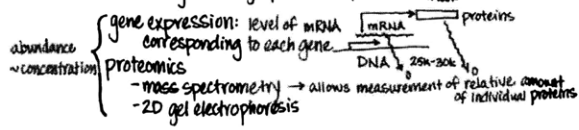
2.23 High-Throughput Data and Analysis

6.581 / BE.482
LECTURE 24: HIGH-THROUGHPUT DATA & ITS ANALYSIS

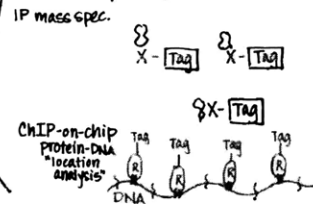
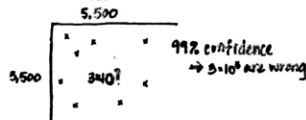
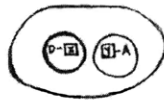
THURS
11 MAY 2006

Data Mining

- multichannel - massively parallel
- system wide
- possibility of closely spaced time points

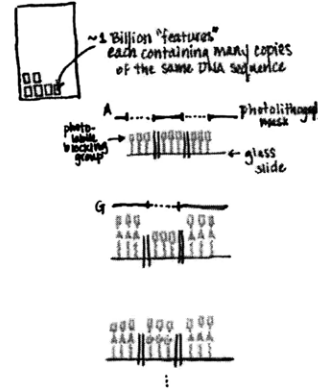


- activity - enzyme activity assay
- network structure - yeast 2-hybrid experiment



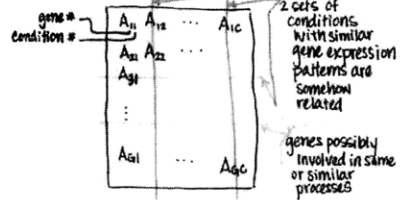
GFP (green fluorescent protein) tagged proteins
↑ localization

Affymetrix



R. Hughes... SH Friend Cell 102, 109-126 (2000)

- 300 sets of "conditions"
- some were chemically treated
- most were deletion strains

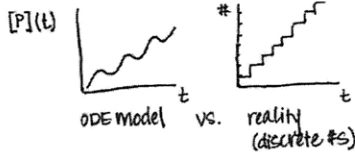


2.24 Inference and Statistics

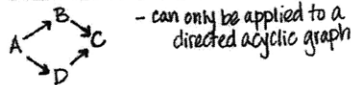
6.531/BE.482

LECTURE 25: INFERENCE & STATISTICS

TUESDAY
16 MAY 2006



BAYESIAN NETWORK ANALYSIS



Example:



P(C) T F
0.5 0.5

P(S C) T F
T 0.1 0.9
F 0.5 0.5

P(R C) T F
T 0.8 0.2
F 0.2 0.8

P(W S,R) T F
T T 0.99 0.01
F T 0.9 0.1
T F 0.9 0.1
F F 0.0 1.0

W=T
The chain for probability
 $P(C,S,R,W) = P(C) \cdot P(S|C) \cdot P(R|S,C) \cdot P(W|S,R)$
 current model topology $\rightarrow P(C) \cdot P(S|C) \cdot P(R|C) \cdot P(W|S,R)$

BAYES THEOREM: $P(A,B) = P(A|B)P(B)$
 $P(A|B) = \frac{P(A,B)}{P(B)}$

$$P(S=T|W=T) = \frac{P(S=T, W=T)}{P(W=T)} = \frac{\sum_{C,R} P(C,S=T,R,W=T)}{P(W=T)} = 0.4298$$

$$P(R=T|W=T) = \frac{P(R=T, W=T)}{P(W=T)} = \frac{\sum_{C,S} P(C,S,R=T,W=T)}{P(W=T)} = 0.7079$$

$$P(W=T) = 0.6471$$

Experiments give P(Data)
 Models can produce P(Data|Model)
 What we want to compare P(Model|Data)

$$P(\text{Model}|\text{Data}) = \frac{P(\text{Data}|\text{Model})P(\text{Model})}{P(\text{Data})}$$

$$\text{Score}_i = \log P(\text{Model}_i|\text{Data}) = \log P(\text{Data}|\text{Model}_i) + \log P(\text{Model}_i) - \log P(\text{Data})$$

← Prior knowledge about system & consistent across all models

Sachs, Gifford, Jaakkola, Sorger, & Lauffenburger
<http://www.stke.org/cgi/content/full/sigtrans;2002/148/p1>

stimulates (cue) - fibronectin (fn)

(F) Focal Adhesion Kinase

(E) Extracellular signal-related kinase

Experiment was carried out twice

M0: Cue → (F) → (E)	score	
	Data Set I	Data Set II (fold-change)
M1: Cue → (F) → (E)	M0: -50.5 (258)	-60.8 (1700)
M2: Cue → (E) → (F)	M1: -42.8* (1)	-58.0 (104)
M3: Cue → (F) → (E)	M2: -43.8 (3)	-57.6 (73)
M4: Cue → (E) → (F)	M3: -44.3 (4)	-55.6 (10)
	M4: -44.3 (4)	-53.4* (1)

Simulation of Reaction Kinetics:

Continuous vs Stochastic (ODE)

$$\frac{dx}{dt} = f(x, u)$$

Alternative framework

$$\vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \leftarrow \begin{array}{l} \# \text{ of molecules} \\ \text{of species } i \end{array}$$

Think about dynamics as:

- 1) system remains unchanged for time τ
- 2) system will change by a single chemical reaction, μ

Let $P(\tau, \mu) \delta\tau$ = probability that \vec{y} at time t will evolve such that no reaction occurs in $(t, t+\tau)$, but that the next reaction occurs in $(t+\tau, t+\tau+\delta\tau)$ and, it is Reaction μ = reaction probability density function

discrete list of reactions

$$= P_0(\tau) \cdot a_\mu \delta\tau$$

probability that waiting time is τ probability that μ is the reaction

It is shown that

$$P(\tau, \mu) = a_\mu e^{-a_0 \tau} \text{ for } 0 \leq \tau < \infty$$

where $a_\mu \in \{1, \dots, M\}$

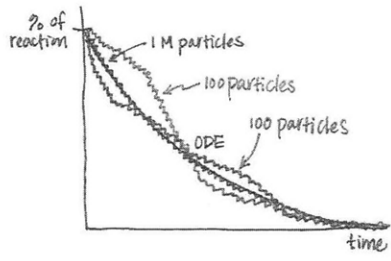
$$a_0 = \sum_{i=1}^M a_i \text{ individual reaction propensities}$$

$$a_\mu = h_\mu c_\mu \text{ reaction rate} \times \# \text{ of collision types}$$

DT Gillespie J. Phys Chem 81. 2340-61 (1977)
 J. Phys Chem 115. 1716-33 (2001)

Pick 2 random numbers R_1, R_2 uniformly on $[0, 1]$

$$\tau = \frac{1}{a_0} \ln \frac{1}{R_1} \text{ select } \mu \text{ such that } \sum_{i=1}^{\mu} a_i < R_2 a_0 < \sum_{i=1}^{\mu+1} a_i$$



Chapter 3

Introduction to Numerical Simulation

Introduction to Numerical Simulation teaches an immense amount of material, it's a general overview for graduate level numerical analysis and covers everything from how to formulate equations based on a model to procedures one should use when solving linear equations. This chapter offers in depth analysis as to why these procedures work and allows the reader to look at problems from several different perspectives. The material is presented in an easy to follow manner with several examples to elucidate each point.

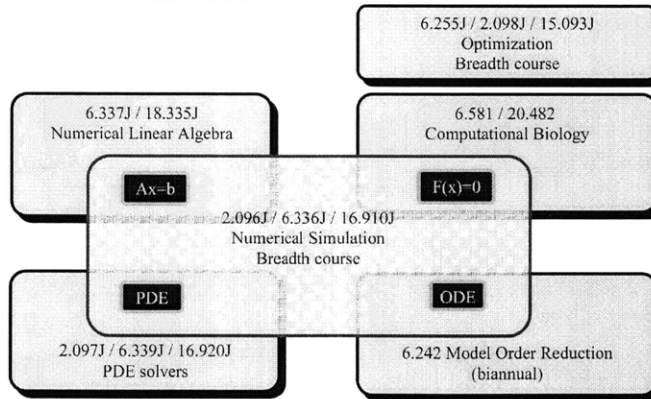
INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 1. Example Problems and Basic Equations

COURSE OUTLINE:

- Assembling system of equations automatically [3 lec]
- Steady state solution
 - Linear Problems: $Ax=b$ [7 lec]
 - Non-Linear Problems: $F(x)=0$ [3 lec]
- Dynamics (ODE solvers)
 - Time domain integration [3 lec]
 - Periodic steady state [2 lec]
- PDE solvers
 - Integral Equation Methods [2 lec]
 - Finite Element Methods (FEM) [1 lec]
 - Finite Difference Methods (FD) [2 lec]
 - Preconditioners for PDE solvers [1 lec]
- Model Order Reduction [2 lec]

RELATION TO OTHER COURSES

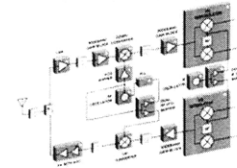


TODAY'S OUTLINE:

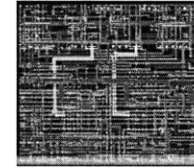
- Uses For Simulation
 - Engineering Design
 - Virtual Environments
 - Model Verification
- Course Philosophy
- Example Problems
 - Power distribution on an Integrated Circuit [Electrical]
 - Load bearing on a space frame [Structural]
 - Temperature distribution in an engine [Thermal]

USES FOR SIMULATION

- Circuit Analysis (e.g. cell phone)



From www.maxim.com



- Equations
 - Current-voltage relations for circuit elements (resistors, capacitors, transistors, inductors), current balance equations
- Recent Developments
 - Matrix-Implicit Krylov Subspace methods.

- Electromagnetic Analysis of Packages

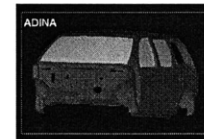
- Equations
 - Maxwell's Partial Differential Equations
- Recent Developments
 - Fast Solvers for Integral Formulations



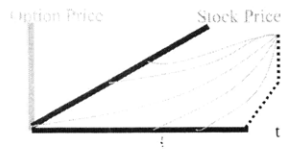
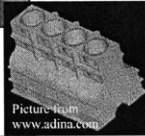
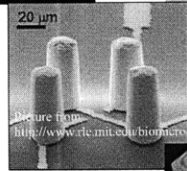
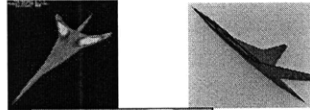
Thanks to Microcosm Inc. now Coventor

- Structural Analysis of Automobiles

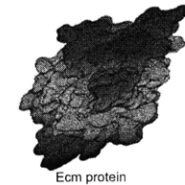
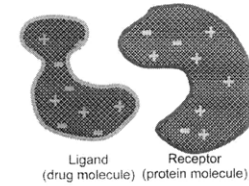
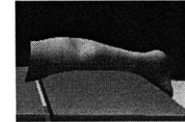
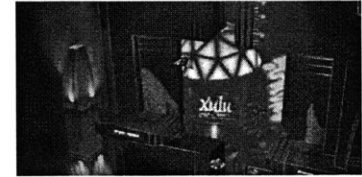
- Equations
 - Force-displacement relationships for mechanical elements (plates, beams, shells) and sum of forces = 0.
 - Partial Differential Equations of Continuum Mechanics
- Recent Developments
 - Meshless Methods, Iterative Methods, Automatic Error Control



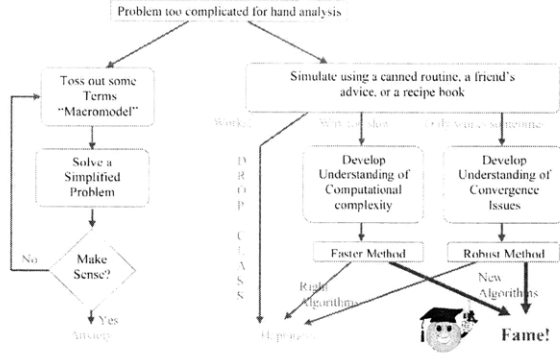
- Drag Force Analysis of Aircraft
 - Equations
 - Navier-Stokes Partial Differential Equations
 - Recent Developments
 - Multigrid Methods for Unstructured Grids
- Analysis of Cell Traps for Sorting Cytometry
 - Equations
 - Navier-Stokes Partial Differential Equations
 - Recent Developments
 - Multigrid Methods for Unstructured Grids
- Engine Thermal Analysis
 - Equations
 - Poisson Partial Differential Equations
 - Recent Developments
 - Fast Integral Equation Solvers, Monte-Carlo Methods
- Micromachined Device Performance Analysis
 - Equations
 - Elastomechanics, Electrostatics, Stokes Flow
 - Recent Developments
 - Fast Integral Equation Solvers, Matrix-Implicit Multi-level Newton Methods for coupled domain problems
- Stock Option Pricing for Hedge Funds
 - Equations
 - Black-Scholes Partial Differential Equation
 - Recent Developments
 - Financial Service Companies are hiring engineers, mathematicians, and physicists



- Virtual Environments for Computer Games
 - Equations
 - Multibody Dynamics, Elastic Collision Equations
 - Recent Developments
 - Multirate integration methods, parallel simulation
- Virtual Surgery
 - Equations
 - Partial Differential Equations of Elastomechanics
 - Recent Developments
 - Parallel Computing, Fast Methods
- Biomolecule Electrostatic Optimization
 - Equations
 - The Poisson Partial Differential Equation
 - Recent Developments
 - Matrix-Implicit Iterative Methods, Fast Integral Equation Solvers



THE COMPUTER SIMULATION SCENARIO

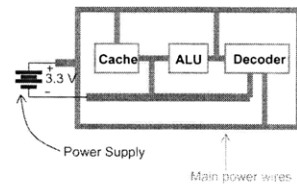


COURSE PHILOSOPHY

- Examine Several Modern Techniques
- Understand, practically and theoretically, how the techniques perform on representative, but real, applications
- Why Prove Theorems?
 - Guarantees, given assumptions, that the method will always work.
 - Can help debug programs.
 - The theorem proof can tell you what to do in practice.

EXAMPLE PROBLEMS

- Power Distribution on an Integrated Circuit**



- Is there at least 3V across the ALU?
- Design Objectives for the VLSI Problem
 - Select topology and metal widths and lengths so that

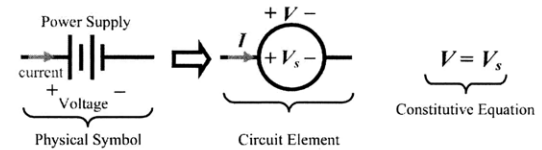
- a) Voltage across every function block > 3V
- b) Minimize the area used for the metal wires
- First Step – Analysis Tools
 - Given the topology and metal widths and lengths determine the voltage across the ALU, Cache and Decoder.
- Who uses VLSI Tools?
 - Several big companies
 - IBM, Motorola, TI, Intel, Compaq, Sony, Hitachi
 - Nonfunctional prototype costs:
 - Increases time to market
 - Design rework costs millions
 - 1000's of small companies
 - Small companies make application circuits disk drives, graphics accelerators, CD players, cell phones
 - What is the cost of nonfunctional prototypes?
 - Out of business

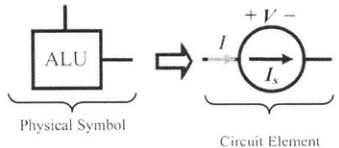
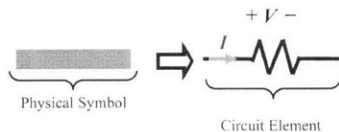
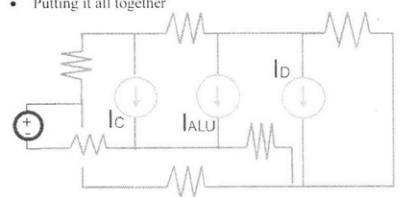
Who makes VLSI Tools?

Company	Employees	Sales	Market capital
Cadence	4,000	1.3 billion	3.8 billion
Synopsys/Avanti	5,000	1.5 billion	6.9 billion
Mentor Graphics	2,600	0.6 billion	1.4 billion

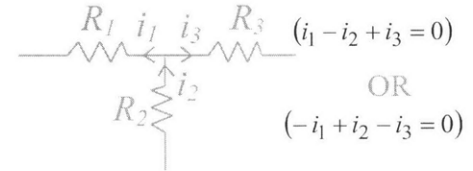
(Data from 2003)

- Companies compete by improving analysis efficiency.
- Modeling VLSI Circuit Power Distribution
 - Power Supply provide current at a certain voltage
 - Functional blocks draw current
 - The wire resistance generates losses
- Modeling the Circuit
 - Supply becomes a Voltage Source

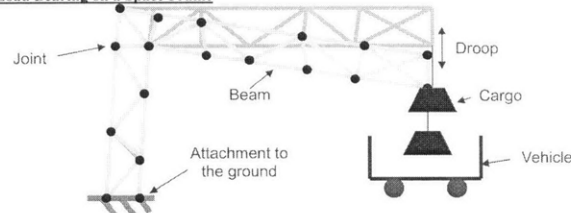


- Functional Blocks become Current Sources
 - 
 - Physical Symbol
 - Circuit Element
 - Constitutive Equation: $I = I_s$
- Metal Lines become Resistors
 - 
 - Physical Symbol
 - Circuit Element
 - Constitutive Equation (Ohm's Law): $IR - V = 0$
 - $R = \frac{\text{length}}{\text{area}} \cdot \text{resistivity}$
 - Design Parameter: length, area
 - Material Property: resistivity
- Putting it all together
 - 
 - Power Supply → Voltage Source
 - Functional Blocks → Current Sources
 - Wires → Resistors
 - Result is a schematic
- Formulating Equations from Schematics
 - Two Types of Unknowns
 - Node voltages, element currents
 - Constitutive Equations

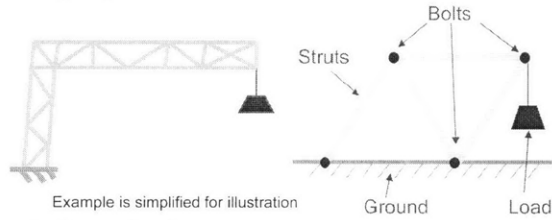
- Element current is related to voltage across the element
- Two Types of Equations
 - Sum of currents at each node = 0
 - Conservation/Balance Law Equations



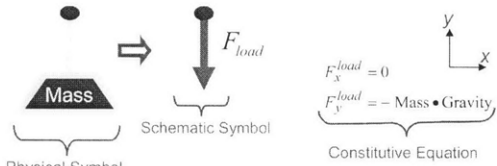
Note: here $i_1 = -3$ means it goes to the right i.e. opposite to the arrow

- **Load Bearing on a Space Frame**
 - 
 - Does the Space Frame droop too much under the load?
 - Design Objectives for the Space Frame
 - Select topology and strut widths and lengths so that
 - a) Droop is small enough
 - b) Minimize the metal used
 - First Step – Analysis Tools
 - Given the topology and metal widths and lengths determine the droop of the space frame under load

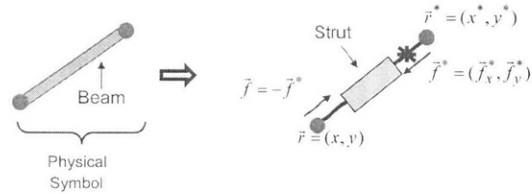
o Modeling the Space Frame



- Load becomes Force Source

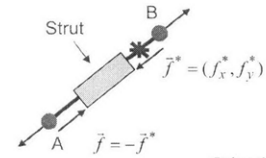


- Beam becomes Strut

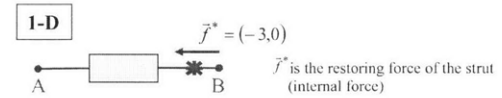
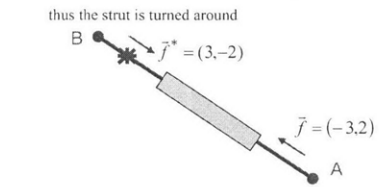
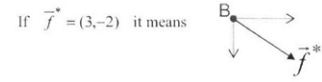


$l_0 = \text{Unstretched Length}$
 $A_c = \text{Cross-Sectional Area}$ } Design Parameters
 $E = \text{Young's Modulus}$ } Material Property
 $f^* = EA_c \frac{\Delta L}{L_0}$

Abstraction – leads to simplifications, such as:
 No bending/buckling
 No twisting
 No breaking



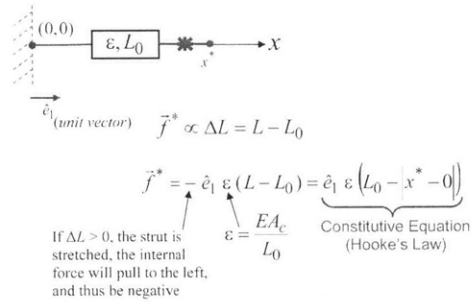
Only solve for the force, f^* , on one side
 The force on the other side is $f = -f^*$



Is this strut compressed or stretched?

The force of the strut is pulling the end in this means that the strut is being stretched.

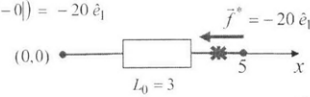
Think of what happens when you pull on the ends of a spring. Is the force from the spring pushing your hands away from one another or pulling them towards one another?



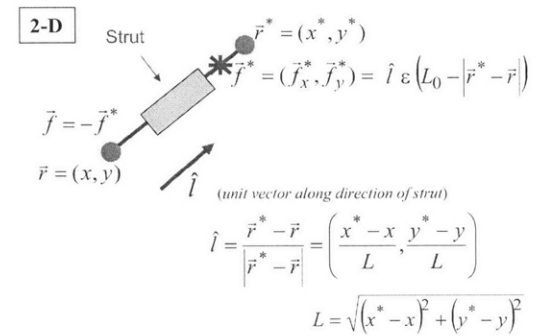
Example.

$\epsilon = 10, L_0 = 3, x^* = 5$

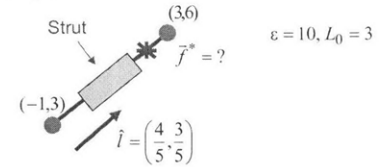
$\vec{f}^* = \hat{e}_1 10 (3 - |5 - 0|) = -20 \hat{e}_1$



In this case $\begin{cases} f_y = 0 \\ f_x = f^* \end{cases}$



Example.

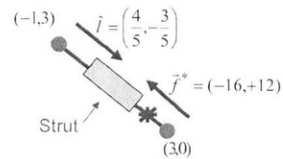


$L = |\vec{r}^* - \vec{r}| = \sqrt{[3 - (-1)]^2 + (6 - 3)^2} = 5$ stretched strut

$f_x^* = \frac{x^* - x}{L} \epsilon (L_0 - L) = \frac{3 - (-1)}{5} 10 (3 - 5) = -16$

$f_y^* = \frac{y^* - y}{L} \epsilon (L_0 - L) = \frac{6 - 3}{5} 10 (3 - 5) = -12$

Example. $\epsilon = 10, L_0 = 3$

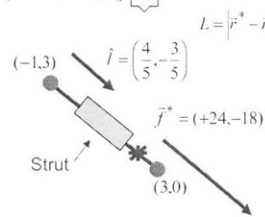


$$L = |\vec{r}^* - \vec{r}| = \sqrt{(3 - (-1))^2 + (0 - 3)^2} = 5 \quad \text{stretched strut}$$

$$f_x^* = \frac{x^* - x}{L} \epsilon (L_0 - L) = \frac{3 - (-1)}{5} 10 (3 - 5) = -16$$

$$f_y^* = \frac{y^* - y}{L} \epsilon (L_0 - L) = \frac{0 - 3}{5} 10 (3 - 5) = +12$$

Example. $\epsilon = 10, L_0 = 8$

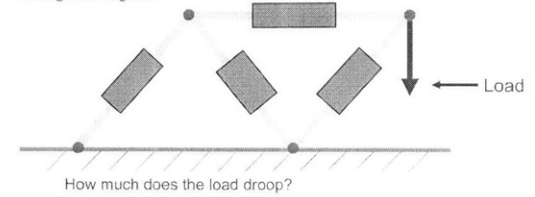


$$L = |\vec{r}^* - \vec{r}| = \sqrt{(3 - (-1))^2 + (0 - 3)^2} = 5 \quad \text{compressed strut}$$

$$f_x^* = \frac{x^* - x}{L} \epsilon (L_0 - L) = \frac{3 - (-1)}{5} 10 (8 - 5) = +24$$

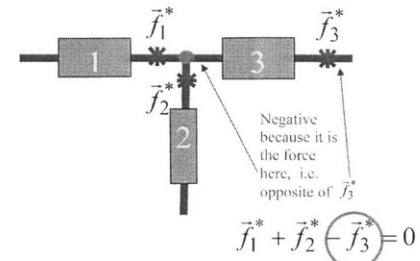
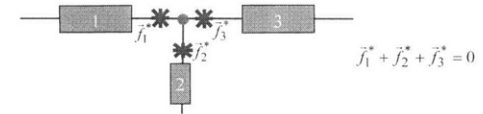
$$f_y^* = \frac{y^* - y}{L} \epsilon (L_0 - L) = \frac{0 - 3}{5} 10 (8 - 5) = -18$$

Putting It all Together



Formulating Equations from Schematics

- Two Types of Unknowns
 - Joint Positions, Strut Forces
- Constitutive Equations
 - Element Force is related to the change in Element Length
- Two Types of Equations
 - Sum of Forces at each joint = 0
- Conservation/Balance Law Equations

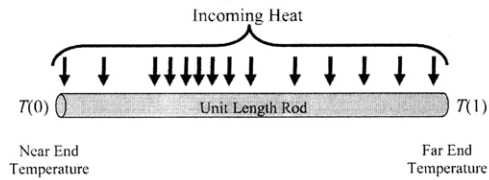
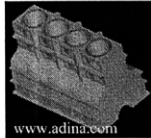


Note: $\vec{f}_1^* = (-3, 0)$ means the force goes to the left.

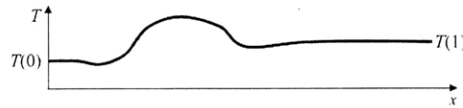
• **Temperature Distribution in an Engine**

Does the engine get too hot?

- Design Objectives for the Engine
 - Select the shape so that
 - a) The temperature does not get too high
 - b) Minimize the metal used
- Heat Flow 1-D Example
 - Conservation Laws and Constitutive Equations

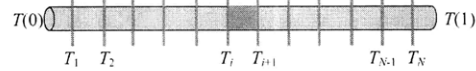


Question: What is the temperature distribution along the bar?



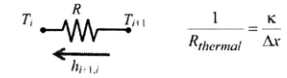
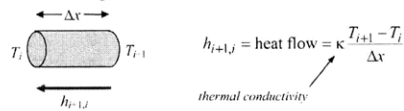
▪ Discrete Representation

1. Cut the bar into short sections
2. Assign each cut a temperature

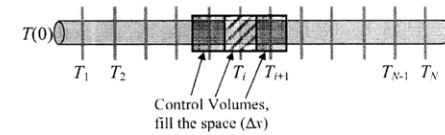
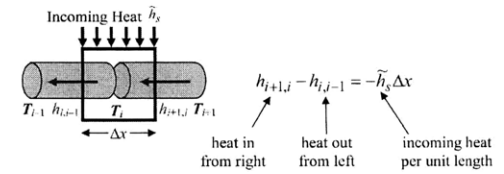


▪ Constitutive Relation

Heat Flow through one section



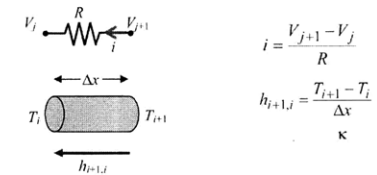
▪ Conservation Law
Net Heat Flow into Control Volume = 0



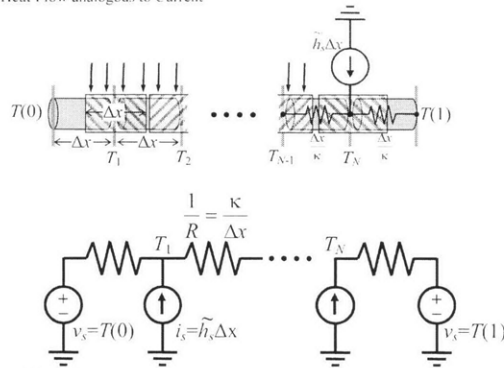
Limit as the sections become vanishingly small:

$$\lim_{\Delta x \rightarrow 0} \frac{h_{i+1,j} - h_{i,j-1}}{\Delta x} = -\tilde{h}_s$$

▪ Circuit Analogy

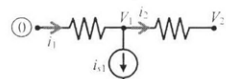


Temperature analogous to Voltage
Heat Flow analogous to Current

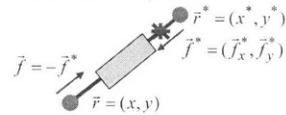


SUMMARY OF KEY POINTS:

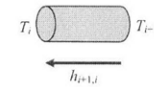
- Many Applications of simulation
 - Picked Three Representative Examples
Circuits, Struts and Joints, Heat Flow in Bar
- Two Types of Unknowns
 - Circuit – Node Voltages, Element Currents



- Struts – Joint Positions, Strut Forces



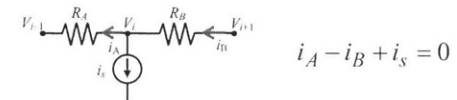
- Conducting Bar – Temperature, Section Heat Flows



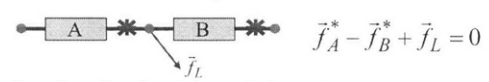
- Two Types of Equations

- Conservation/Balance Laws

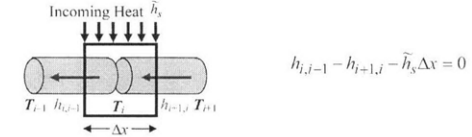
- Circuit – Sum of Currents at each node = 0



- Struts – Sum of Forces at each joint = 0



- Bar – Sum of heat flows into control volume = 0



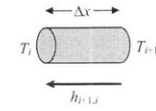
- Constitutive Equation
- Circuit – current-voltage relationship



- Struts – force-displacement relationship



- Bar – temperature drop-heat flow relationship



$$\dot{h}_{i+1,j} = \kappa \frac{T_{i+1} - T_i}{\Delta x}$$

INTRODUCTION TO NUMERICAL SIMULATION

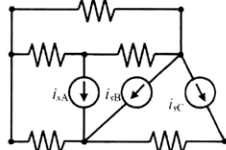
LECTURE 2. Equation Formulation & Node-Branch Stamping

TODAY'S OUTLINE:

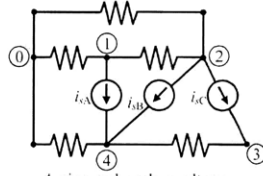
- Formulating Equations
 - Circuit Example
 - Struts and Joints Example
- Matrix Construction From Schematics
 - Node-Branch "Stamping Procedure"
 - Circuits
 - Struts and Joints

FORMULATING EQUATIONS FROM SCHEMATICS

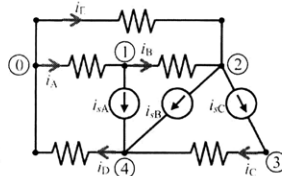
Circuit Example



• Step 1: Identifying Unknowns

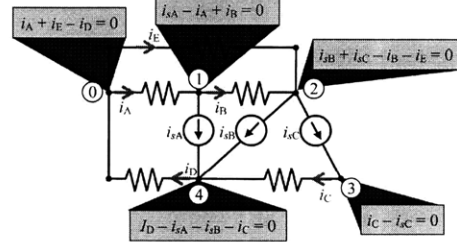


Assign each node a voltage, with one node as 0.

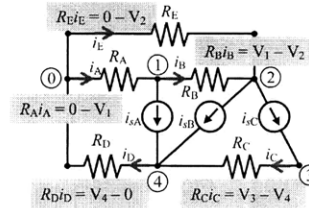


Assign each element a current.

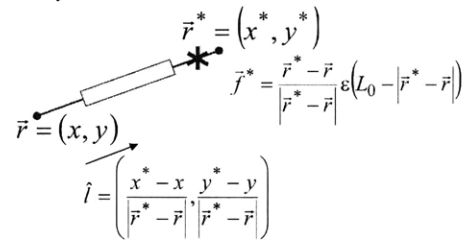
- Step 2: Conservation Laws
 - Sum of currents = 0 (Kirchhoff's Current Law)

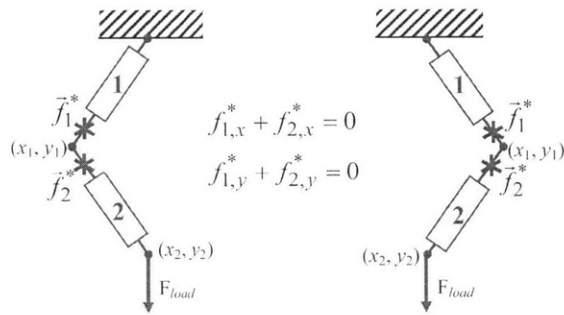


- Step 3: Constitutive Equations
 - Use Constitutive Equations to relate branch currents to node voltages (Currents flow from plus node to minus node)



Struts Example

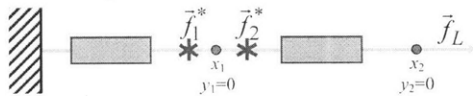




Will the solutions be the same?
 Will the set of conservation law equations be different?

Conservation laws for the two examples will be exactly the same. The perceived force "direction" is inconsequential, it is the adjacent forces that matter.

Two struts aligned with the x axis



Conservation Law

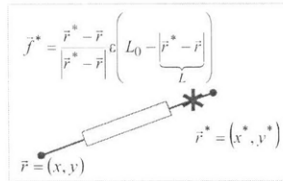
At node 1: $f_{1,x}^* + f_{2,x}^* = 0$

At node 2: $-f_{2,x}^* + f_L = 0$

Constitutive Equations

$$f_{1,x}^* = \frac{x_1 - 0}{|x_1 - 0|} \epsilon(L_0 - |x_1 - 0|)$$

$$f_{2,x}^* = \frac{x_1 - x_2}{|x_1 - x_2|} \epsilon(L_0 - |x_1 - x_2|)$$



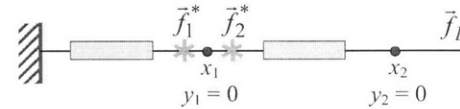
Reduced (Nodal) Equations

$$f_{1,x}^* + f_{2,x}^* = 0$$

$$\frac{x_1}{|x_1|} \epsilon(L_0 - |x_1|) + \frac{x_1 - x_2}{|x_1 - x_2|} \epsilon(L_0 - |x_1 - x_2|) = 0$$

$$-f_{2,x}^* + f_L = 0$$

$$-\frac{x_1 - x_2}{|x_1 - x_2|} \epsilon(L_0 - |x_1 - x_2|) + f_L = 0$$



Example.) $f_L = 10\epsilon_1$ (force in positive x direction)

Solution of Nodal Equations

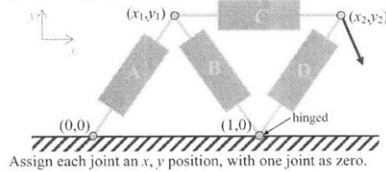
$$x_1 = L_0 + \frac{10}{\epsilon} \quad x_2 = x_1 + L_0 + \frac{10}{\epsilon}$$

Notice the signs of the forces

$$f_{2,x}^* = 10 \quad \text{force in positive x direction}$$

$$f_{1,x}^* = -10 \quad \text{force in negative x direction}$$

• Step 1: Identifying Unknowns



Assign each strut an x and y force component.

• Step 2: Conservation Laws

$$\begin{aligned} f_{Ax}^* + f_{Bx}^* + f_{Cx}^* &= 0 \\ f_{Ay}^* + f_{By}^* + f_{Cy}^* &= 0 \end{aligned}$$

$$\begin{aligned} -f_{Cx}^* + f_{Dx}^* + f_{load,x} &= 0 \\ -f_{Cy}^* + f_{Dy}^* + f_{load,y} &= 0 \end{aligned}$$

Force Equilibrium
Sum of x-directed forces at a joint = 0
Sum of y-directed forces at a joint = 0

• Step 3: Constitutive Equations

$$f_{Ax}^* = \frac{x_1 - 0}{L_A} e(L_{A,0} - L_A)$$

$$f_{Ay}^* = \frac{y_1 - 0}{L_A} e(L_{A,0} - L_A)$$

$$L_A = \sqrt{(x_1 - 0)^2 + (y_1 - 0)^2}$$

$$f_{Bx}^* = \frac{x_1 - 0}{L_B} e(L_{B,0} - L_B)$$

$$f_{By}^* = \frac{y_1 - 0}{L_B} e(L_{B,0} - L_B)$$

$$f_{Cx}^* = \frac{x_2 - x_1}{L_C} e(L_{C,0} - L_C)$$

$$f_{Cy}^* = \frac{y_2 - y_1}{L_C} e(L_{C,0} - L_C)$$

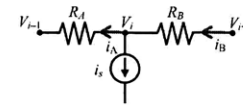
$$f_{Dx}^* = \frac{x_2 - 1}{L_D} e(L_{D,0} - L_D)$$

$$f_{Dy}^* = \frac{y_2 - 0}{L_D} e(L_{D,0} - L_D)$$

Use Constitutive Equations to relate strut forces to joint positions.

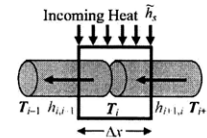
Comparing Conservation Laws

- Circuit



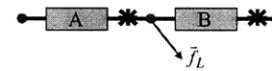
$$i_A - i_B + i_s = 0$$

- Heat Flow



$$h_{i,i-1} - h_{i+1,i} - \tilde{h}_s \Delta x = 0$$

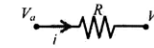
- Struts



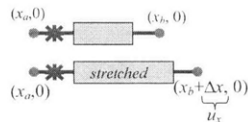
$$\vec{f}_A - \vec{f}_B + \vec{f}_L = 0$$

Summary of Key Points

- Two Types of Unknowns
 - Circuit: Node Voltages, Element Currents
 - Struts: Joint Positions, Strut Forces
 - Bar: Node Temperatures, Heat Flows
- Two Types of Equations
 - Conservation/Balance Laws
 - Circuit: Sum of Currents at each node = 0
 - Struts: Sum of Forces at each joint = 0
 - Bar: Sum of heat flows into control volume = 0
 - Constitutive Equation
 - Circuit: current-voltage relationship
 - Struts: force-displacement relationship



$$i = \frac{1}{R}(V_a - V_b)$$



$$f_x^* = \frac{x_a - (x_b + \Delta x)}{x_a - (x_b + \Delta x_s)} \varepsilon [x_a - x_b - |x_a - (x_b + \Delta x_s)|]$$

f_x^* is a positive number – force is to the right, this makes sense for a stretched strut.

- Bar: temperature drop-heat flow relationship



GENERATING MATRICES FROM SCHEMATICS

Number of Columns = Number of Unknowns

Number of Rows = Number of Equations

$$\begin{bmatrix} A_{1,j} \\ \vdots \\ A_{i,j} \\ \vdots \\ A_{n,j} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$

row column

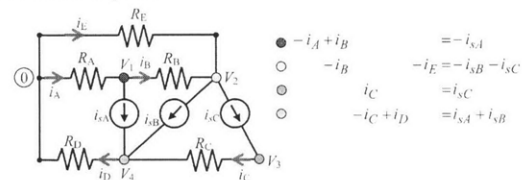
Circuit Example

Assume linear constitutive equations...

$$\begin{matrix} \text{N KCL Eqns} \\ \text{B Constitutive Eqns} \end{matrix} \begin{bmatrix} \text{B} \\ \text{branch} \\ \text{currents} \end{bmatrix} \begin{bmatrix} \text{N} \\ \text{node} \\ \text{voltages} \end{bmatrix} \begin{bmatrix} i_B \\ V_N \end{bmatrix} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

- One matrix column for each unknown
 - N columns for the Node voltage
 - B columns for the Branch currents
- One matrix row for each equation
 - N rows for KCL
 - B rows for element constitutive equations (linear and square system!)

o Conservation Equation



Matrix Form for the Equations

The matrix A is usually not square

$$\begin{matrix} \text{one row} \\ \text{for each} \\ \text{KCL} \\ \text{equation} \end{matrix} \begin{matrix} \text{N} \\ \downarrow \\ \text{1} \\ \text{2} \\ \text{3} \\ \text{4} \end{matrix} \begin{bmatrix} -1 & 1 & & & \\ & -1 & & & \\ & & 1 & & \\ & & & -1 & 1 \\ & & & & & \end{bmatrix} \begin{bmatrix} i_A \\ i_B \\ i_C \\ i_D \\ i_E \end{bmatrix} = \begin{bmatrix} -i_{SA} \\ -i_{SB} - i_{SC} \\ i_{SC} \\ i_{SA} + i_{SB} \end{bmatrix}$$

one column for each branch current

right hand side for source currents

Four Nodes

Do we know how many rows the A matrix will have?
 3 because there are 3 nonzero nodes

Do we know how many columns the A matrix will have?
 Could have any number, depends on the number of elements

Input file of a circuit simulation has one line per resistor:

Rname	node1	node2	value
RnameA	n1	n2	10
RnameB	n1	n3	5
Iname	n2	n3	5

Spice → circuit simulator code...

Node 1 • KCL • V_1
 Node 2 • KCL • V_2
 Node 3 • KCL • V_3
 ⋮

	R_A	R_B	I
Node 1	+1	+1	
Node 2	-1		+1
Node 3		-1	
⋮			

How does each resistor contribute to the matrix?

A has no more than two non-zeros per column

KCL at...

$$n_1: \sum i_{other} + i_k = \sum i_s$$

$$n_2: \sum i_{other} - i_k = \sum i_s$$

	k
n_1	1
n_2	-1

A

How does each current source contribute to the Conservation Law Equation?

Affects the Right Hand Side

KCL at...

$$n_1: \sum i_b^s = \sum i_{s_{other}} - i_{sb}$$

$$n_2: \sum i_b^s = \sum i_{s_{other}} + i_{sb}$$

	RHS
n_1	$\sum i_{s_{other}} - i_{sb}$
n_2	$\sum i_{s_{other}} + i_{sb}$

Conservation Matrix Equation Generation Algorithm

For each resistor

if ($n_1 \neq 0$) then $A(n_1, k) \leftarrow 1$

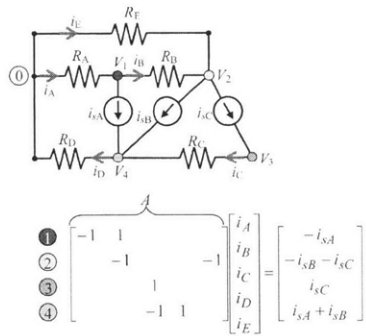
if ($n_2 \neq 0$) then $A(n_2, k) \leftarrow -1$

Set $I_s =$ zero vector

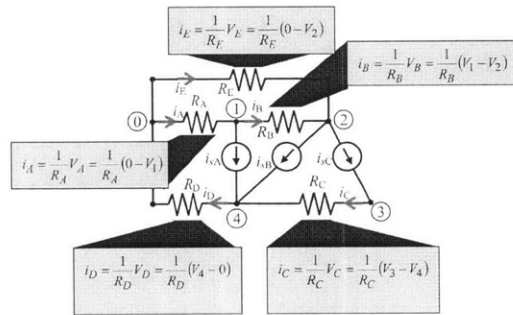
For each current source

if ($n_1 \neq 0$) then $I_s(n_1) \leftarrow I_s(n_1) - i_{sb}^n$

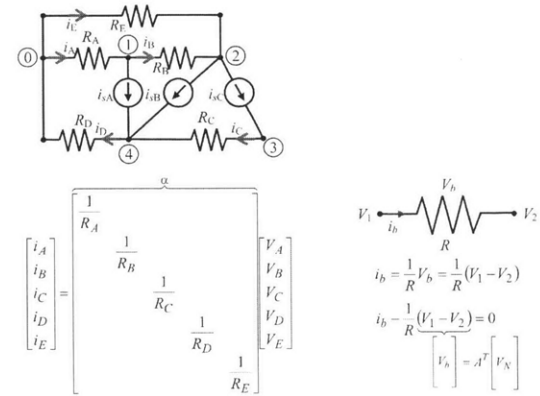
if ($n_2 \neq 0$) then $I_s(n_2) \leftarrow I_s(n_2) + i_{sb}$



o Constitutive Equation



First relate branch currents to branch voltages
Second determine voltages across resistors (Branch Voltages)



The k^{th} resistor contributes $\frac{1}{R_k}$ to α (k,k)

The matrix α relates branch voltages to branch currents.

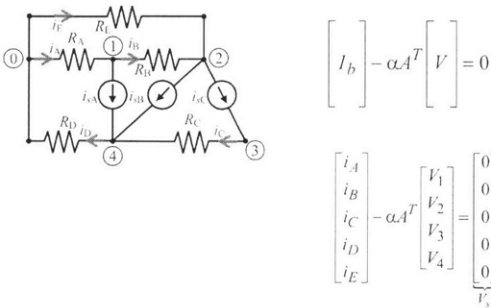
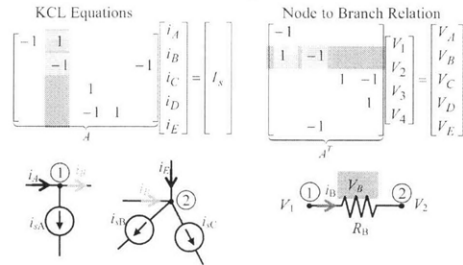
- One row for each unknown current.
- One column for each associated branch voltage.

The matrix α is square and diagonal.

Relationship between branch voltages and node voltages:

$$\begin{matrix} V_A = 0 - V_1 \\ V_B = V_1 - V_2 \\ V_C = V_3 - V_4 \\ V_D = V_4 - 0 \\ V_E = 0 - V_2 \end{matrix} \quad \begin{matrix} \text{Examine} \\ \text{Matrix} \\ \text{Construction} \end{matrix} \quad \begin{bmatrix} V_A \\ V_B \\ V_C \\ V_D \\ V_E \end{bmatrix} = \begin{bmatrix} -1 & & & & \\ 1 & -1 & & & \\ & & 1 & -1 & \\ & & & 1 & \\ & & & & -1 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \end{bmatrix}$$

The node-to-Branch matrix is the transpose of the KCL Matrix.



The node voltages can be related to branch currents
 - A^T relates node voltages to branch voltages.
 - α relates branch voltages to branch currents.
 α is square and diagonal

- Node-Branch Form
 $I_b - \alpha A^T V_N = 0$ Constitutive Relation
 $A I_b = I_s$ Conservation Law

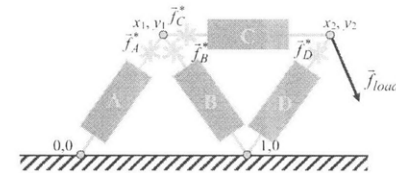
$$B \Downarrow \begin{bmatrix} I & -\alpha A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} I_b \\ V_N \end{bmatrix} = \begin{bmatrix} 0 \\ I_s \end{bmatrix}$$

$\begin{matrix} \leftarrow & \rightarrow \\ B & N \end{matrix}$

N = number of Nodes with unknown voltages
 B = number of Branches with unknown currents

Struts Example

- In 2-D
 - One pair of columns for each unknown
 - J pairs of columns for the joint positions
 - S pairs of columns for the strut forces
 - One pair of matrix rows for each equation
 - J pairs of rows for the force equilibrium equations
 - S pairs of rows for the **linearized** constitutive relations
- Follow Approach Parallel to Circuits
 - Form an "Incidence Matrix," A , from Conservation Law.
 - Determine strut deformation using A^T .
 - Use linearized constitutive equations to relate strut deformation.
 - Combine (1), (2), and (3) to generate a node-branch form.
- Conservation Laws

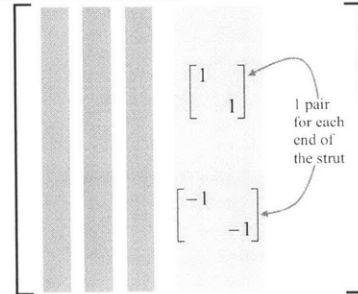


$$\begin{aligned} \vec{f}_{A,x} + \vec{f}_{B,x} + \vec{f}_{C,x} + \vec{f}_{D,x} &= 0 \\ \vec{f}_{A,y} + \vec{f}_{B,y} - \vec{f}_{C,y} + \vec{f}_{D,y} &= 0 \\ -\vec{f}_{C,x} + \vec{f}_{D,x} + \vec{f}_{D,y} - \vec{f}_{C,y} &= -\vec{f}_{load,x} \\ -\vec{f}_{C,y} + \vec{f}_{D,y} &= -\vec{f}_{load,y} \end{aligned}$$

Note that struts A & B only contribute one pair of entries into the A matrix and strut C contributes two pairs of entries into the A matrix. This is because struts A & B are connected to the wall (ground) and strut C has two free ends.

Stamping Approach

Matrix – load 1 column at a time



Load pair of columns per strut
Load right side for load

$$\begin{matrix} x_1 \\ y_1 \\ x_2 \\ y_2 \end{matrix} \underbrace{\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ & & & 1 \\ & & & -1 \end{bmatrix}}_A \begin{matrix} f_{A,x}^* \\ f_{A,y}^* \\ f_{B,x}^* \\ f_{B,y}^* \\ f_{C,x}^* \\ f_{C,y}^* \\ f_{D,x}^* \\ f_{D,y}^* \end{matrix} = \begin{matrix} 0 \\ 0 \\ -f_{load,x} \\ -f_{load,y} \end{matrix}$$

Conservation Matrix Generation Algorithm

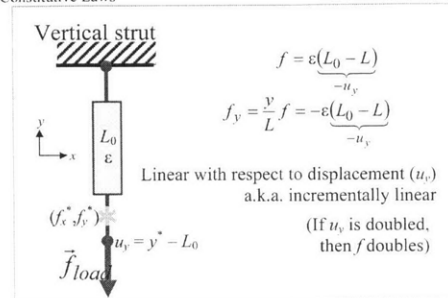
For each strut
 If (j_1 is not fixed) $A(j_{1x}, b_x) = 1$ $A(j_{1y}, b_y) = 1$
 If (j_2 is not fixed) $A(j_{2x}, b_x) = -1$ $A(j_{2y}, b_y) = -1$

For each load
 If (j_1 is not fixed) $F_I(j_{1x}) = F_I(j_{1x}) - f_{load,x}$
 $F_I(j_{1y}) = F_I(j_{1y}) - f_{load,y}$

A has at most 2 non-zeros per column



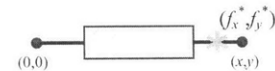
Constitutive Laws



Linearize the Constitutive Equations.

$$f_x = \frac{x}{L} EA_c \frac{L_0 - L}{L_0} = \frac{x}{L} \epsilon(L_0 - L)$$

$$f_y = \frac{y}{L} EA_c \frac{L_0 - L}{L_0} = \frac{y}{L} \epsilon(L_0 - L)$$



Determine the first derivatives of the non-linear constitutive equations as a first step to linearize the equations.

$$\frac{\partial f_x}{\partial x} = \frac{\partial}{\partial x} \left[\epsilon x \left(\frac{L_0}{L} - 1 \right) \right] = \frac{\partial}{\partial x} \left[\epsilon x \left(\frac{L_0}{\sqrt{x^2 + y^2}} - 1 \right) \right] = \epsilon \left(\frac{L_0 y^2}{(x^2 + y^2)^{3/2}} - 1 \right)$$

$$\frac{\partial f_x}{\partial y} = \frac{\partial}{\partial y} \left[\epsilon x \left(\frac{L_0}{\sqrt{x^2 + y^2}} - 1 \right) \right] = \epsilon \frac{-L_0 x y}{(x^2 + y^2)^{3/2}}$$

$$\frac{\partial f_y}{\partial x} = \frac{\partial}{\partial x} \left[\epsilon y \left(\frac{L_0}{\sqrt{x^2 + y^2}} - 1 \right) \right] = \epsilon \frac{-L_0 x y}{(x^2 + y^2)^{3/2}}$$

$$\frac{\partial f_y}{\partial y} = \frac{\partial}{\partial y} \left[\epsilon y \left(\frac{L_0}{\sqrt{x^2 + y^2}} - 1 \right) \right] = \epsilon \left(\frac{L_0 x^2}{(x^2 + y^2)^{3/2}} - 1 \right)$$

Evaluate these first derivatives about the point (x_0, y_0) .

Linearization becomes

$$f_x \approx \frac{\partial f_x}{\partial x} \Big|_{(x_0, y_0)} (x - x_0) + \frac{\partial f_x}{\partial y} \Big|_{(x_0, y_0)} (y - y_0)$$

$$f_y \approx \frac{\partial f_y}{\partial x} \Big|_{(x_0, y_0)} (x - x_0) + \frac{\partial f_y}{\partial y} \Big|_{(x_0, y_0)} (y - y_0)$$

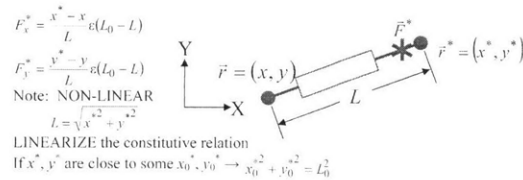
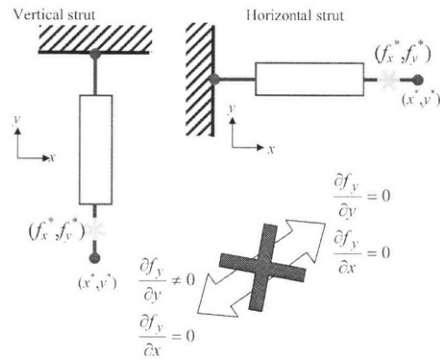
Linearization is only a valid estimate if the strut is very close to its original configuration - x and y are very close to x_0 and y_0 .

What if the strut is *rotated*, is this linearization still valid?

It seems like it should be because the strut is not stretched (only rotated), but it isn't. Because the non-linearity is *not* in the relationship between how much the strut is stretched and the force going through the strut but the non-linearity is in the projection of that force onto the x and y axis.

Referred to as the "geometric" nonlinearity in finite element literature.

Example.



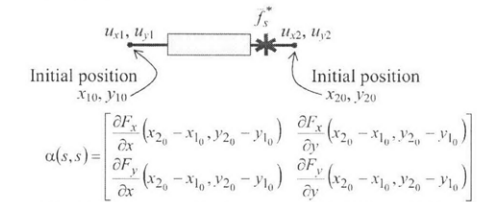
$$\begin{bmatrix} f_x^* \\ f_y^* \end{bmatrix} = \begin{bmatrix} \frac{\partial F_x^*}{\partial x_1} (x_0^*, y_0^*) & \frac{\partial F_x^*}{\partial y_1} (x_0^*, y_0^*) \\ \frac{\partial F_y^*}{\partial x_1} (x_0^*, y_0^*) & \frac{\partial F_y^*}{\partial y_1} (x_0^*, y_0^*) \end{bmatrix} \begin{bmatrix} u_x^* \\ u_y^* \end{bmatrix} + \Delta \vec{f}$$

$$\begin{aligned} \vec{u}^* &= (x^* - x_0^*, y^* - y_0^*) \\ \vec{u} &= (x - x_0, y - y_0) \\ \vec{u}_b &= \vec{u}^* - \vec{u} \end{aligned}$$

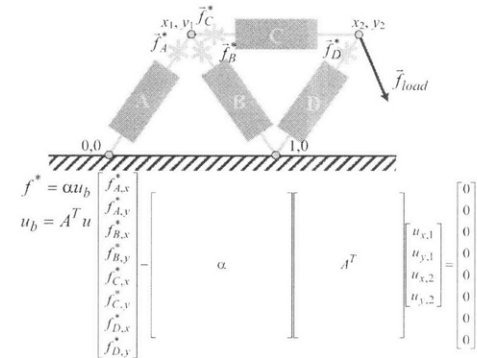
Jacobian Matrix

We will learn more about linearizing equations later on when we study Newton's Method.

The $\alpha(s,s)$ block



Note that the arguments in the α matrix are the difference in the original position between node 1 and node 2. This is because the force of the strut is dependent upon the relative position of the two joints on either side.



- o Node-Branch Form
 - $f_s - \alpha A^T u = 0$ Constitutive Equation
 - $A f_s = f_L$ Conservation Law

$$\begin{array}{c}
 2 \cdot S \downarrow \\
 2 \cdot J \downarrow \\
 \begin{array}{c} \leftarrow \\ \leftarrow \\ \leftarrow \end{array} \\
 \begin{array}{c} \leftarrow \\ \leftarrow \\ \leftarrow \end{array}
 \end{array}
 \begin{bmatrix} I & -\alpha A^T \\ A & 0 \end{bmatrix}
 \begin{bmatrix} f_s \\ u \end{bmatrix}
 =
 \begin{bmatrix} 0 \\ f_L \end{bmatrix}$$

S = Number of Struts
 J = Number of Unfixed Joints

Comparison

$$\begin{array}{c}
 \text{struts} \\
 2 \cdot S \downarrow \\
 2 \cdot J \downarrow \\
 \begin{array}{c} \leftarrow \\ \leftarrow \\ \leftarrow \end{array} \\
 \begin{array}{c} \leftarrow \\ \leftarrow \\ \leftarrow \end{array}
 \end{array}
 \begin{bmatrix} I & -\alpha A^T \\ A & 0 \end{bmatrix}
 \begin{bmatrix} f_s \\ u \end{bmatrix}
 =
 \begin{bmatrix} 0 \\ f_L \end{bmatrix}$$

$$\begin{array}{c}
 \text{circuit} \\
 B \downarrow \\
 N \downarrow \\
 \begin{array}{c} \leftarrow \\ \leftarrow \\ \leftarrow \end{array} \\
 \begin{array}{c} \leftarrow \\ \leftarrow \\ \leftarrow \end{array}
 \end{array}
 \begin{bmatrix} I & -\alpha A^T \\ A & 0 \end{bmatrix}
 \begin{bmatrix} I_b \\ V_N \end{bmatrix}
 =
 \begin{bmatrix} 0 \\ I_s \end{bmatrix}$$

Summary of Key points...

- o Developed algorithms for automatically constructing matrix equations from schematics using
 - Conservation law
 - Constitutive equations
- o Looked at one formulation: node-branch
- o Next time: nodal formulation

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 3. Equation Formulation – Nodal Analysis

TODAY'S OUTLINE:

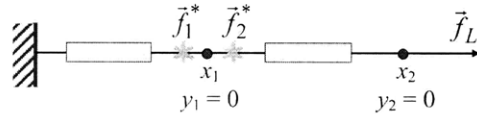
- Matrix Construction from Schematics
 - Nodal "Stamping Procedure"
 - Struts and Joints
 - Circuits
 - Comparing Node-Branch vs. Nodal
- Solution of Linear Systems
 - Existence and Uniqueness

MATRIX CONSTRUCTION FROM SCHEMATICS

Nodal "Stamping Procedure"

❖ Struts and Joints

Two struts aligned with the X axis



Conservation Law

At node 1: $f_{1,x}^* + f_{2,x}^* = 0$

At node 2: $-f_{2,x}^* + f_L = 0$

Constitutive Equations

$$f_{1,x}^* = \frac{x_1 - 0}{x_1 - 0} \epsilon(L_0 - x_1 - 0)$$

$$f_{2,x}^* = \frac{x_1 - x_2}{x_1 - x_2} \epsilon(L_0 - |x_1 - x_2|)$$

Reduced (Nodal) Equations

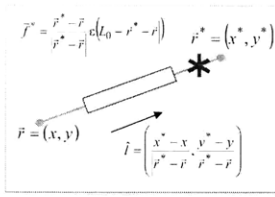
$$f_{1,x}^* + f_{2,x}^* = 0$$

$$\frac{x_1}{x_1} \epsilon(L_0 - |x_1|) + \frac{x_1 - x_2}{x_1 - x_2} \epsilon(L_0 - |x_1 - x_2|) = 0$$

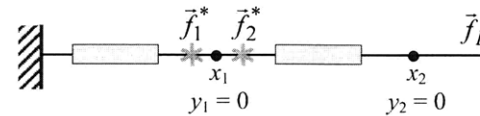
$$\underbrace{f_{1,x}^*}_{f_{1,x}^*} + \underbrace{f_{2,x}^*}_{f_{2,x}^*} = 0$$

$$-f_{2,x}^* + f_L = 0$$

$$\underbrace{-\frac{x_1 - x_2}{x_1 - x_2} \epsilon(L_0 - |x_1 - x_2|)}_{-f_{2,x}^*} + f_L = 0$$



Two struts aligned with the X axis



Example.) $\vec{f}_L = 10\hat{e}_1$ (force in positive x direction)

Solution of Nodal Equations

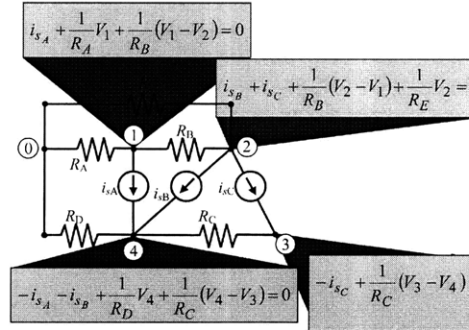
$$x_1 = L_0 + \frac{10}{\epsilon} \quad x_2 = x_1 + L_0 + \frac{10}{\epsilon}$$

Notice the signs of the forces:

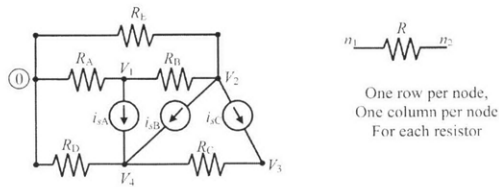
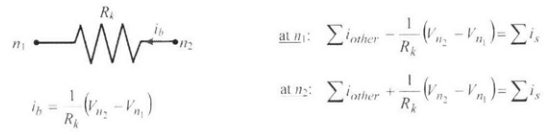
$$f_{2,x}^* = 10 \quad \text{force in positive } x \text{ direction}$$

$$f_{1,x}^* = -10 \quad \text{force in negative } x \text{ direction}$$

❖ Circuits



- (1) Number the nodes with one node as 0.
- (2) Write a conservation law at each node except (0) in terms of the node voltages!



KCL equations

$$\begin{bmatrix} \frac{1}{R_A} + \frac{1}{R_B} & -\frac{1}{R_B} \\ -\frac{1}{R_B} & \frac{1}{R_B} + \frac{1}{R_E} \\ & & \frac{1}{R_C} & -\frac{1}{R_C} \\ -\frac{1}{R_C} & \frac{1}{R_C} + \frac{1}{R_D} \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \end{bmatrix} = \begin{bmatrix} -i_{sA} \\ -i_{sB} - i_{sC} \\ i_{sC} \\ i_{sA} + i_{sB} \end{bmatrix}$$

G

Notice that the contributions are positive on the diagonal and negative on the off-diagonal.
 G is square.

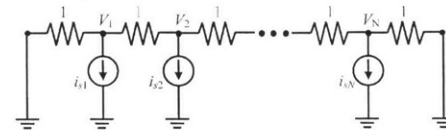
Node-Branch Matrix

$$\begin{matrix} N+B \\ \downarrow \\ \begin{bmatrix} I & -\alpha A^T \\ A & 0 \end{bmatrix} \\ \leftarrow N+B \end{matrix}$$

Nodal Matrix

$$\begin{matrix} N \\ \downarrow \\ [G] \\ \leftarrow N \end{matrix}$$

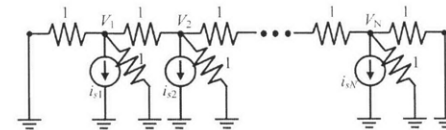
Circuit Example.)



What does the nodal form for the equations look like?

$$\begin{bmatrix} 1+1 & -1 \\ -1 & 1+1 & -1 \\ & -1 & 1+1 & -1 \\ & & \ddots & \ddots \\ & & & -1 & 1+1 & -1 \\ & & & & -1 & 1+1 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \\ \vdots \\ V_N \end{bmatrix} = \begin{bmatrix} i_{s1} \\ i_{s2} \\ \vdots \\ i_{sN} \end{bmatrix}$$

Diagonally dominant



What does the nodal form for the equations look like?

$$\begin{bmatrix} 1+1+1 & -1 \\ -1 & 1+1+1 & -1 \\ & -1 & 1+1+1 & -1 \\ & & \ddots & \ddots \\ & & & -1 & 1+1+1 & -1 \\ & & & & -1 & 1+1+1 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \\ \vdots \\ V_N \end{bmatrix} = \begin{bmatrix} i_{s1} \\ i_{s2} \\ \vdots \\ i_{sN} \end{bmatrix}$$

Strictly Diagonally Dominant

Nodal Matrix Generation Algorithm

if $(n_1 > 0)$ & $(n_2 > 0)$

$$G(n_1, n_1) = G(n_1, n_1) + \frac{1}{R} \quad G(n_2, n_1) = G(n_2, n_1) - \frac{1}{R}$$

$$G(n_1, n_2) = G(n_1, n_2) - \frac{1}{R} \quad G(n_2, n_2) = G(n_2, n_2) + \frac{1}{R}$$

else if $(n_1 > 0)$

$$G(n_1, n_1) = G(n_1, n_1) + \frac{1}{R}$$

else

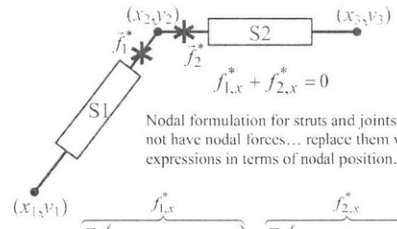
$$G(n_2, n_2) = G(n_2, n_2) + \frac{1}{R}$$

Resistor Networks

$$N \times N \quad \mathbf{G} \quad \mathbf{V}_n = \mathbf{I}_s$$

Struts and Joints

$$2 \times J \times 2 \times J \quad \mathbf{G} \quad \mathbf{u}_j = \mathbf{F}_l$$

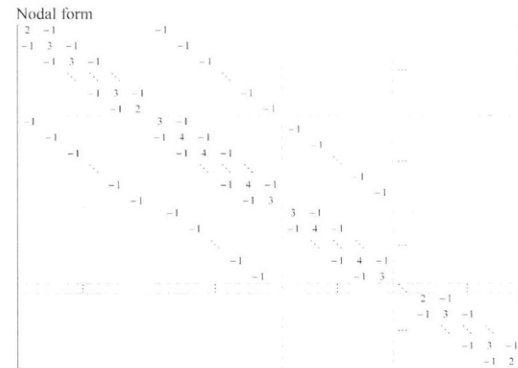
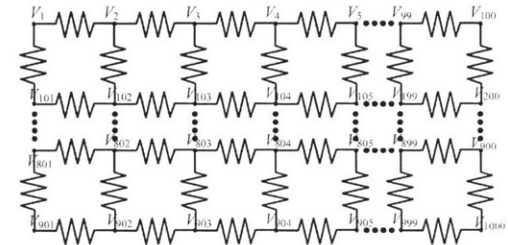


Nodal formulation for struts and joints will not have nodal forces... replace them with expressions in terms of nodal position.

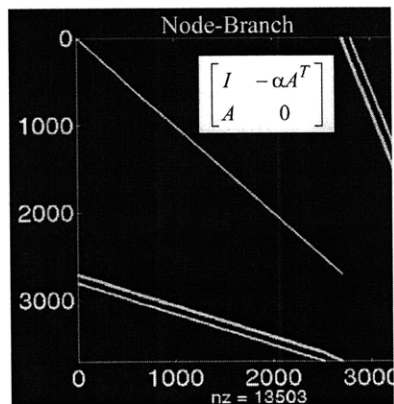
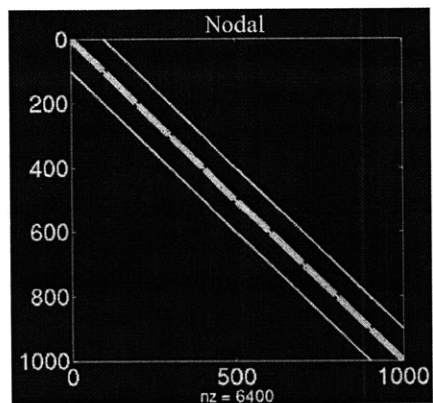
$$F_x(x_2 - x_1, y_2 - y_1) + F_x(x_2 - x_3, y_2 - y_3) = 0$$

Comparing Node-Branch vs. Nodal

- ❖ Comparing Matrix Sparsity
- Example.) Resistor Grid



Matrix non-zero locations for 100 × 10 Resistor Grid.



Node-Branch Matrix

$$\begin{matrix} \text{Constitutive} \\ \text{Conservation Law} \end{matrix} \begin{bmatrix} I & -\alpha A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} I_b \\ V_N \end{bmatrix} = \begin{bmatrix} 0 \\ I_s \end{bmatrix}$$

Nodal Matrix

$$\text{Conservation Law } [G][V_N] = [I_s]$$

(constitutive pre-substituted)

- ❖ G matrix properties
 - Smaller $N \times N \ll (N+B) \times (N+B)$
 - Symmetric $G_{ij} = G_{ji}$
 - Diagonally Dominant $|G_{ii}| \geq \sum_{j \neq i} |G_{ij}|$

❖ Node-Branch Form

$$\begin{bmatrix} I & -\alpha A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} I_b \\ V_n \end{bmatrix} = \begin{bmatrix} 0 \\ I_s \end{bmatrix}$$

$\underbrace{\hspace{1.5cm}}_M \quad \underbrace{\hspace{1.5cm}}_x \quad \underbrace{\hspace{1.5cm}}_b$

- Not Symmetric, nor Diagonally Dominant
- Matrix is $(n+b) \times (n+b)$

❖ Deriving Formulation from Node-Branch constitutive equations

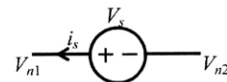
$$I_b - \alpha A^T V_N = 0 \Rightarrow I_b = \alpha A^T V_N$$

substitute into conservation equations:

$$A I_b = I_s \Rightarrow \underbrace{A \alpha A^T}_G V_N = I_s$$

❖ Problem Element

- Voltage Source
- Constitutive Equation $0 \cdot i_s + V_{n1} - V_{n2} = V_s$



SOLUTION OF LINEAR SYSTEMS

$$M\vec{x} = \vec{b}$$

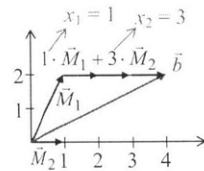
$$\begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ \vec{M}_1 & \vec{M}_2 & \dots & \vec{M}_N \\ \downarrow & \downarrow & \dots & \downarrow \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_N \end{bmatrix} \Rightarrow x_1\vec{M}_1 + x_2\vec{M}_2 + \dots + x_N\vec{M}_N = \vec{b}$$

Find a set of weights, x , so that the weighted sum of the columns of the matrix M is equal to the right hand side b .

$$\overbrace{\begin{bmatrix} 1 \\ 2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}}^M \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 4 \\ 2 \end{bmatrix}$$

$$x_1\vec{M}_1 + x_2\vec{M}_2 = \vec{b}$$

$$x_1 \begin{bmatrix} 1 \\ 2 \end{bmatrix} + x_2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 4 \\ 2 \end{bmatrix}$$



Key Questions

- ❖ Given $Mx = b$
 - Is there a solution?
 - Is the solution unique?

If $M = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix}$ Can $b = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$?

No – so change M so that b can be a solution $\Leftrightarrow M = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \end{bmatrix}$

What if $M = \begin{bmatrix} 1 & 1 & 2 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \end{bmatrix}$ Note that the third column is the sum of the first two columns in this case. The columns of M are not linearly independent.

Can $b = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$? $x = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$

Can $b = \begin{bmatrix} 4 \\ 3 \\ 1 \end{bmatrix}$? $x = \begin{bmatrix} 1 \\ 3 \\ 0 \end{bmatrix}$ or $x = \begin{bmatrix} 0 \\ 2 \\ 1 \end{bmatrix}$

$My = 0 \Leftrightarrow y = \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix}$

Existence and Uniqueness

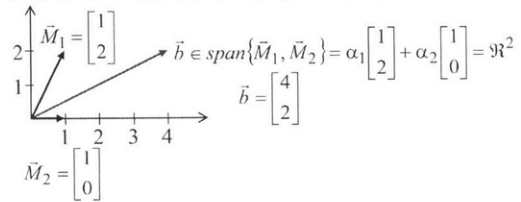
❖ Is there a solution?

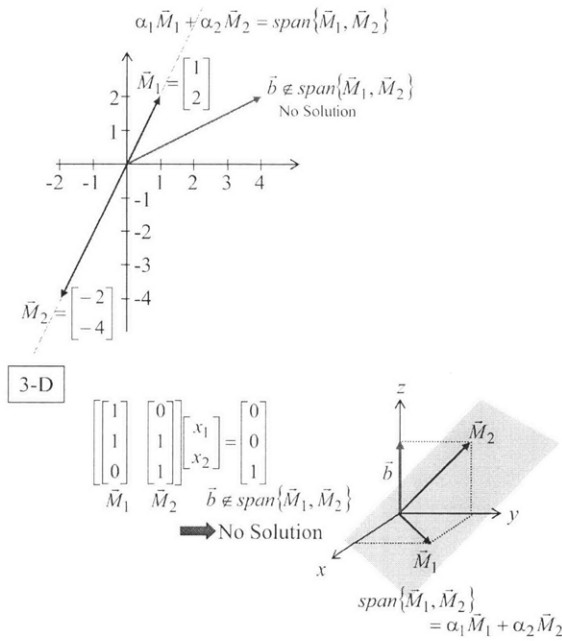
There exist weights x_1, \dots, x_N , such that

$$x_1\vec{M}_1 + x_2\vec{M}_2 + \dots + x_N\vec{M}_N = \vec{b}$$

A solution exists when b is in the span of the columns of M

$$\text{span}\{\vec{M}_1, \vec{M}_2, \dots, \vec{M}_N\} = \{\alpha_1\vec{M}_1 + \alpha_2\vec{M}_2 + \dots + \alpha_N\vec{M}_N, \alpha_i \in \mathbb{R}\}$$

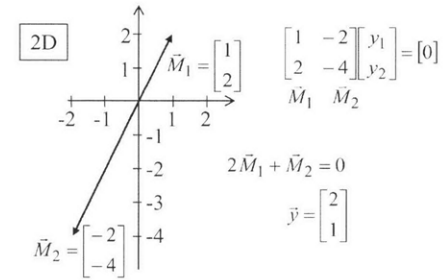


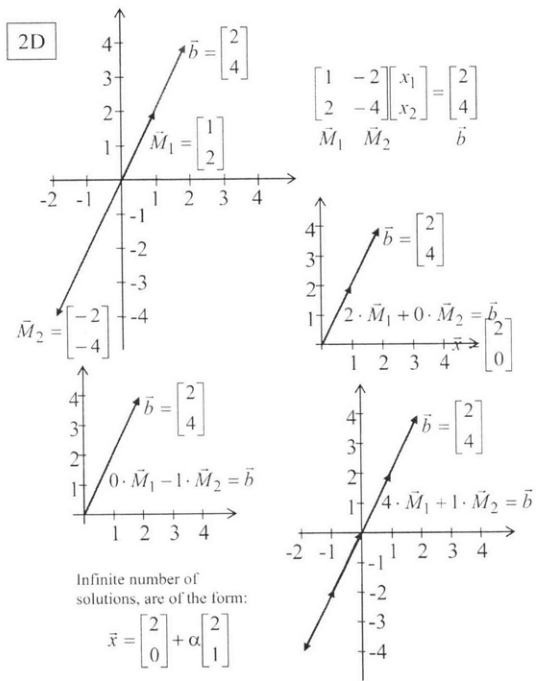


- ❖ Is the solution unique?
 Suppose there exist weights, y_1, \dots, y_N , not all zero
 $y_1 \vec{M}_1 + y_2 \vec{M}_2 + \dots + y_N \vec{M}_N = 0 \rightarrow \text{[nullspace]}$
 Then if $\mathbf{M}\vec{x} = \vec{b}$, then $\mathbf{M}(\vec{x} + \vec{y}) = \vec{b}$
 A solution is unique only if the columns of \mathbf{M} are linearly independent.
 Linearly independent columns means: $\sum_{i \neq k} \alpha_i \vec{M}_i \neq \vec{M}_k$

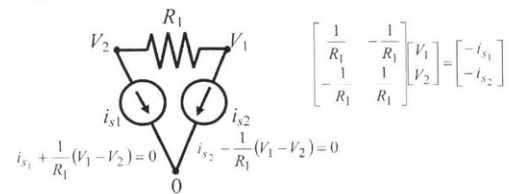
$\mathbf{M}\vec{x} = \vec{b}$	\vec{x} is a solution
$\mathbf{M}\vec{y} = \mathbf{0}$	linearly dependent columns in \mathbf{M}
$\mathbf{M}(\vec{x} + \vec{y}) = \vec{b}$	$\vec{x} + \vec{y}$ is a solution

If $\mathbf{M}\vec{y} = \mathbf{0}$ then $\alpha\mathbf{M}\vec{y} = \mathbf{0} \Rightarrow \mathbf{M}(\vec{x} + \alpha\vec{y}) = \vec{b} \Rightarrow$ infinite solutions





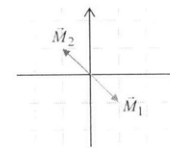
Physical Circuit Example.)



Let $R_1=1$, is there a right hand side where the system can be solved?

$$\begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} = \begin{bmatrix} -i_{s1} \\ -i_{s2} \end{bmatrix} = k \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

Physically this means that the current values must be equal and opposite. This makes sense as there is only one path for the current. Any k volt difference, where $V_1 = V_2 + k$ will give the solution.



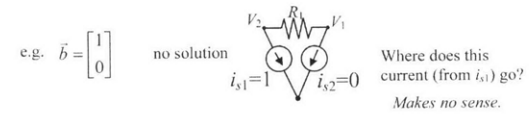
L.D. $\begin{cases} \text{no solution} \\ \text{Infinite sols} \end{cases}$

e.g. $\vec{b} = \begin{bmatrix} -2 \\ +2 \end{bmatrix}$

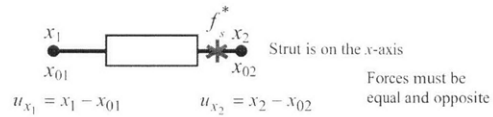
$$\begin{aligned} 0\vec{M}_1 + 2\vec{M}_2 &= \vec{b} & \rightarrow & V_1 = 0 & V_2 = 2 \\ \text{or } -\vec{M}_1 + \vec{M}_2 &= \vec{b} & \rightarrow & V_1 = -1 & V_2 = 1 \\ \text{or } 8\vec{M}_1 + 10\vec{M}_2 &= \vec{b} & \rightarrow & V_1 = 8 & V_2 = 10 \\ \text{or } -2\vec{M}_1 + 0\vec{M}_2 &= \vec{b} & \rightarrow & V_1 = -2 & V_2 = 0 \\ \text{or } -103\vec{M}_1 + 105\vec{M}_2 &= \vec{b} & \rightarrow & V_1 = -103 & V_2 = 105 \end{aligned}$$

etc...

Physically we can shift V_1 and V_2 by the same amount!



Struts and Joints Example.)



$$f_2 = f_1^* = \epsilon[(x_{02} - x_{01}) - (x_2 - x_1)] = \epsilon(u_{x_1} - u_{x_2})$$

$$f_1 = -f_2^*$$

$$\begin{bmatrix} -\epsilon & +\epsilon \\ +\epsilon & -\epsilon \end{bmatrix} \begin{bmatrix} u_{x_1} \\ u_{x_2} \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 \end{bmatrix} \quad \begin{aligned} -f_1 - \epsilon(u_{x_1} - u_{x_2}) &= 0 \\ -f_2 + \epsilon(u_{x_1} - u_{x_2}) &= 0 \end{aligned}$$

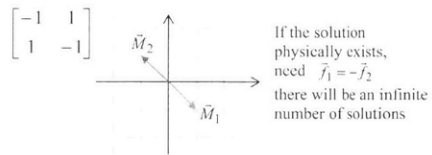
$$u_{x_1} \begin{bmatrix} \epsilon \\ -\epsilon \end{bmatrix} + u_{x_2} \begin{bmatrix} -\epsilon \\ \epsilon \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 \end{bmatrix}$$

The solution is in the span of the matrix columns if $f_1 = -f_2$.

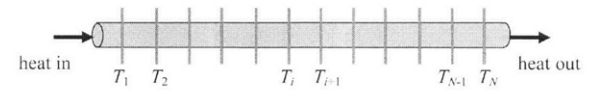
Solution is not unique because for $y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$$1 \cdot \begin{bmatrix} \epsilon \\ -\epsilon \end{bmatrix} + 1 \cdot \begin{bmatrix} -\epsilon \\ \epsilon \end{bmatrix} = 0!$$

Singular System.



Heat Conducting Bar Example.)



Boundary conditions

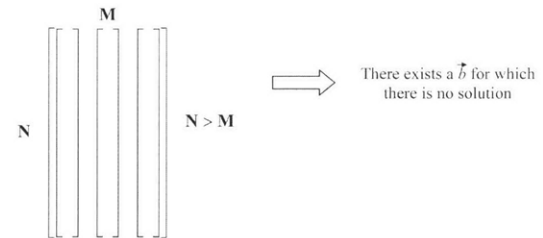
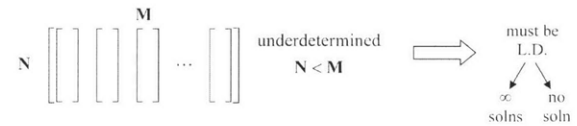
if heat in = heat out \rightarrow no unique solution

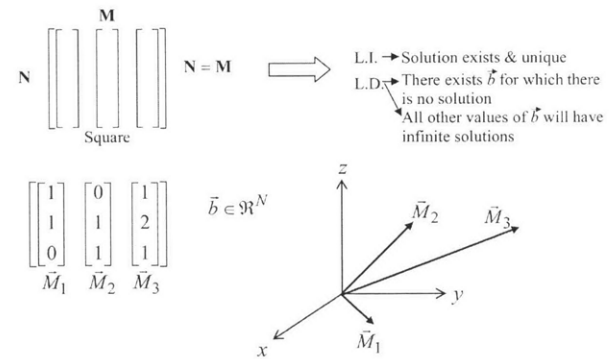
Could displace solution by a factor of 100, and it would still solve the system

if heat in \neq heat out \rightarrow no solution

Summary Table

	$b \in \text{range}\{M\}$	$b \notin \text{range}\{M\}$
L.I. columns	Solution exists and is unique	No solution
L.D. columns	Infinite solutions	No solution





Square Matrices

Given $\mathbf{M}\mathbf{x} = \mathbf{b}$, where \mathbf{M} is square

If a solution exists for all \mathbf{b} , then the solution for a specific \mathbf{b} is unique.

For a solution to exist for any \mathbf{b} , the columns of \mathbf{M} must span all N -length vectors.

Since there are only N columns of the matrix \mathbf{M} to span this space, these vectors must be linearly independent.

A square matrix with linearly independent columns is said to be **nonsingular**.

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 4. Linear Systems – LU Decomposition

TODAY'S OUTLINE:

- Solution of Linear Systems
 - Gaussian Elimination Basics
 - LU factorization
 - Computational Complexity
 - Pivoting for Growth Control

SOLUTION OF LINEAR SYSTEMS

Gaussian Elimination Basics

- ❖ LU Factorization
 - o Important Properties
 - Gaussian Elimination Method for solving $M\bar{x} = \bar{b}$
 - A "Direct" Method
 - Finite Termination for exact result (ignoring roundoff error)
 - Produces accurate results for a broad range of matrices
 - Computationally expensive

o Remainder by Example

3x3 Example

$$\begin{bmatrix} M_{11} & M_{12} & M_{13} \\ M_{21} & M_{22} & M_{23} \\ M_{31} & M_{32} & M_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \Rightarrow \begin{cases} M_{11}x_1 + M_{12}x_2 + M_{13}x_3 = b_1 \\ M_{21}x_1 + M_{22}x_2 + M_{23}x_3 = b_2 \\ M_{31}x_1 + M_{32}x_2 + M_{33}x_3 = b_3 \end{cases}$$

Use equation 1 to eliminate x_1 from equation 2 and 3

$$M_{11}x_1 + M_{12}x_2 + M_{13}x_3 = b_1$$

$$\left(M_{21} - \frac{M_{21}}{M_{11}} M_{11} \right) x_1 + \left(M_{22} - \frac{M_{21}}{M_{11}} M_{12} \right) x_2 + \left(M_{23} - \frac{M_{21}}{M_{11}} M_{13} \right) x_3 = b_2 - \frac{M_{21}}{M_{11}} b_1$$

$$\left(M_{31} - \frac{M_{31}}{M_{11}} M_{11} \right) x_1 + \left(M_{32} - \frac{M_{31}}{M_{11}} M_{12} \right) x_2 + \left(M_{33} - \frac{M_{31}}{M_{11}} M_{13} \right) x_3 = b_3 - \frac{M_{31}}{M_{11}} b_1$$

Pivot: M_{11}

Multipliers: $\frac{M_{21}}{M_{11}}$, $\frac{M_{31}}{M_{11}}$

$$\begin{bmatrix} M_{11} & M_{12} & M_{13} \\ 0 & M_{22} - \frac{M_{21}}{M_{11}} M_{12} & M_{23} - \frac{M_{21}}{M_{11}} M_{13} \\ 0 & M_{32} - \frac{M_{31}}{M_{11}} M_{12} & M_{33} - \frac{M_{31}}{M_{11}} M_{13} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 - \frac{M_{21}}{M_{11}} b_1 \\ b_3 - \frac{M_{31}}{M_{11}} b_1 \end{bmatrix}$$

Pivot should NOT BE = 0.

Simplify the notation

$$\begin{bmatrix} M_{11} & M_{12} & M_{13} \\ 0 & \tilde{M}_{22} & \tilde{M}_{23} \\ 0 & \tilde{M}_{32} & \tilde{M}_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ \tilde{b}_2 \\ \tilde{b}_3 \end{bmatrix}$$

Use updated equation 2 to eliminate x_2 from the updated equation 3

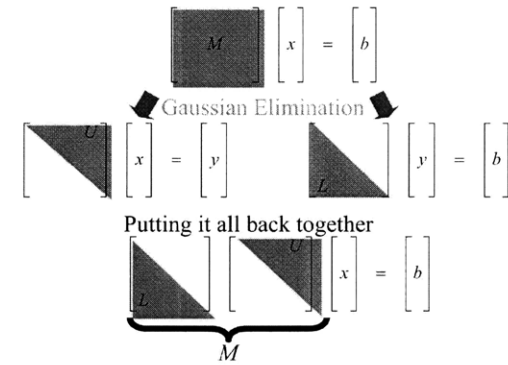
$$\begin{bmatrix} M_{11} & M_{12} & M_{13} \\ 0 & \tilde{M}_{22} & \tilde{M}_{23} \\ 0 & 0 & \tilde{M}_{33} - \frac{\tilde{M}_{32}}{\tilde{M}_{22}} \tilde{M}_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ \tilde{b}_2 \\ \tilde{b}_3 - \frac{\tilde{M}_{32}}{\tilde{M}_{22}} \tilde{b}_2 \end{bmatrix}$$

Right-hand side

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 - \frac{M_{21}}{M_{11}} b_1 \\ b_3 - \frac{M_{31}}{M_{11}} b_1 - \frac{\tilde{M}_{32}}{\tilde{M}_{22}} \tilde{b}_2 \end{bmatrix} \Rightarrow \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

$$M\bar{x} = \bar{b} \Rightarrow \underline{M} \underline{L} \bar{x} = \bar{b} \quad M \leftrightarrow LU$$

$$\Rightarrow \underline{L} \underline{U} \bar{x} = \bar{b} \quad \begin{cases} \underline{U} \bar{x} = \bar{y} \\ \underline{L} \bar{y} = \bar{b} \end{cases}$$



Fitting the pieces together

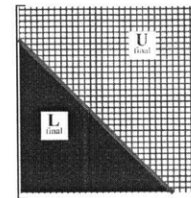
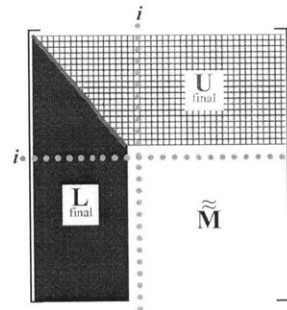
$$\begin{bmatrix} 1 & & & \\ \frac{M_{21}}{M_{11}} & 1 & & \\ \frac{M_{31}}{M_{11}} & \frac{\tilde{M}_{32}}{\tilde{M}_{22}} & 1 & \\ \frac{M_{31}}{M_{11}} & \frac{\tilde{M}_{32}}{\tilde{M}_{22}} & & 1 \end{bmatrix} \begin{bmatrix} M_{11} & M_{12} & M_{13} \\ & \tilde{M}_{22} & \tilde{M}_{23} \\ & & \tilde{M}_{33} \end{bmatrix}$$

$$\begin{matrix} \mathbf{L} & & \mathbf{U} \\ & \mathbf{L} & \mathbf{U} \\ & & \mathbf{L} & \mathbf{U} \\ & & & \mathbf{L} & \mathbf{U} \end{matrix}$$

$$\begin{bmatrix} M_{11} & M_{12} & M_{13} \\ \frac{M_{21}}{M_{11}} & \tilde{M}_{22} & \tilde{M}_{23} \\ \frac{M_{31}}{M_{11}} & \frac{\tilde{M}_{32}}{\tilde{M}_{22}} & \tilde{M}_{33} \end{bmatrix}$$

$$\mathbf{L} \mathbf{U}$$

Store all the data in one matrix – no extra memory needed.



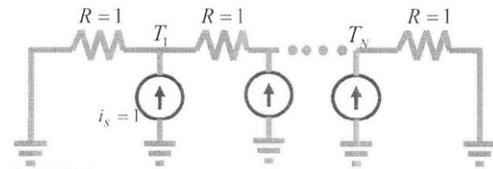
The L matrix will contain all subdiagonal elements with a diagonal of all 1's

- o Factoring
An "in place" implementation

$$\begin{bmatrix} M_{11} & M_{12} & M_{13} \\ M_{21} & M_{22} & M_{23} \\ M_{31} & M_{32} & M_{33} \end{bmatrix} \Rightarrow \begin{bmatrix} M_{11} & M_{12} & M_{13} \\ \frac{M_{21}}{M_{11}} & \tilde{M}_{22} & \tilde{M}_{23} \\ \frac{M_{31}}{M_{11}} & \tilde{M}_{32} & \tilde{M}_{33} \end{bmatrix}$$

$$\begin{bmatrix} M_{11} & M_{12} & M_{13} \\ \frac{M_{21}}{M_{11}} & \tilde{M}_{22} & \tilde{M}_{23} \\ \frac{M_{31}}{M_{11}} & \frac{\tilde{M}_{32}}{\tilde{M}_{22}} & \tilde{M}_{33} \end{bmatrix} \leftarrow \begin{bmatrix} M_{11} & M_{12} & M_{13} \\ \frac{M_{21}}{M_{11}} & \tilde{M}_{22} & \tilde{M}_{23} \\ \frac{M_{31}}{M_{11}} & \tilde{M}_{32} & \tilde{M}_{33} \end{bmatrix}$$

- Example – Heat Flow
 Temperature analogous to Voltage
 Heat Flow analogous to Current



Nodal Matrix

$$\begin{bmatrix} 2 & -1 & & & & \\ -1 & 2 & -1 & & & \\ & -1 & 2 & -1 & & \\ & & -1 & 2 & -1 & \\ & & & -1 & 2 & -1 \\ & & & & -1 & 2 & -1 \\ & & & & & -1 & 2 & -1 \\ & & & & & & -1 & 2 & -1 \\ & & & & & & & -1 & 2 & -1 \end{bmatrix} \begin{bmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \\ \vdots \\ T_{N-1} \\ T_N \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ \vdots \\ 1 \\ 1 \end{bmatrix}$$

$$G = \begin{bmatrix} 2 & -1 & & & & \\ -1 & 2 & -1 & & & \\ & -1 & 2 & -1 & & \\ & & -1 & 2 & -1 & \\ & & & -1 & 2 & -1 \\ & & & & -1 & 2 & -1 \\ & & & & & -1 & 2 & -1 \\ & & & & & & -1 & 2 & -1 \\ & & & & & & & -1 & 2 & -1 \end{bmatrix} \Rightarrow \begin{bmatrix} 2 & -1 & & & & \\ -1 & 2 & -1 & & & \\ & -1 & 2 & -1 & & \\ & & -1 & 2 & -1 & \\ & & & -1 & 2 & -1 \\ & & & & -1 & 2 & -1 \\ & & & & & -1 & 2 & -1 \\ & & & & & & -1 & 2 & -1 \\ & & & & & & & -1 & 2 & -1 \end{bmatrix}$$

1st step multiplier

$$\Rightarrow \begin{bmatrix} 2 & -1 & & & & \\ -1 & 2 & -1 & & & \\ & -1 & 2 & -1 & & \\ & & -1 & 2 & -1 & \\ & & & -1 & 2 & -1 \\ & & & & -1 & 2 & -1 \\ & & & & & -1 & 2 & -1 \\ & & & & & & -1 & 2 & -1 \\ & & & & & & & -1 & 2 & -1 \end{bmatrix}$$

2nd step multiplier

$$\Rightarrow \begin{bmatrix} 2 & -1 & & & & \\ -1 & 2 & -1 & & & \\ & -1 & 2 & -1 & & \\ & & -1 & 2 & -1 & \\ & & & -1 & 2 & -1 \\ & & & & -1 & 2 & -1 \\ & & & & & -1 & 2 & -1 \\ & & & & & & -1 & 2 & -1 \\ & & & & & & & -1 & 2 & -1 \end{bmatrix}$$

In place LU decomposition

$$M = L \cdot U = \begin{bmatrix} 1 & & & & & \\ -1/2 & 1 & & & & \\ & -2/3 & 1 & & & \\ & & -3/4 & 1 & & \\ & & & -4/5 & 1 & \\ & & & & -5/6 & 1 \\ & & & & & -6/7 & 1 \\ & & & & & & -7/8 & 1 \\ & & & & & & & -8/9 & 1 \\ & & & & & & & & -9/10 & 1 \end{bmatrix} \begin{bmatrix} 2 & -1 & & & & \\ -1 & 2 & -1 & & & \\ & -1 & 2 & -1 & & \\ & & -1 & 2 & -1 & \\ & & & -1 & 2 & -1 \\ & & & & -1 & 2 & -1 \\ & & & & & -1 & 2 & -1 \\ & & & & & & -1 & 2 & -1 \\ & & & & & & & -1 & 2 & -1 \\ & & & & & & & & -1 & 2 & -1 \end{bmatrix}$$

- Computational Complexity
 - Three Steps to Solving a System
 Solve $Mx = b$
 Step 1: Decomposition
 $M = L \cdot U$
 $\begin{bmatrix} \blacksquare & & \\ & \blacksquare & \\ & & \blacktriangledown \end{bmatrix} = \begin{bmatrix} \blacksquare & & \\ & \blacksquare & \\ & & \blacktriangledown \end{bmatrix} \cdot \begin{bmatrix} \blacksquare & & \\ & \blacksquare & \\ & & \blacktriangledown \end{bmatrix}$
 Step 2: Forward Elimination
 Solve $Ly = b$
 Step 3: Backward Substitution
 Solve $Ux = y$

- Step 2: Solving Lower Triangular Systems

$$\begin{bmatrix} l_{11} & & & & 0 \\ l_{21} & l_{22} & & & \\ l_{31} & & l_{33} & & \\ \vdots & & & \ddots & \\ l_{N1} & & & & l_{NN} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_N \end{bmatrix}$$

$$y_1 = \frac{1}{l_{11}} b_1$$

$$y_2 = \frac{1}{l_{22}} (b_2 - l_{21} y_1)$$

$$y_3 = \frac{1}{l_{33}} (b_3 - l_{31} y_1 - l_{32} y_2)$$

$$\vdots$$

$$y_N = \frac{1}{l_{NN}} (b_N - \sum_{i=1}^{N-1} l_{Ni} y_i)$$

Number of multiplications:

$$1 + 2 + 3 + 4 + \dots + (n-1) = \sum_{i=1}^{n-1} i = \frac{(n-1)n}{2} = O(n^2)$$

- Step 3: Solving Upper Triangular Systems

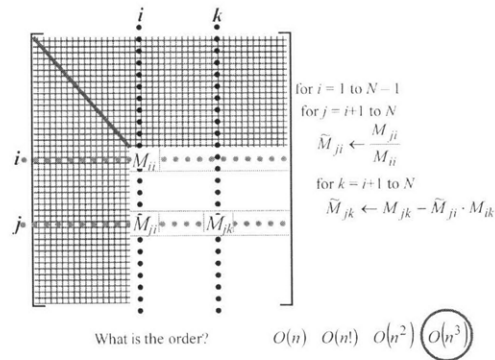
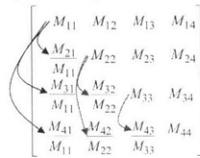
$$\begin{bmatrix} U_{11} & U_{12} & U_{13} \\ 0 & U_{22} & U_{23} \\ 0 & 0 & U_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$$\begin{aligned} x_3 &= \frac{y_3}{U_{33}} \\ x_2 &= \frac{y_2 - U_{23}x_3}{U_{22}} \\ x_1 &= \frac{y_1 - U_{12}x_2 - U_{13}x_3}{U_{11}} \end{aligned}$$

Number of multiplications:

$$\sum_{i=1}^{n-1} \sum_{j=i+1}^n (n-j) = O(n^2)$$

- Step 1: Decomposition



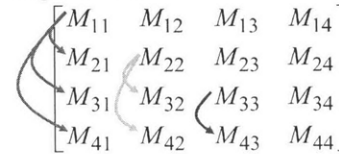
- Algorithm

```

For i = 1 to n - 1 {
  For j = i + 1 to n {
    "For each Row"
    "For each target Row below the source"
    M_ji = M_ji / M_ii -> Pivot
    "sum_{l=i+1}^{n-1} (n-l) = n^2/2 multipliers"
    For k = i + 1 to n {
      "For each Row element beyond Pivot"
      M_jk <- M_jk - M_ji * M_ik
    }
  }
}
Multiplier
sum_{i=1}^{n-1} (n-i)^2 = 2/3 n^3
Multiply-adds

```

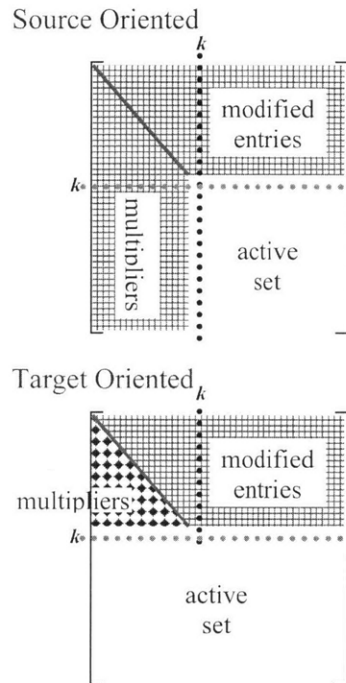
"Target Oriented"

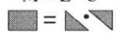


```

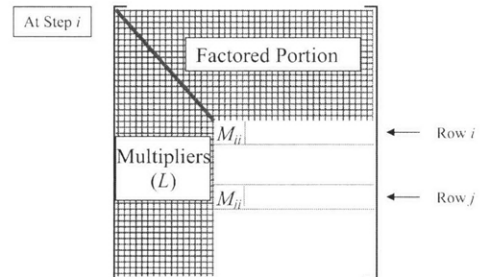
for i = 2 to N
  for j = 1 to i - 1
    "for each target"
    "for each source"
    M-tilde_ji <- M_ji / M_ii
    target source
  for k = j+1 to N
    M-tilde_jk <- M_jk - M-tilde_ji * M_jk

```



- o Summary
 - Solve $Mx = b$
 - Step 1 Decomposition
 - $M = L \cdot U$
 - 
 - $O(n^3)$
 - Step 2 Forward Elimination
 - Solve $Ly = b$
 - $O(n^2)$
 - Step 3 Backward Substitution
 - Solve $Ux = y$
 - $O(n^2)$

- ❖ Pivoting for Growth Control
 - o Zero Pivots



What if $M_{ii} = 0$? $\frac{M_{ji}}{M_{ii}}$
 Cannot form $\frac{M_{ji}}{M_{ii}}$

Simple Fix (Partial Pivoting)

If $M_{ii} = 0$
 Find $M_{ij} \neq 0, j > i$
 Swap Row j with Row i
 Pick row j whose M_{ij} value is large

Swap rows 2 and 4

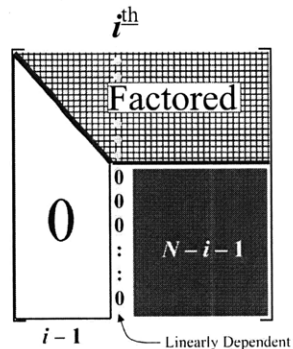
$$\begin{bmatrix} M_{11} & M_{12} & M_{13} & M_{14} \\ M_{21} & M_{22} & M_{23} & M_{24} \\ M_{31} & M_{32} & M_{33} & M_{34} \\ M_{41} & M_{42} & M_{43} & M_{44} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix}$$

$$\begin{bmatrix} M_{11} & M_{12} & M_{13} & M_{14} \\ M_{41} & M_{42} & M_{43} & M_{44} \\ M_{31} & M_{32} & M_{33} & M_{34} \\ M_{21} & M_{22} & M_{23} & M_{24} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_4 \\ b_3 \\ b_2 \end{bmatrix}$$

Two Important Theorems

1. Partial pivoting (swapping rows) always succeeds if **M** is nonsingular.

Creating U



(N - 1) Linearly Independent Columns

2. LU factorization applied to a strictly diagonally dominant matrix will never produce a zero pivot.

Small Pivots

Example [Singular Matrix]

$$\begin{bmatrix} 1 & -1 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Singular matrix produces a row of zeros

If **M** is singular, will not be able to solve for some unknowns.

Contrived Example

$$\begin{bmatrix} 10^{-17} & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

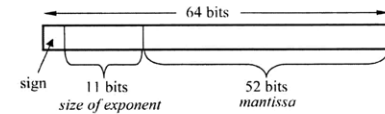
$$L = \begin{bmatrix} 1 & 0 \\ 10^{17} & 1 \end{bmatrix} \quad U = \begin{bmatrix} 10^{-17} & 1 \\ 0 & 2 - 1 \cdot 10^{17} \end{bmatrix}$$

Can we represent this?

An Aside on Floating Point Arithmetic

Double Precision Number

$$X.XXXX...X \cdot 10^{\text{exponent}}$$



How large a number can we represent with 52 bits? How many decimal digits is it?

- 3 bits $\rightarrow 2^3 = 8 \approx 10 \rightarrow 1$ digit
 - 7 bits $\rightarrow 2^7 = 128 \approx 100 \rightarrow 2$ digits
 - 10 bits $\rightarrow 2^{10} = 1024 \approx 1000 \rightarrow 3$ digits
- [Remember cubits \leftrightarrow 3 decimal digits]

52 bits $\rightarrow 52 \text{ bits} \cdot \frac{3 \text{ digits}}{10 \text{ bits}} \approx 15 - 16$ decimal digits!

52 bits is the precision used in MATLAB, for example

Basic Problem

A.) $1.0000001 - 0.000000000000001 = 1.0$

B.) $(\pi \cdot 10^{-7}) + 1 - 1.0000001 = ?$

Look at order of procedure:

Ex #1.) $-1.0000001 + 1 = -1 \times 10^{-7}$

$$-1 \times 10^{-7} + \pi \cdot 10^{-7} = \underbrace{2.141592653589793}_{\text{correct?}} \cdot 10^{-7}$$

Ex #2.) All 15 decimal digits are correct in this case

$$1 + \pi \cdot 10^{-7} = \underbrace{1.000000314159265}_{\text{lost 7 digits of } \pi}$$

$$-1.0000001 + 1.000000314159265 = 2.141592652105118 \cdot 10^{-7}$$

Only the first 8 decimal digits are correct in this case, the remainder are garbage. Lose digits of precision

Key Issue
 Avoid additions and subtractions between large and small numbers!
 EVEN BETTER: AVOID GENERATING LARGE NUMBERS AT ALL IF POSSIBLE!!!

Back to the contrived Example

$$\begin{bmatrix} 10^{-17} & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

$$LU_{Exact} = \begin{bmatrix} 1 & 0 \\ 10^{17} & 1 \end{bmatrix} \begin{bmatrix} 10^{-17} & 1 \\ 0 & 2 - 10^{17} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

$$LU_{Rounded} = \begin{bmatrix} 1 & 0 \\ 10^{17} & 1 \end{bmatrix} \begin{bmatrix} 10^{-17} & 1 \\ 0 & -10^{17} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_{Exact} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_{Rounded} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad \times$$

Original Problem:

$$\begin{bmatrix} 10^{-17} & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

Swap the rows

$$\begin{bmatrix} 1 & 2 \\ 10^{-17} & 1 \end{bmatrix}$$

Partial Pivoting for Roundoff Reduction

If $|M_{ii}| < \max_{j>i} |M_{ji}|$ swap row i with $\arg(\max_{j>i} |M_{ji}|)$

$$LU_{reordered} = \begin{bmatrix} 1 & 0 \\ 10^{-17} & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 0 & 1 - 2 \cdot 10^{-17} \end{bmatrix}$$

This multiplier is small This term still gets rounded, but this time the multiplier does not 'overpower' the row.

$$L = \begin{bmatrix} 1 & 0 \\ 10^{-17} & 1 \end{bmatrix} \quad U = \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix}$$

Let's solve $LUx = b$.

First, solve $Ly = b$:

$$Ly = b \rightarrow \begin{bmatrix} 1 & 0 \\ 10^{-17} & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

$$\Rightarrow y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 - 10^{-17} \end{bmatrix} \approx \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

Next, solve $Ux = y$:

$$Ux = y \rightarrow \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix} \Rightarrow x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -5 \\ 3 \end{bmatrix} \quad \times \text{ STILL WRONG}$$

PROBLEM: Need to swap the elements of b as well!

Solve $Ly = b$:

$$Ly = b \rightarrow \begin{bmatrix} 1 & 0 \\ 10^{-17} & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$$

$$\Rightarrow y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 3 \\ 1 - 3 \cdot 10^{-17} \end{bmatrix} \approx \begin{bmatrix} 3 \\ 1 \end{bmatrix}$$

Next, solve $Ux = y$:

$$Ux = y \rightarrow \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \end{bmatrix} \Rightarrow x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \text{ Correct Answer.}$$

Example.

$$\begin{bmatrix} 1 & 0 & 0 & \dots & 0 & 1 \\ 10 & 1 & 0 & \dots & 0 & 0 \\ 0 & 10 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 10 & 1 & 0 \\ 0 & 0 & \dots & 0 & 10 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_n \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 1 & 0 & 0 & \dots & 1 \\ 10 & 1 & 0 & \dots & -10 \\ & 10 & 1 & \dots & 0 \\ & & \ddots & \ddots & \vdots \\ & & & 10 & 1 \\ & & & & x_N \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_n \end{bmatrix}$$

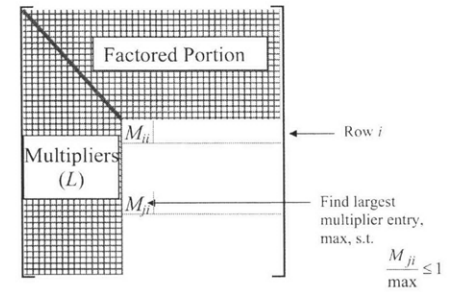
$$\Rightarrow \begin{bmatrix} 1 & 0 & 0 & \dots & 1 \\ 10 & 1 & 0 & \dots & -10 \\ 0 & 10 & 1 & \dots & -100 \\ & \ddots & \ddots & \ddots & \vdots \\ & & & 10 & 1 \\ & & & & x_N \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_n \end{bmatrix}$$

$$\Rightarrow \dots \Rightarrow \begin{bmatrix} 1 & 0 & 0 & \dots & 1 \\ 10 & 1 & 0 & \dots & -10 \\ 0 & 10 & 1 & \dots & -100 \\ & \ddots & \ddots & \ddots & \vdots \\ & & & 10 & -10^{N-1} \\ & & & & x_N \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_n \end{bmatrix}$$

The last entry is very large! Use partial pivoting – swap first and second rows...

- If the matrix is strictly diagonally dominant
- or if use partial pivoting for round-off reduction:
 1. The multipliers will always be smaller than one in magnitude.
 2. The maximum magnitude entry in the LU factors will never be larger than $2^{(n-1)}$ times the maximum magnitude entry in the original matrix.

For 1000 nodes, what is 2^{1000} ?
 Know that $2^{10} \approx 10^3$.
 So, $2^{1000} = (2^{10})^{100} = 10^{300} \rightarrow$ Very large number!
 Might not be a very useful theorem \rightarrow generally not this large.



INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 5.

Linear Systems – Conditioning

TODAY'S OUTLINE:

- Solution of Dense Linear Systems
 - Hard to Solve Problems
 - Perturbation Analysis and Conditioning
- Solution of Sparse Linear Systems
 - LU Factorization Reminder.
 - Example of Problems with Sparse Matrices
 - Struts and joints, resistor grids, 3-D heat flow
 - Tridiagonal Matrix Factorization
 - General Sparse Factorization
 - Fill-in and Reordering
 - Graph Based Approach
 - Sparse Matrix Data Structures
 - Scattering

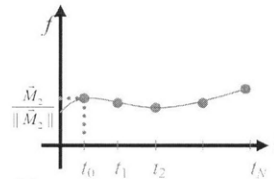
SOLUTION OF DENSE LINEAR SYSTEMS

Hard to Solve Problems

- ❖ Fitting example
- Polynomial Interpolation

Table of Data

t_0	$f(t_0)$
t_1	$f(t_1)$
\vdots	\vdots
t_N	$f(t_N)$



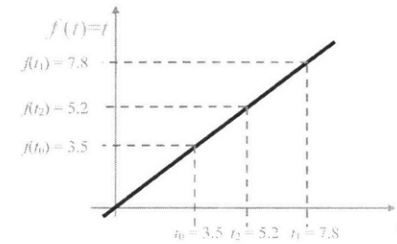
Problem: fit data with an N^{th} order polynomial

$$f(t) = \alpha_0 + \alpha_1 t + \alpha_2 t^2 + \alpha_3 t^3 + \dots + \alpha_N t^N$$

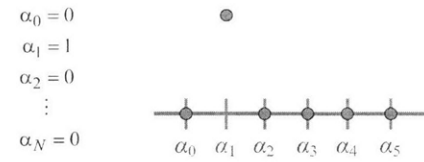
Matrix Form

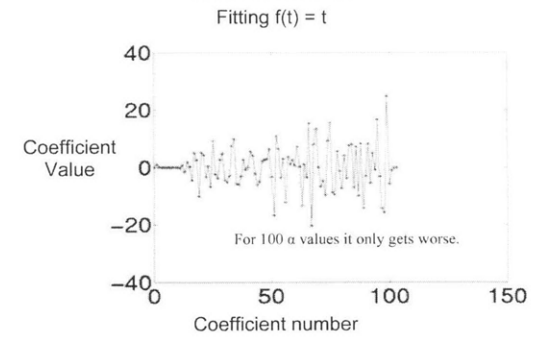
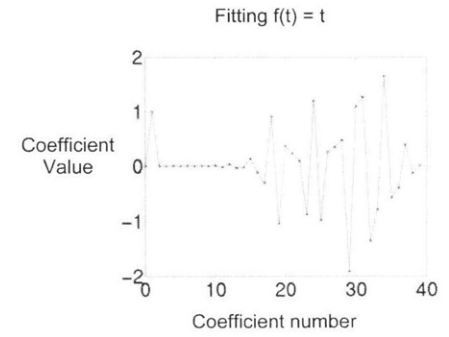
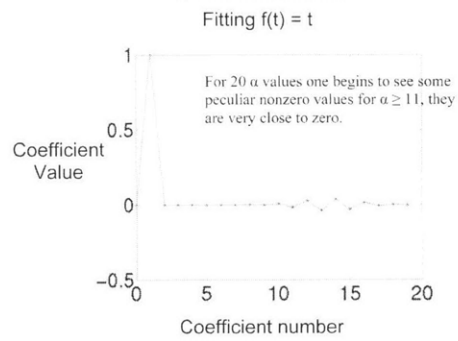
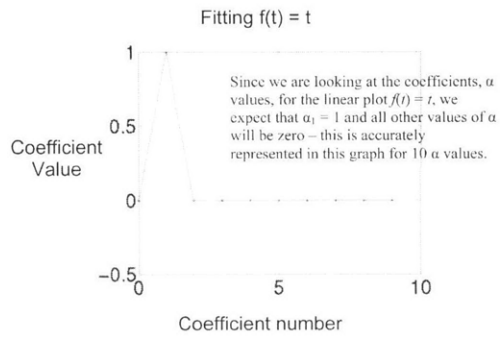
$$\underbrace{\begin{bmatrix} 1 & t_0 & t_0^2 & \dots & t_0^N \\ 1 & t_1 & t_1^2 & \dots & t_1^N \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & t_N & t_N^2 & \dots & t_N^N \end{bmatrix}}_{M_{\text{interp}}} \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_N \end{bmatrix} = \begin{bmatrix} f(t_0) \\ f(t_1) \\ \vdots \\ f(t_N) \end{bmatrix}$$

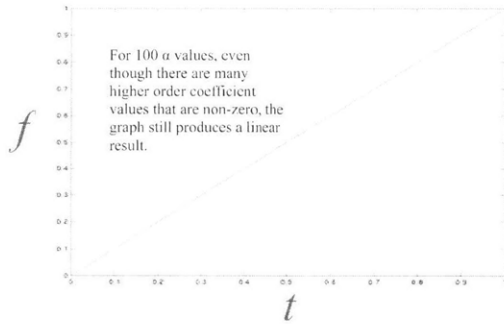
$$\begin{aligned} \alpha_0 + \alpha_1 t_0 + \alpha_2 t_0^2 + \dots + \alpha_N t_0^N &= f(t_0) \\ \alpha_0 + \alpha_1 t_1 + \alpha_2 t_1^2 + \dots + \alpha_N t_1^N &= f(t_1) \\ &\vdots \\ \alpha_0 + \alpha_1 t_N + \alpha_2 t_N^2 + \dots + \alpha_N t_N^N &= f(t_N) \end{aligned}$$



$$\alpha_0 + \alpha_1 t + \alpha_2 t^2 + \dots = f(t) = t$$

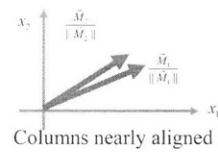
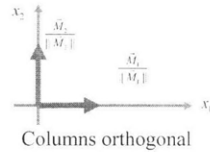






❖ Geometric Approach is clearer

$M = [\vec{M}_1, \vec{M}_2]$ Solving $Mx = b$ is finding $x_1\vec{M}_1 + x_2\vec{M}_2 = \vec{b}$



When vectors are nearly aligned, difficult to determine how much of \vec{M}_1 versus how much of \vec{M}_2

$M = \begin{bmatrix} 1 & 0 \\ 0 & 10^{-6} \end{bmatrix}$
orthogonal columns

$M = \begin{bmatrix} 1 & 1-10^{-6} \\ 1-10^{-6} & 1 \end{bmatrix}$
columns nearly aligned

Example. Orthogonal vectors - exactly one solution.

$$\begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 4 \\ 1 \end{bmatrix} \quad \vec{b} = 2 \cdot \vec{M}_1 + 1 \cdot \vec{M}_2$$

Physical meaning of above problem in the form of a circuit.

Example.

$$\begin{bmatrix} +1 & -1 \\ -1 & +1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 4 \\ 1 \end{bmatrix}$$

No possible solution.

Physical meaning of above problem in the form of a circuit.

Example.

$$\begin{bmatrix} +1 & -1 \\ -1 & +1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -2 \\ 2 \end{bmatrix}$$

Physical meaning of above problem in the form of a circuit.

Solution(s):

$x_1 = -1$	$x_2 = 1$	$-1 \cdot \vec{M}_1 + 1 \cdot \vec{M}_2$
$x_1 = -2$	$x_2 = 0$	$-2 \cdot \vec{M}_1$
$x_1 = 0$	$x_2 = 2$	$2 \cdot \vec{M}_2$
\vdots	\vdots	\vdots

Infinite number of solutions

Example.

Node 1: $f_{L1} - \varepsilon(u_{x1} - u_{x2}) = 0$

Node 2: $f_{L2} + \varepsilon(u_{x1} - u_{x2}) = 0$

$$\begin{bmatrix} -\varepsilon & \varepsilon \\ \varepsilon & -\varepsilon \end{bmatrix} \begin{bmatrix} u_{x1} \\ u_{x2} \end{bmatrix} = \begin{bmatrix} -f_{L1} \\ -f_{L2} \end{bmatrix}$$

either infinite number of solutions or none

Same problem as circuit with two current sources (above)

Linearly dependent vectors.

$$\begin{bmatrix} 2 & 2 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 4 \\ 1 \end{bmatrix}$$

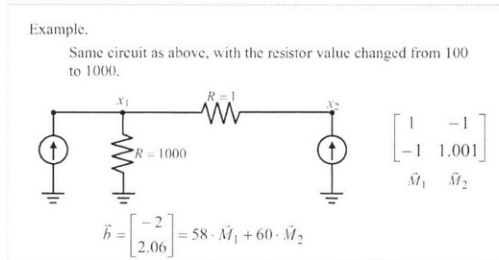
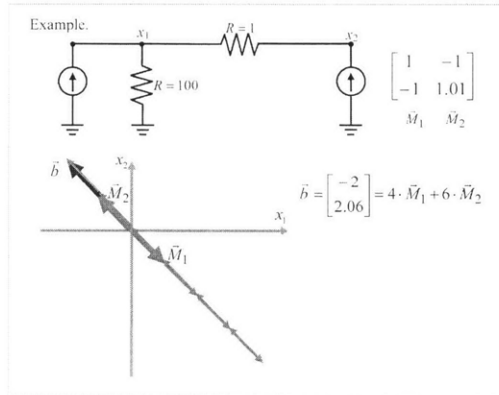
Infinite number of solutions

$\vec{b} = 2 \cdot \vec{M}_1 = 2 \cdot \vec{M}_2 = \vec{M}_1 + \vec{M}_2 = \text{etc.}$

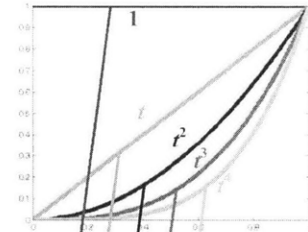
Example.

$$\vec{b} = \underset{x_1}{5} \cdot \vec{M}_1 - \underset{x_2}{3} \cdot \vec{M}_2$$

Columns are very close to being linearly dependent (closely aligned) – have a larger condition number / columns are not close to orthogonal



❖ Polynomial Interpolation



The power series polynomials are nearly linearly dependent. (i.e. if viewed as vectors they are almost aligned)

$$\begin{bmatrix} 1 & t_1 & t_1^2 & t_1^3 & t_1^4 & \dots \\ 1 & t_2 & t_2^2 & t_2^3 & t_2^4 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & t_N & t_N^2 & t_N^3 & t_N^4 & \dots \end{bmatrix}$$

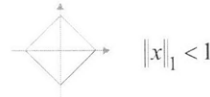
Perturbation Analysis and Conditioning

- ❖ Induced Norms
- Matrix Magnification Question
 - Suppose $y = Mx$
 - How much larger is y than x ? OR How much does M magnify x ?

o Vector Norm Review

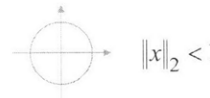
- L₁ norm:

$$\|x\|_1 = \sum_{i=1}^n |x_i|$$



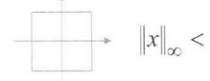
- L₂ (Euclidean) norm:

$$\|x\|_2 = \sqrt{\sum_{i=1}^n |x_i|^2}$$



- L_∞ norm:

$$\|x\|_\infty = \max_i |x_i|$$



If $\|x\|_2 < 1$

yes $\|x\|_\infty < 1$

no $\|x\|_1 < 1$

Try $x = \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$

$\|x\|_2 = \sqrt{\left(\frac{1}{\sqrt{2}}\right)^2 + \left(\frac{1}{\sqrt{2}}\right)^2} = 1$

$\|x\|_1 = \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{2}} = \sqrt{2} > 1$

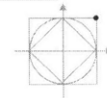
Which one? $\|x\|_1 \leq \|x\|_2 \leq \|x\|_\infty$
 $\|x\|_\infty \leq \|x\|_2 \leq \|x\|_1$ *

Example.) Try $x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$$\|x\|_1 = 1 + 1 = 2$$

$$\|x\|_2 = \sqrt{1^2 + 1^2} = \sqrt{2}$$

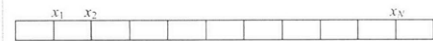
$$\|x\|_\infty = \max(1, 1) = 1$$



Is there an x such that $\|x\|_\infty = \|x\|_2 = \|x\|_1$?

$$x = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \Rightarrow \|x\|_\infty = \|x\|_2 = \|x\|_1 = 1$$

Example.) Two heat conducting bars.



$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

Vector of temperatures

Given:

Bar A: $|x|_x < 100^\circ$

Bar B: $|x|_1 < 100^\circ$

Which one do you touch?

Bar B will likely be cooler to the touch since all nodes on bar B will sum to 100° , whereas on bar A, the maximum node value is 100° .

Worst-case-scenario, bar A will be 100° everywhere and bar B will be 100° on one node and 0° elsewhere

o Matrix Norms [Standard Induced l-norms]

- Definition:

$$\|M\|_l = \max_x \frac{\|Mx\|_l}{\|x\|_l} = \max_{\|x\|_l=1} \|Mx\|_l$$

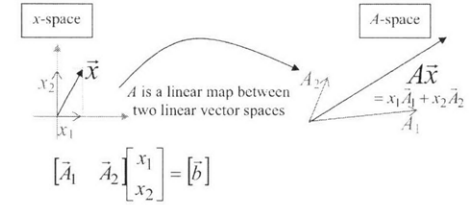
- $\|M\|_1 = \max_i \sum_{j=1}^N |M_{ij}| = \text{max abs column sum}$
- $\|M\|_\infty = \max_j \sum_{i=1}^N |M_{ij}| = \text{max abs row sum}$
- $\|M\|_2 = \max_i \sqrt{\lambda_i(M^H M)} = \max_i |\lambda_i(M)| = \text{max eigenvalue}$ for symmetric, real matrices

$\|M\|_1 = \max_j \sum_{i=1}^n |M_{ij}| = \text{max abs column sum}$

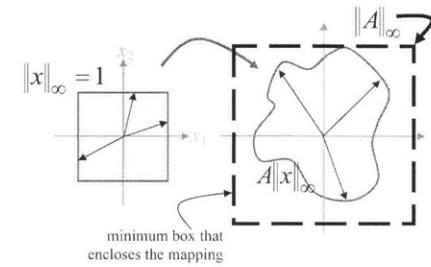
$$\begin{bmatrix} \mathbf{M} & \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \end{bmatrix} = \begin{bmatrix} m_{11} \\ m_{21} \\ \vdots \\ m_{n1} \end{bmatrix} \quad \begin{bmatrix} \mathbf{M} & \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \end{bmatrix} = \begin{bmatrix} m_{12} \\ m_{22} \\ \vdots \\ m_{n2} \end{bmatrix}$$

Which unit vector input will produce the largest output?

$\|M\|_\infty = \max_i \sum_{j=1}^n |M_{ij}| = \text{max abs row sum}$

$$\begin{bmatrix} \mathbf{M} & \begin{bmatrix} \pm 1 \\ \pm 1 \\ \vdots \\ \pm 1 \end{bmatrix} \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^n \pm m_{1j} \\ \sum_{j=1}^n \pm m_{2j} \\ \vdots \\ \sum_{j=1}^n \pm m_{ij} \end{bmatrix}$$


$\|A\|$ is a measure of the largest gain or magnification power



Example.) Given the following A , find $\|A\|_\infty$

$$A = \begin{bmatrix} 2 & 7 & -4 \\ 5 & -2 & 3 \\ 1 & -2 & 4 \end{bmatrix}$$

$$\|A\|_\infty = \max \begin{pmatrix} 2+7+4 \\ 5+2+3 \\ 1+2+4 \end{pmatrix} = \max\{13, 10, 7\} = 13$$

❖ Hard to Solve Systems

○ Perturbation Analysis

$$\frac{\|\delta x\|}{\|x + \delta x\|} \stackrel{?}{\rightarrow} \frac{\|\delta M\|}{\|M\|} \quad \begin{matrix} [M + \delta M] & [x + \delta x] & = & b \\ M & x & = & b \end{matrix}$$

Perturbation Equation

$$\underbrace{(M + \delta M)}_{\text{models LU roundoff}} \underbrace{(x + \delta x)}_{\text{models solution perturbation}} = Mx + \delta Mx + M\delta x + \delta M\delta x = \underbrace{\tilde{b}}_{\text{unperturbed right hand side}}$$

Since $Mx - b = 0$

$$M\delta x = -\delta M(x + \delta x) \Rightarrow \delta x = -M^{-1}\delta M(x + \delta x)$$

Taking Norms

$$\|\delta x\| \leq \|M^{-1}\| \|M\| \frac{1}{\|M\|} \|\delta M\| \|x + \delta x\|$$

Relative Error Relation

$$\frac{\|\delta x\|}{\|x + \delta x\|} \leq \underbrace{\|M^{-1}\| \|M\|}_{\text{"Condition Number"}} \frac{\|\delta M\|}{\|M\|}$$

$$M\tilde{x} = \tilde{b}$$

Assume $\frac{\|\delta M\|}{\|M\|} = 10^{-16}$ What is $\frac{\|\delta x\|}{\|x\|} = \frac{\|\delta x\|}{\|x + \delta x\|} = ?$

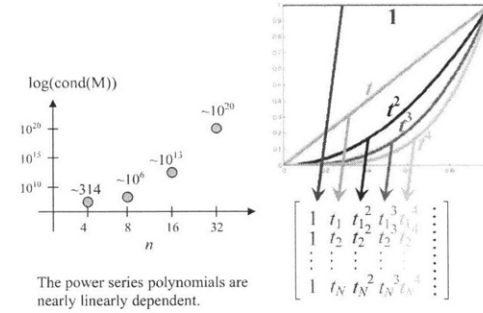
e.g. $\|M^{-1}\| \|M\| = 100 \quad \frac{\|\delta x\|}{\|x + \delta x\|} \leq 100 \cdot 10^{-16} = 10^{-14}$

e.g. $\|M^{-1}\| \|M\| = 10^{13} \quad \frac{\|\delta x\|}{\|x + \delta x\|} \leq 10^{13} \cdot 10^{-16} = 10^{-3} = 0.1\% \text{ error}$

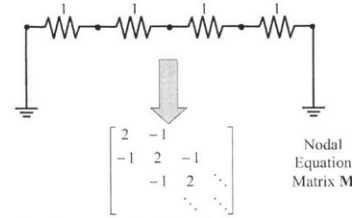
e.g. $\|M^{-1}\| \|M\| = 10^{20} \quad \frac{\|\delta x\|}{\|x + \delta x\|} \leq 10^{20} \cdot 10^{-16} = 10^4 = 10^6\% \text{ error!!}$

Large condition number \rightarrow Big error

○ Geometric Analysis / Polynomial Interpolation

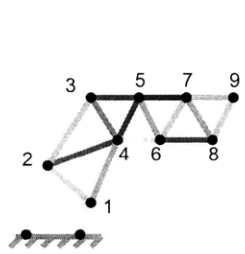


SOLUTION OF SPARSE LINEAR SYSTEMS
Examples of Problems with Sparse Matrices



number of nonzeros $\approx 3N - 2 = O(N)$
 $\frac{3N}{N^2} = \frac{3000}{10^6} = 0.003 = 0.3\%$ of the entries are nonzero for a 1000×1000 matrix

❖ Space Frame



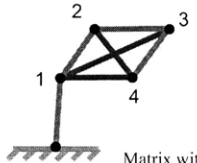
Nodal Matrix

$$\begin{bmatrix} X & & & & & & & & \\ X & X & X & & & & & & \\ X & X & X & X & & & & & \\ X & X & X & X & X & & & & \\ & X & X & X & X & X & & & \\ & & X & X & X & X & X & & \\ & & & X & X & X & X & X & \\ & & & & X & X & X & X & X \\ & & & & & X & X & X & X \end{bmatrix}$$

Unknowns : Joint positions
Equations : $\sum \text{forces} = 0$

$$X = \begin{bmatrix} \bullet & & & \\ \bullet & \bullet & & \\ \bullet & \bullet & \bullet & \\ \bullet & \bullet & \bullet & \bullet \end{bmatrix}$$

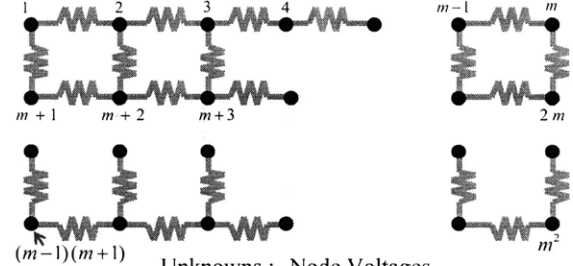
2x2 block



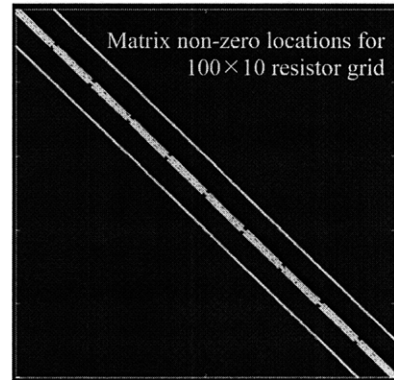
$$\begin{bmatrix} X & X & X & X \\ X & X & X & X \\ X & X & X & X \\ X & X & X & X \end{bmatrix}$$

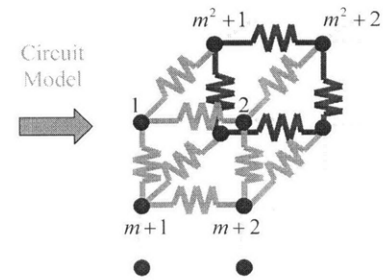
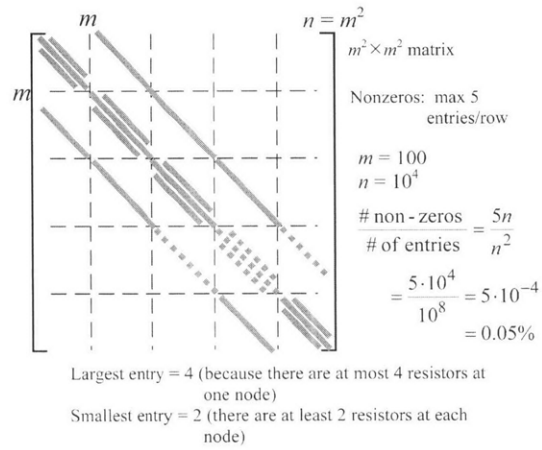
Matrix with all non-zero entries.
Struts connect all every joint in the frame.

❖ Resistor Grid

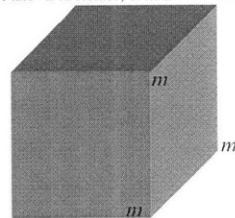


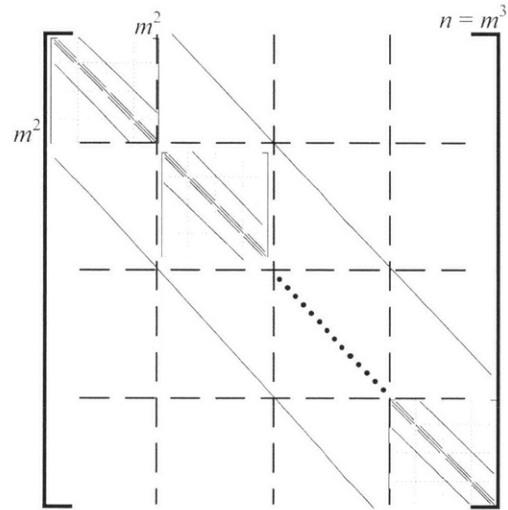
Unknowns : Node Voltages
Equations : $\sum \text{currents} = 0$





- ❖ Temperature in a cube
 Temperature known on surface, determine interior temperature





$m = 100$
 $n = m^3 = 10^6$
 $\frac{\# \text{ non-zeros}}{\# \text{ of entries}} = \frac{7n}{n^2} = \frac{7 \cdot 10^6}{10^{12}} = 7 \cdot 10^{-6} = 0.0007\%$
 Largest entry = 6
 Smallest entry = 3

Tridiagonal Matrix Factorization

Matrix Form

$\begin{bmatrix} x & x & & & \\ x & x & x & & \\ & x & x & x & \\ & & x & x & x \\ & & & x & x & x \\ & & & & x & x & x \\ & & & & & x & x & x \\ & & & & & & x & x & x \\ & & & & & & & x & x & x \end{bmatrix}$

Tridiagonal Matrix

$$\begin{bmatrix} 1 & -1 & & & \\ -1 & 2 & -1 & & \\ & -1 & 2 & -1 & \\ & & -1 & 2 & -1 \\ & & & \ddots & \ddots & \ddots \end{bmatrix}$$

* updated entries

$$\begin{bmatrix} 1 & -1 & & & \\ * & * & -1 & & \\ & * & * & -1 & \\ & * & * & * & -1 \\ & & * & * & * & -1 \\ & & & * & * & * & -1 \\ & & & & * & * & * & -1 \\ & & & & & * & * & * & -1 \\ & & & & & & * & * & * & -1 \end{bmatrix}$$

GE Algorithm for tridiagonal matrix

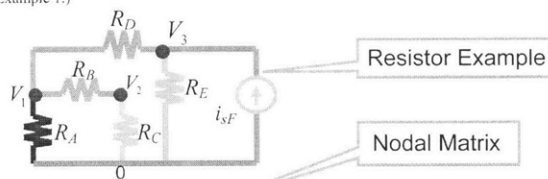
```

For i = 1 to n-1 {
  "For each Row"
  For j = i+1 to n {
    "For each target Row below the source"
     $M_{ji} = \frac{M_{ji}}{M_{ii}}$  Pivot
    For k = i+1 to n {
      "For each Row element beyond Pivot"
       $M_{jk} \leftarrow M_{jk} - M_{ji} M_{ik}$ 
      Multiplier
    }
  }
}
  
```

Order N Operations!

General Sparse Factorization

- ❖ Fill-in and Reordering
 - Example 1.)



$$\begin{bmatrix} \frac{1}{R_A} + \frac{1}{R_B} + \frac{1}{R_D} & -\frac{1}{R_B} & 0 \\ -\frac{1}{R_B} & \frac{1}{R_B} + \frac{1}{R_C} & 0 \\ 0 & 0 & \frac{1}{R_D} + \frac{1}{R_E} \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \\ V_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ i_{sF} \end{bmatrix}$$

Symmetric Diagonally Dominant

Matrix Non Zero Structure

$$\begin{bmatrix} X & X & X \\ X & X & 0 \\ X & 0 & X \end{bmatrix}$$

X = Non zero

Matrix after one LU step

$$\begin{bmatrix} X & X & X \\ X & X & X \\ X & X & X \end{bmatrix}$$

- Example 2.)

Fill-ins Propagate

$$\begin{bmatrix} X & X & X & X \\ X & X & 0 & 0 \\ 0 & X & X & 0 \\ 0 & X & 0 & 0 \end{bmatrix}$$

Fill-ins from Step 1 result in Fill-ins in step 2

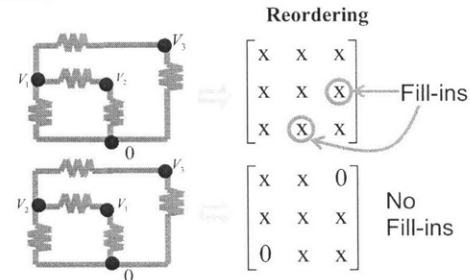
$$\begin{bmatrix} X & X & X & X \\ X & X & X & X \\ 0 & X & X & X \\ 0 & X & X & X \end{bmatrix}$$

Use to 64 bits ~ 8 bytes represent each coefficient

Store all numbers (even zeros)...

1000 × 1000	→	8 MB	10 min
10,000 × 10,000	→	800 MB	10,000 min
100,000 × 100,000	→	80 GB	"forever"

- Reordering



Node Reordering Can Reduce Fill-In

- Preserves Properties (symmetry, diagonal dominance)
- Equivalent to swapping rows and columns

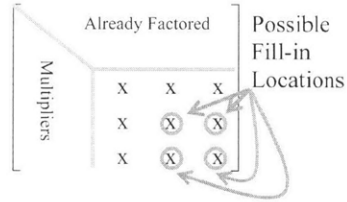
• Numeric Example

$$\begin{matrix} & V_1 & V_2 & V_3 \\ KCL_1 & \begin{bmatrix} 2 & -1 & -1 \end{bmatrix} \\ KCL_2 & \begin{bmatrix} -1 & 3 & 0 \end{bmatrix} \\ KCL_3 & \begin{bmatrix} -1 & 0 & 1 \end{bmatrix} \end{matrix}$$

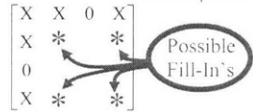
→ Swap KCL_1 and KCL_2 → Swap variable order V_1 and V_2

$$\begin{matrix} & V_1 & V_2 & V_3 \\ KCL_2 & \begin{bmatrix} -1 & 3 & 0 \end{bmatrix} \\ KCL_1 & \begin{bmatrix} 2 & -1 & -1 \end{bmatrix} \\ KCL_3 & \begin{bmatrix} -1 & 0 & 1 \end{bmatrix} \end{matrix} \quad \begin{matrix} & V_2 & V_1 & V_3 \\ KCL_2 & \begin{bmatrix} 3 & -1 & 0 \end{bmatrix} \\ KCL_1 & \begin{bmatrix} -1 & 2 & -1 \end{bmatrix} \\ KCL_3 & \begin{bmatrix} 0 & -1 & 1 \end{bmatrix} \end{matrix}$$

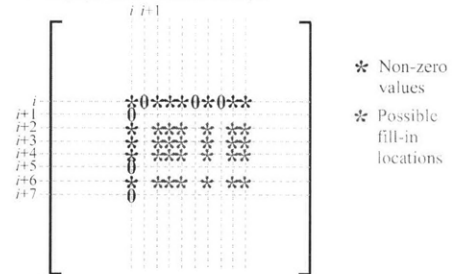
Where can fill-in occur?



Fill-in Estimate = (Non zeros in unfactored part of Row i) · (Non zeros in unfactored part of Col j) = Markowitz product

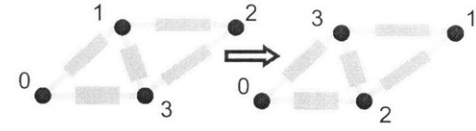


Fill-ins propagate down and to the right



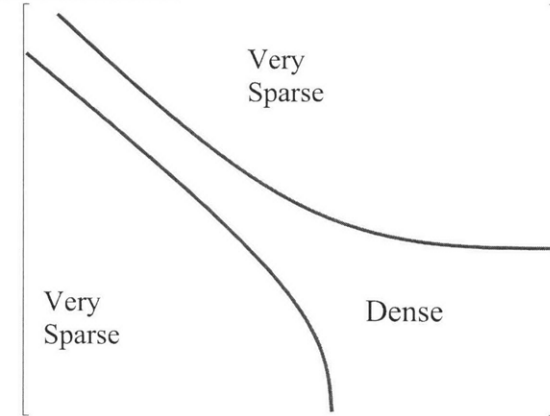
Markowitz Reordering

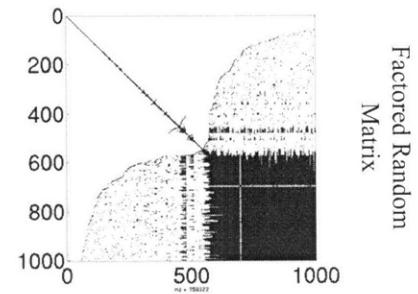
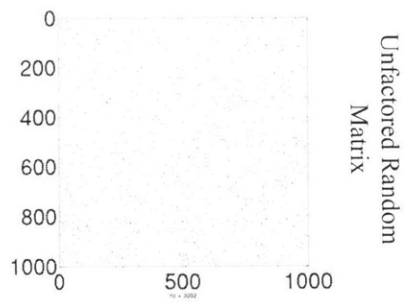
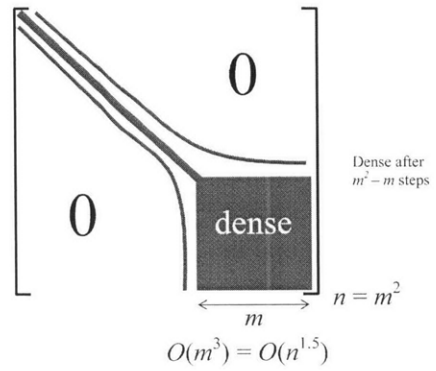
- For $i = 1$ to n
 - Find diagonal $j \geq i$ with min Markowitz Product
 - Swap rows $j \neq i$ and columns $j \neq i$
 - Factor the new row i and determine fill-ins
 - End
- Greedy Algorithm!
Why only try diagonals?
- Corresponds to node reordering in Nodal formulation



- Reduces search cost
- Preserves Matrix Properties
 - Diagonal Dominance
 - Symmetry

○ Pattern of a Filled-in Matrix





INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 6.A.

Solution of Sparse Linear Systems

TODAY'S OUTLINE:

- ❖ Solution of Sparse Linear Systems
 - General Sparse Factorization
 - Graph Based Approach
 - Sparse Matrix Data Structures
 - Scattering

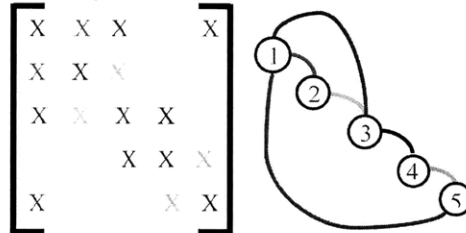
SOLUTION OF SPARSE LINEAR SYSTEMS

General Sparse Factorization

- ❖ Graph Based Approach
 - o Construction

Structurally Symmetric Matrices and Graphs

Note that *structurally* symmetric does not imply that the values of the matrix are symmetric (i.e., $a_{12} \neq a_{21}$)



- One Node Per Matrix Row
- One Edge Per Off-Diagonal Pair

Can one apply these graph-based techniques to the following matrix?

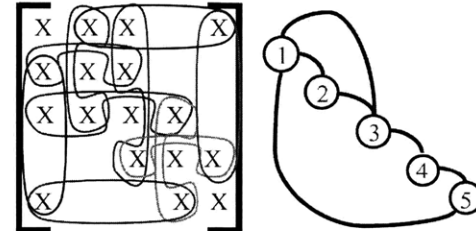
$$\begin{bmatrix} 1 & 3 & 0 & 0 \\ 0 & 1 & 2 & 0 \\ 0 & 2 & 1 & 5 \\ 0 & 0 & 5 & 1 \end{bmatrix}$$

This matrix is not "structurally symmetric" in that there is a zero in a_{21} and a non-zero value in a_{12} .

Can still use this approach, just treat it as if there is a nonzero in the a_{21} place and use the graphs to do the analysis; there will be some efficiency loss, but the methods will still work.

Thus, this technique may be applied to "mildly" structurally symmetric matrices by assuming there is structural symmetry with some loss of efficiency in treating some of the zeros as if they are non-zeros.

o Markowitz Products

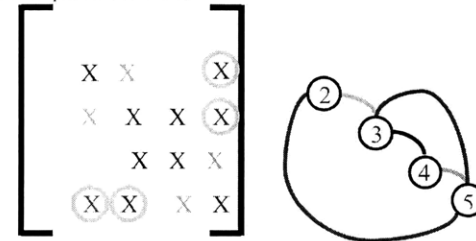


Markowitz Products = (Node Degree)²

$M_{11} \Rightarrow 3 \times 3 = 9$	(degree 1) ² = 3 ² = 9
$M_{22} \Rightarrow 2 \times 2 = 4$	(degree 2) ² = 2 ² = 4
$M_{33} \Rightarrow 3 \times 3 = 9$	(degree 3) ² = 3 ² = 9
$M_{44} \Rightarrow 2 \times 2 = 4$	(degree 4) ² = 2 ² = 4
$M_{55} \Rightarrow 2 \times 2 = 4$	(degree 5) ² = 2 ² = 4

o Factorization

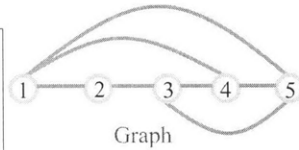
One Step of LU Factorization



- Delete the node associated with pivot row
- "Tie together" the graph edges

When node 1 is removed, in the matrix, the non-zero entries in the first row create fill-ins that connect up the other nodes that can be seen in the graph.

o Example

$$\begin{bmatrix} X & X & & X & X \\ X & X & X & & \\ & X & X & X & X \\ X & & X & X & X \\ X & & X & X & X \end{bmatrix}$$


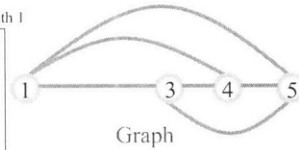
Graph

Col Row

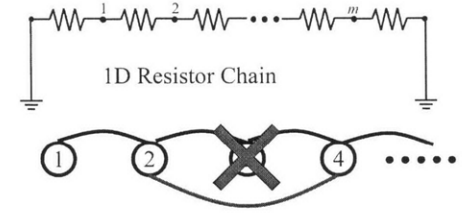
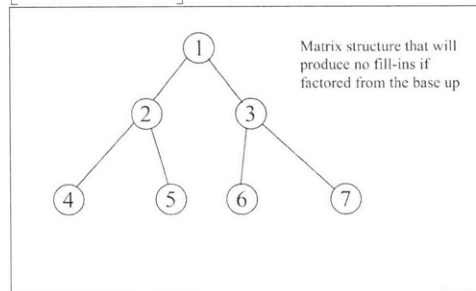
- (1) 3 3 = 9
- (2) 2 2 = 4
- (3) 3 3 = 9
- (4) 3 3 = 9
- (5) 3 3 = 9

Markowitz Products
(= node degree)

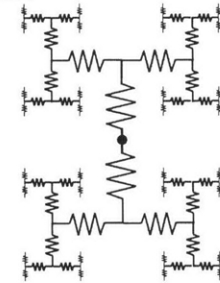
Swap (row & column) 2 with 1

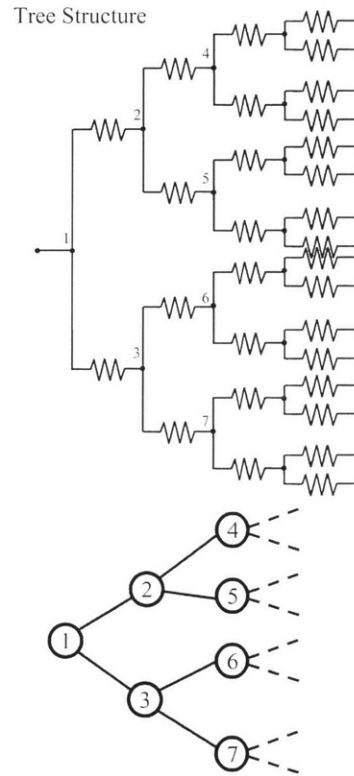
$$\begin{bmatrix} X & X & X \\ X & X & X & X & X \\ X & X & X & X & X \\ X & X & X & X \\ X & X & X & X \end{bmatrix}$$


Graph

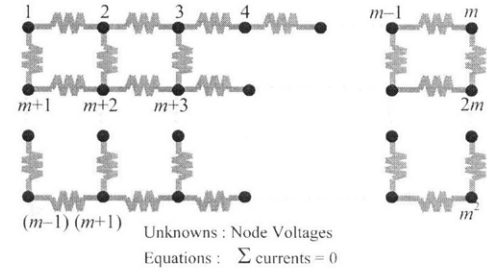


Clock Tree



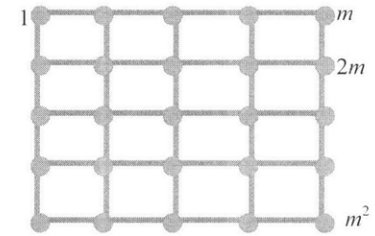


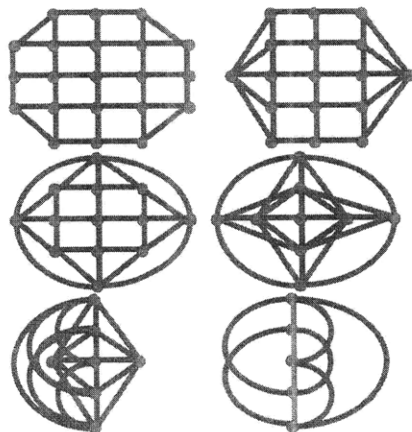
Resistor Grid Example



Grid Example

How long does it take to factor an $m \times m$ grid?





Suppose the center column is eliminated last

Factoring an $m \times m$ grid.

Dense LU: cost is $O(m^2)^3 = O(m^6)$

Sparse LU: each row has approximately 5 nonzeros
matrix has approximately $5m^2$ nonzeros
cost is at least $O(m^3)$

Separator is of size $2m$

Factoring: cost is $O(2m)^3 = O(8m^3)$ {dense matrix}

Subgraph: size is $m/2$

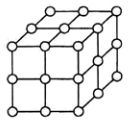
Create 4 separators: size is $2 \times m/2 = m$

Cost is $O(m^3)$ per separator $\rightarrow O(4m^3)$

As one continues to subdivide, the cost keeps halving.

Adding together these costs, the final cost will be roughly $O(m^3)$

3-D Example.



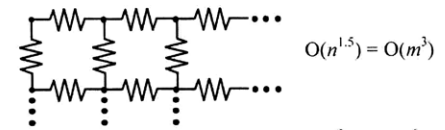
$3 \times 3 \times 3$ grid $\rightarrow 27$ nodes
in general: $m \times m \times m$ grid $\rightarrow m^3$ nodes
What does the separator look like for this grid?
 $m \times m$ nodes
 m^2 nodes
Factoring cost = $O(m^6)$

Structure & Dimension

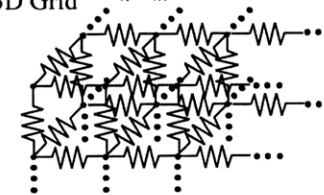
1D Row $n = m$ tridiagonal



2D Grid $n = m^2$ 5-diagonals

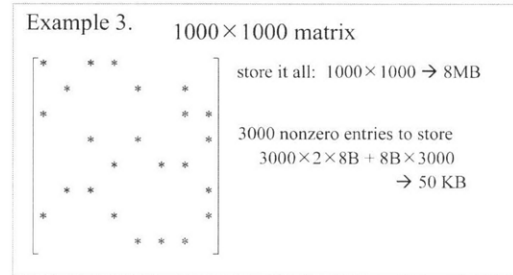
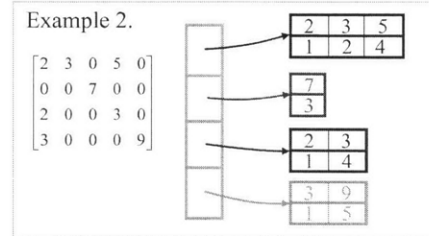
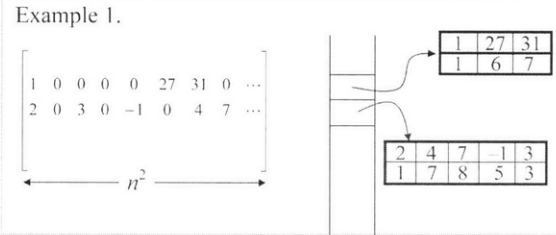


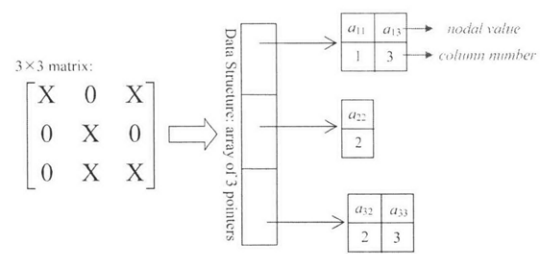
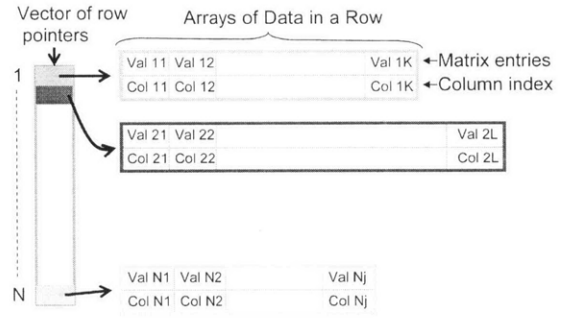
3D Grid $n = m^3$ $O(n^2) = O(m^6)$



- o What should you pivot for?
Growth control?
or to Avoid Fill-ins?
- A.) LU factorization applied to a strictly diagonally dominant matrix will never produce a zero pivot.
- B.) The matrix entries produced by LU factorization applied to a strictly diagonally dominant matrix will never increase by more than a factor of $2^{(n-1)}$, [which is the best you can do by pivoting for growth control]
- Bottom Line:
If your matrix is strictly diagonally dominant no need for numerical pivot for growth control so just pivot for sparsity control!
- o Sparse Factorization Approach
 1. Assume matrix requires NO numerical pivoting.
Diagonally dominant or symmetric positive definite.
 2. Use graphs to determine matrix ordering.
Many graph manipulation tricks used.
 3. Form data structures for storing filled-in matrix.
Lots of additional nonzeros added.
 4. Put numerical values in data structure and factor
Computation must be organized carefully!

❖ Sparse Matrix Data Structures





- Why store sparse matrix information in a data structure array?
- Too much storage space storing all the zero entries
 - Avoid floating point computation on all the zeros (minimal in comparison to the memory cost)
 - Memory reference cost

Eliminating Source Row *i* from Target Row *j*:

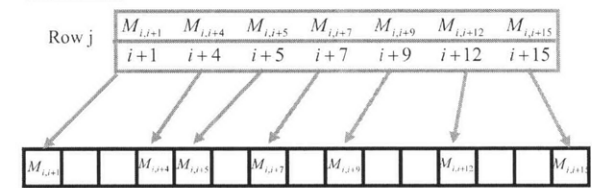
Row <i>i</i>	$M_{i,i+1}$	$M_{i,i+7}$	$M_{i,i+15}$				
	$i+1$	$i+7$	$i+15$				
Row <i>j</i>	$M_{j,i+1}$	$M_{j,i+4}$	$M_{j,i+5}$	$M_{j,i+7}$	$M_{j,i+9}$	$M_{j,i+12}$	$M_{j,i+15}$
	$i+1$	$i+4$	$i+5$	$i+7$	$i+9$	$i+12$	$i+15$

Must read in **all** the row *j* entries to find the three that match row *i*.

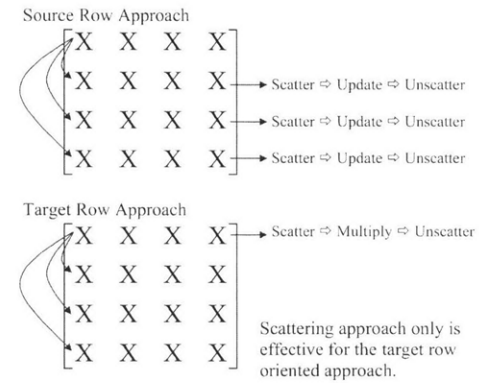
Rows	Ops	Misses	
Res	300	904,387	248,967
RAM	2,806	1,017,289	3,817,587
Grid	4,356	3,180,726	3,597,746

More misses than operations!
Every miss is an unneeded memory reference!

Scattering



- 1) Read all the elements in Row *j*, and scatter them in an *n*-length vector
- 2) Access only the needed elements using array indexing!



- ❖ Summary of Sparse Systems
 - Sparse Matrices
 - Struts, resistor grids, 3-D heat flow $\rightarrow O(N)$ nonzeros
 - Tridiagonal Matrix Factorization
 - Factor in $O(N)$ operations
 - General Sparse Factorization
 - Markowitz Reordering to minimize fill
 - Graph Based Approach
 - Factorization and Fill-in
 - Useful for estimating Sparse GE complexity
 - Sparse Data Structures
 - Scattering

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 6.B.

QR Factorization

TODAY'S OUTLINE:

- ❖ Singular Problems
- ❖ Projection Formulas
- ❖ Modified Gram-Schmidt Algorithm

QR FACTORIZATION

Singular Problems - LU Factorization Fails

- ❖ Struts Example

$f_x^* + f_{L_x} = 0$
 $f_y^* + f_{L_y} = 0$

$u_x = x - x_0$
 $u_y = y - y_0$

$$\begin{bmatrix} 0 & 0 \\ 0 & -\varepsilon \end{bmatrix} \begin{bmatrix} u_x \\ u_y \end{bmatrix} = \begin{bmatrix} 0 \\ +3 \end{bmatrix}$$

$$\begin{bmatrix} \vec{M}_1 & \vec{M}_2 \end{bmatrix} \vec{b}$$

load force $\vec{f}_L = (0, -3)$

$\vec{b} \in \text{span}\{\vec{M}_1, \vec{M}_2\} \Rightarrow$ infinite solutions

The resulting nodal matrix is SINGULAR:

- LU decomposition fails
- But a solution exists! Actually, many....

$$\begin{bmatrix} +1 & -1 & 0 \\ -1 & +1 & 0 \\ 0 & 0 & +1 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix} = \begin{bmatrix} -1 \\ +1 \\ -1 \end{bmatrix}$$

$$\begin{bmatrix} \vec{M}_1 & \vec{M}_2 & \vec{M}_3 \end{bmatrix} \vec{b}$$

$\vec{b} \in \text{span}\{\vec{M}_1, \vec{M}_2, \vec{M}_3\} \Rightarrow$ solution exists plane

\vec{M}_1 and \vec{M}_2 are L.D. \Rightarrow infinite solutions

e.g. $\vec{b} = 0 \cdot \vec{M}_1 + 1 \cdot \vec{M}_2 - 1 \cdot \vec{M}_3$

$\vec{b} = -1 \cdot \vec{M}_1 + 0 \cdot \vec{M}_2 - 1 \cdot \vec{M}_3$

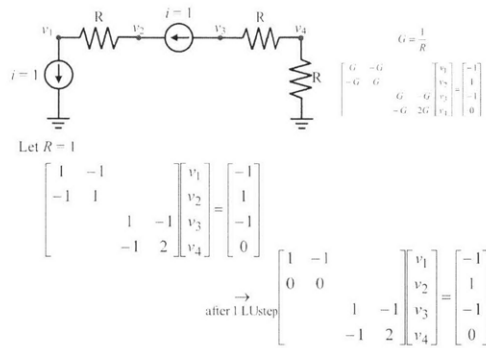
LU will fail

$$\begin{bmatrix} +1 & -1 & 0 \\ -1 & +1 & 0 \\ 0 & 0 & +1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

No pivots available

projection of \vec{b} onto $\text{span}\{\vec{M}_1, \vec{M}_2, \vec{M}_3\}$

$\vec{b} \notin \text{span}\{\vec{M}_1, \vec{M}_2, \vec{M}_3\}$



Recall weighted sum of columns view of systems of equations

$$\begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ M_1 & M_2 & \dots & M_N \\ \downarrow & \downarrow & \dots & \downarrow \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_N \end{bmatrix} \Rightarrow x_1 M_1 + x_2 M_2 + \dots + x_N M_N = \vec{b}$$

M is singular but **b** is in the span of the columns of **M** so there is a solution, actually lots of them. How do we find them?

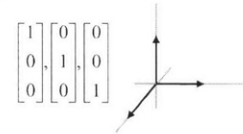
❖ Systems of Linear Equations – Summary Table

	$b \in \text{range}\{M\}$	$b \notin \text{range}\{M\}$
L.I. columns	Solution exists and is unique Use LU	No solutions Find the "closest"
L.D. columns	Infinite solutions exist Find one ... or all	No solutions Find the "closest"

Projection Formulas – Orthogonal Projection

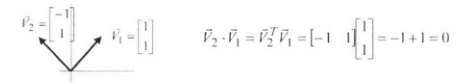
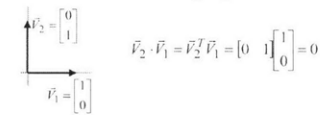
3-D Orthogonal Vectors

Example of orthogonal vectors



❖ Orthonormal Q Picture

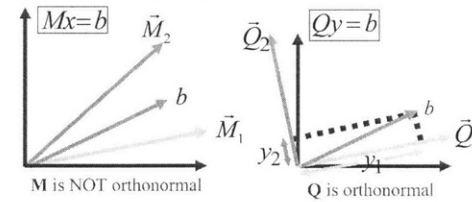
$$\vec{Q}_i \cdot \vec{Q}_j = \vec{Q}_i^T \vec{Q}_j = \begin{bmatrix} \vec{Q}_i^T \\ \vec{Q}_j \end{bmatrix}$$



Definition of orthonormal matrix **Q**:

$$\vec{Q}_i \cdot \vec{Q}_j = 0 \text{ if } i \neq j \text{ and } \vec{Q}_i \cdot \vec{Q}_i = 1$$

Picture for the two-dimensional case



$$\mathbf{Q} = \begin{bmatrix} \uparrow & \uparrow \\ \bar{Q}_1 & \bar{Q}_2 \\ \downarrow & \downarrow \end{bmatrix} \quad \bar{Q}_1 \cdot \bar{Q}_2 = \bar{Q}_1^T \bar{Q}_2 \quad y_1 = \bar{Q}_1^T \bar{b} \\
 \bar{y} = \mathbf{Q}^T \bar{b} \quad y_2 = \bar{Q}_2^T \bar{b}$$

❖ Projection Formula
 $\mathbf{Q}\bar{y} = \bar{b}$ solving an orthonormal system is easy...

$$y_1 \bar{Q}_1 + y_2 \bar{Q}_2 + \dots + y_N \bar{Q}_N = \bar{b}$$

Example.

$$Qy = b \Rightarrow y_1 \begin{bmatrix} \bar{Q}_1 \\ \vdots \\ \bar{Q}_1 \end{bmatrix} + y_2 \begin{bmatrix} \bar{Q}_2 \\ \vdots \\ \bar{Q}_2 \end{bmatrix} = \begin{bmatrix} \bar{b} \\ \vdots \\ \bar{b} \end{bmatrix}$$

Multiplying the weighted columns equation by i^{th} column:

$$\bar{Q}_i \cdot (y_1 \bar{Q}_1 + y_2 \bar{Q}_2 + \dots + y_N \bar{Q}_N) = \bar{Q}_i \cdot \bar{b}$$

Example.

$$\bar{Q}_1 \cdot \left(y_1 \begin{bmatrix} \bar{Q}_1 \\ \vdots \\ \bar{Q}_1 \end{bmatrix} + y_2 \begin{bmatrix} \bar{Q}_2 \\ \vdots \\ \bar{Q}_2 \end{bmatrix} = \begin{bmatrix} \bar{b} \\ \vdots \\ \bar{b} \end{bmatrix} \right) \Rightarrow y_1 \underbrace{(\bar{Q}_1 \cdot \bar{Q}_1)}_1 + y_2 \underbrace{(\bar{Q}_1 \cdot \bar{Q}_2)}_0 = \bar{Q}_1 \cdot \bar{b}$$

$$\bar{Q}_2 \cdot \left(y_1 \begin{bmatrix} \bar{Q}_1 \\ \vdots \\ \bar{Q}_1 \end{bmatrix} + y_2 \begin{bmatrix} \bar{Q}_2 \\ \vdots \\ \bar{Q}_2 \end{bmatrix} = \begin{bmatrix} \bar{b} \\ \vdots \\ \bar{b} \end{bmatrix} \right) \Rightarrow y_1 \underbrace{(\bar{Q}_2 \cdot \bar{Q}_1)}_0 + y_2 \underbrace{(\bar{Q}_2 \cdot \bar{Q}_2)}_1 = \bar{Q}_2 \cdot \bar{b}$$

Simplifying using orthonormality we get the projection formula:

$$y_i = \bar{Q}_i \cdot \bar{b}$$

Example.

$$\begin{aligned}
 y_1 &= \bar{Q}_1 \cdot \bar{b} = \bar{Q}_1^T \bar{b} \\
 y_2 &= \bar{Q}_2 \cdot \bar{b} = \bar{Q}_2^T \bar{b} \\
 y_3 &= \bar{Q}_3 \cdot \bar{b} = \bar{Q}_3^T \bar{b} \\
 &\vdots \\
 y_N &= \bar{Q}_N \cdot \bar{b} = \bar{Q}_N^T \bar{b}
 \end{aligned}
 \quad \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} \left[\leftarrow \bar{Q}_1^T \rightarrow \right] \\ \left[\leftarrow \bar{Q}_2^T \rightarrow \right] \\ \left[\leftarrow \bar{Q}_3^T \rightarrow \right] \\ \vdots \\ \left[\leftarrow \bar{Q}_N^T \rightarrow \right] \end{bmatrix} \begin{bmatrix} \bar{b} \\ \vdots \\ \bar{b} \end{bmatrix}$$

Know $\mathbf{Q}\bar{y} = \bar{b}$; have shown that $\bar{y} = \mathbf{Q}^T \bar{b}$

$$\mathbf{Q}^T \mathbf{Q} = \begin{bmatrix} 1 & 0 & 0 & \dots \\ \bar{Q}_1^T \bar{Q}_1 & \bar{Q}_1^T \bar{Q}_2 & \bar{Q}_1^T \bar{Q}_3 & \dots \\ \bar{Q}_2^T \bar{Q}_1 & \bar{Q}_2^T \bar{Q}_2 & \bar{Q}_2^T \bar{Q}_3 & \dots \\ 0 & 1 & 0 & \dots \\ & & \ddots & \ddots \\ & & & \bar{Q}_N^T \bar{Q}_N \end{bmatrix} = I$$

$$\mathbf{Q}^{-1} \mathbf{Q} = \mathbf{Q}^T \mathbf{Q} = I \Rightarrow \mathbf{Q}^T = \mathbf{Q}^{-1}$$

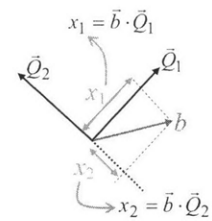
❖ QR Algorithm Key Idea

$$\begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ \bar{M}_1 & \bar{M}_2 & \dots & \bar{M}_N \\ \downarrow & \downarrow & & \downarrow \\ & & & x_N \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_N \end{bmatrix} \Rightarrow \begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ \bar{Q}_1 & \bar{Q}_2 & \dots & \bar{Q}_N \\ \downarrow & \downarrow & & \downarrow \\ & & & y_N \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_N \end{bmatrix}$$

original matrix matrix with orthonormal columns

If we have a set of orthonormal columns that span the same subspace as the columns of \mathbf{M} , then solving the system is very easy:

$$\mathbf{Q}\bar{y} = \bar{b} \Rightarrow \bar{y} = \mathbf{Q}^T \bar{b}$$



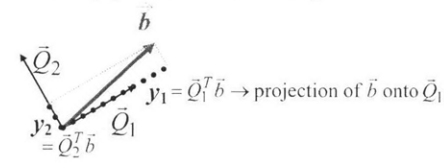
$$\begin{bmatrix} \bar{Q}_1 & \bar{Q}_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = x_1 \bar{Q}_1 + x_2 \bar{Q}_2 = \bar{b}$$

$$\begin{aligned}
 \bar{Q}_1 \cdot \bar{Q}_1 &= 1 \\
 \bar{Q}_2 \cdot \bar{Q}_2 &= 1 \\
 \bar{Q}_1 \cdot \bar{Q}_2 &= \bar{Q}_2 \cdot \bar{Q}_1 = 0
 \end{aligned}$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \bar{Q}_1 \\ \bar{Q}_2 \end{bmatrix} \begin{bmatrix} \bar{b} \end{bmatrix} = \begin{bmatrix} \bar{Q}_1^T \bar{b} \\ \bar{Q}_2^T \bar{b} \end{bmatrix}$$

$$\begin{bmatrix} \uparrow \\ \bar{y} \\ \downarrow \end{bmatrix} = \begin{bmatrix} \leftarrow \bar{Q}_1^T \rightarrow \\ \leftarrow \bar{Q}_2^T \rightarrow \\ \vdots \end{bmatrix} \begin{bmatrix} \uparrow \\ \bar{b} \\ \downarrow \end{bmatrix} = \begin{bmatrix} \bar{Q}_1^T \bar{b} \\ \bar{Q}_2^T \bar{b} \\ \vdots \end{bmatrix} \rightarrow \text{projection of } \bar{b} \text{ onto } \bar{Q}_1$$

$y_i = \bar{Q}_i^T \bar{b}$
Result is the projection of \bar{b} onto the columns of \bar{Q}



But how do we perform the conversion from \mathbf{M} to \mathbf{Q} ?

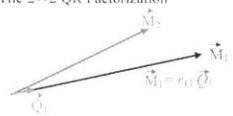
Example.

$$\mathbf{M} = \begin{bmatrix} \uparrow & \uparrow \\ M_1 & M_2 \\ \downarrow & \downarrow \end{bmatrix} = \begin{bmatrix} +1 & 0 \\ -1 & 0 \\ 0 & 0 \end{bmatrix}$$

$$r_{11} = \sqrt{1^2 + (-1)^2} = \sqrt{2}$$

$$\bar{Q}_1 = \frac{\bar{M}_1}{r_{11}} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \\ 0 \end{bmatrix}$$

❖ The 2×2 QR Factorization



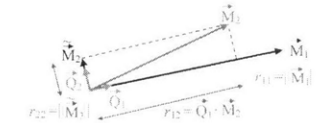
$$r_{11} = \|\bar{M}_1\| = \sqrt{\bar{M}_1 \cdot \bar{M}_1} = \sqrt{\bar{M}_1^T \bar{M}_1}$$

Getting the first orthonormal vector is very easy:

$$\bar{Q}_1 = \frac{\bar{M}_1}{r_{11}} = \frac{\bar{M}_1}{\|\bar{M}_1\|}$$

Given \bar{M}_1, \bar{M}_2 , find $\bar{Q}_2 = \bar{M}_2 - r_{12} \bar{M}_1$ so that

$$\bar{M}_1 \cdot \bar{Q}_2 = \bar{M}_1 \cdot (\bar{M}_2 - r_{12} \bar{M}_1) = 0 \Rightarrow r_{12} = \frac{\bar{M}_1 \cdot \bar{M}_2}{\bar{M}_1 \cdot \bar{M}_1}$$



Using the Projection formula

$$\bar{M}_1 = r_{11} \bar{Q}_1$$

$$\bar{M}_2 = r_{12} \bar{Q}_1 + r_{22} \bar{Q}_2 \Rightarrow \begin{cases} \bar{M}_2 = \bar{M}_2 - r_{12} \bar{Q}_1 \\ \bar{Q}_2 = \frac{\bar{M}_2}{r_{22}} \end{cases}$$

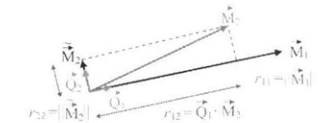
$$r_{12} = \bar{M}_1 \cdot \bar{M}_2 = \bar{Q}_1 \cdot \bar{M}_2 = \bar{Q}_1 \cdot (r_{12} \bar{Q}_1 + r_{22} \bar{Q}_2) = r_{12} \underbrace{\bar{Q}_1 \cdot \bar{Q}_1}_1 + r_{22} \underbrace{\bar{Q}_1 \cdot \bar{Q}_2}_0$$

Formulas simplify if we normalize

$$\bar{Q}_1 = \frac{\bar{M}_1}{\sqrt{\bar{M}_1 \cdot \bar{M}_1}} = \frac{1}{r_{11}} \bar{M}_1 \Rightarrow \bar{Q}_1 \cdot \bar{Q}_1 = 1$$

Now find $\bar{Q}_2 = \bar{M}_2 - r_{12} \bar{Q}_1$ so that $\bar{Q}_2 \cdot \bar{Q}_1 = 0 \Rightarrow r_{12} = \bar{Q}_1 \cdot \bar{M}_2$

$$\text{Finally } \bar{Q}_2 = \frac{1}{\sqrt{\bar{Q}_2 \cdot \bar{Q}_2}} \bar{Q}_2 = \frac{1}{r_{22}} \bar{Q}_2$$



$$\bar{M}_1 = r_{11} \bar{Q}_1 \quad \bar{M}_2 = r_{11} \bar{Q}_1 + r_{22} \bar{Q}_2$$

$$\bar{M}_2 = r_{12} \bar{Q}_1 + r_{22} \bar{Q}_2 \quad \bar{M}_2 = r_{12} \bar{Q}_1 + r_{22} \bar{Q}_2$$

Another way to write this is: \mathbf{M} is factored into \mathbf{Q} and \mathbf{R}

$$\underbrace{\begin{bmatrix} \bar{M}_1 & \bar{M}_2 \end{bmatrix}}_{\mathbf{M}} = \underbrace{\begin{bmatrix} \bar{Q}_1 & \bar{Q}_2 \end{bmatrix}}_{\mathbf{Q}} \underbrace{\begin{bmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{bmatrix}}_{\mathbf{R}}$$

Since \mathbf{Mx} should equal \mathbf{Qy} , we can relate x to y

$$\begin{bmatrix} \bar{M}_1 & \bar{M}_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = x_1 \bar{M}_1 + x_2 \bar{M}_2 = \begin{bmatrix} \bar{Q}_1 & \bar{Q}_2 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = y_1 \bar{Q}_1 + y_2 \bar{Q}_2$$

$$\begin{bmatrix} \bar{M}_1 \\ \bar{M}_2 \end{bmatrix} = \begin{bmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

Three Step Solve Procedure:

- Factor $\mathbf{M} = \mathbf{QR}$, \mathbf{Q} orthonormal, \mathbf{R} upper triangular

$$\mathbf{M}\vec{x} = \vec{b} \quad x_1 \bar{M}_1 + x_2 \bar{M}_2 = \vec{b}$$

Decompose $\mathbf{Q} \cdot \mathbf{R}\vec{x} = \vec{b}$

- Solve $\mathbf{Q}\vec{v} = \vec{b}$. (Very easy: $\vec{v} = \mathbf{Q}^T \vec{b}$)

- Backsolve the triangular system: $\mathbf{R}\vec{x} = \vec{v}$.

❖ The 3×3 QR Factorization

Use previously orthonormalized vectors

$$\begin{bmatrix} \bar{M}_1 & \bar{M}_2 & \bar{M}_3 \\ \uparrow & \uparrow & \uparrow \\ \bar{M}_1 & \bar{M}_2 & \bar{M}_3 \\ \downarrow & \downarrow & \downarrow \end{bmatrix} \Rightarrow \begin{bmatrix} \bar{M}_1 & \bar{M}_2 - r_{12}\bar{Q}_1 & \bar{M}_3 - r_{13}\bar{Q}_1 - r_{23}\bar{Q}_2 \\ \uparrow & \uparrow & \uparrow \\ \bar{M}_1 & \bar{M}_2 - r_{12}\bar{Q}_1 & \bar{M}_3 - r_{13}\bar{Q}_1 - r_{23}\bar{Q}_2 \\ \downarrow & \downarrow & \downarrow \end{bmatrix}$$

To ensure the third column is orthogonal

$$\bar{Q}_1 \cdot (\bar{M}_3 - \bar{Q}_1 r_{13} - \bar{Q}_2 r_{23}) = 0 \Rightarrow r_{13} = \bar{Q}_1 \cdot \bar{M}_3$$

$$\bar{Q}_2 \cdot (\bar{M}_3 - \bar{Q}_1 r_{13} - \bar{Q}_2 r_{23}) = 0 \Rightarrow r_{23} = \bar{Q}_2 \cdot \bar{M}_3$$

$$\bar{M}_1 \cdot (\bar{M}_3 - r_{13}\bar{M}_1 - r_{23}\bar{M}_2) = 0 \Rightarrow \begin{bmatrix} \bar{M}_1 \cdot \bar{M}_1 & \bar{M}_1 \cdot \bar{M}_2 \\ \bar{M}_2 \cdot \bar{M}_1 & \bar{M}_2 \cdot \bar{M}_2 \end{bmatrix} \begin{bmatrix} r_{13} \\ r_{23} \end{bmatrix} = \begin{bmatrix} \bar{M}_1 \cdot \bar{M}_3 \\ \bar{M}_2 \cdot \bar{M}_3 \end{bmatrix}$$

$$\bar{M}_2 \cdot (\bar{M}_3 - r_{13}\bar{M}_1 - r_{23}\bar{M}_2) = 0$$

To Orthogonalize the N^{th} Vector:

$$\begin{bmatrix} \bar{M}_1 \cdot \bar{M}_1 & \dots & \bar{M}_1 \cdot \bar{M}_{N-1} \\ \vdots & \ddots & \vdots \\ \bar{M}_{N-1} \cdot \bar{M}_1 & \dots & \bar{M}_{N-1} \cdot \bar{M}_{N-1} \end{bmatrix} \begin{bmatrix} r_{1,N} \\ \vdots \\ r_{N-1,N} \end{bmatrix} = \begin{bmatrix} \bar{M}_1 \cdot \bar{M}_N \\ \vdots \\ \bar{M}_{N-1} \cdot \bar{M}_N \end{bmatrix}$$

N^2 inner products or N^3 work.

$$\begin{aligned} \bar{M}_3 &= r_{13}\bar{Q}_1 + r_{23}\bar{Q}_2 + r_{33}\bar{Q}_3 & r_{13} &= \bar{Q}_1 \cdot \bar{M}_3 \\ \bar{M}_3 &= \bar{M}_3 - (\bar{Q}_1 \cdot \bar{M}_3)\bar{Q}_1 - (\bar{Q}_2 \cdot \bar{M}_3)\bar{Q}_2 & r_{23} &= \bar{Q}_2 \cdot \bar{M}_3 \\ \Rightarrow \bar{M}_3 &= \bar{M}_3 - r_{13}\bar{Q}_1 - r_{23}\bar{Q}_2 & r_{33} &= \bar{Q}_3 \cdot \bar{M}_3 = \|\bar{M}_3\| \\ \bar{Q}_3 &= \frac{\bar{M}_3}{\|\bar{M}_3\|} \end{aligned}$$

$$\begin{aligned} \bar{M}_1 & \quad \bar{Q}_1 \leftarrow \frac{\bar{M}_1}{\|\bar{M}_1\|} \leftarrow r_{11} \\ \bar{M}_2 & \quad \bar{M}_2 \leftarrow \bar{M}_2 - \frac{(\bar{Q}_1 \cdot \bar{M}_2)\bar{Q}_1}{r_{12}} \quad \bar{Q}_2 \leftarrow \frac{\bar{M}_2}{\|\bar{M}_2\|} \leftarrow r_{22} \\ \bar{M}_3 & \quad \bar{M}_3 \leftarrow \bar{M}_3 - \frac{(\bar{Q}_1 \cdot \bar{M}_3)\bar{Q}_1}{r_{13}} - \frac{(\bar{Q}_2 \cdot \bar{M}_3)\bar{Q}_2}{r_{23}} \quad \bar{Q}_3 \leftarrow \frac{\bar{M}_3}{\|\bar{M}_3\|} \leftarrow r_{33} \\ \bar{M}_4 & \quad \bar{M}_4 \leftarrow \bar{M}_4 - \frac{(\bar{Q}_1 \cdot \bar{M}_4)\bar{Q}_1}{r_{14}} - \frac{(\bar{Q}_2 \cdot \bar{M}_4)\bar{Q}_2}{r_{24}} - \frac{(\bar{Q}_3 \cdot \bar{M}_4)\bar{Q}_3}{r_{34}} \quad \bar{Q}_4 \leftarrow \frac{\bar{M}_4}{\|\bar{M}_4\|} \leftarrow r_{44} \\ & \quad \vdots \\ \bar{M}_N & \quad \bar{M}_N \leftarrow \bar{M}_N - \frac{(\bar{Q}_1 \cdot \bar{M}_N)\bar{Q}_1}{r_{1N}} - \frac{(\bar{Q}_2 \cdot \bar{M}_N)\bar{Q}_2}{r_{2N}} - \dots - \frac{(\bar{Q}_{N-1} \cdot \bar{M}_N)\bar{Q}_{N-1}}{r_{N-1,N}} \end{aligned}$$

$$\begin{bmatrix} \bar{M}_1 & \bar{M}_2 & \bar{M}_3 \end{bmatrix} = \begin{bmatrix} \bar{Q}_1 & \bar{M}_2 & \bar{M}_3 \\ \times & \times & \times \\ \bar{Q}_1 & \bar{M}_2 & \bar{M}_3 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ & r_{22} & r_{23} \\ & & r_{33} \end{bmatrix}$$

$$\begin{bmatrix} \bar{M}_1 & \bar{M}_2 & \bar{M}_3 \end{bmatrix} = \begin{bmatrix} \bar{Q}_1 & \bar{M}_2 & \bar{M}_3 \\ \bar{Q}_1 & \bar{M}_2 & \bar{M}_3 \\ \times & \times & \times \\ \bar{Q}_1 & \bar{M}_2 & \bar{M}_3 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ & r_{22} & r_{23} \\ & & r_{33} \end{bmatrix}$$

$$\begin{bmatrix} \bar{M}_1 & \bar{M}_2 & \bar{M}_3 \end{bmatrix} = \begin{bmatrix} \bar{Q}_1 & \bar{Q}_2 & \bar{Q}_3 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ & r_{22} & r_{23} \\ & & r_{33} \end{bmatrix}$$

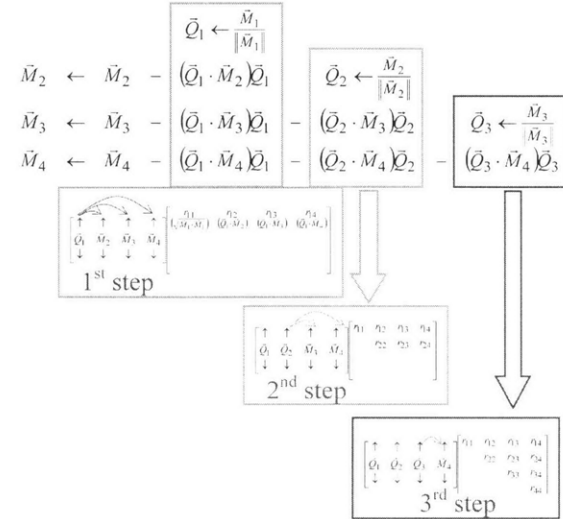
Modified Gram-Schmidt Algorithm
 ❖ "By Picture"

$$\bar{M}_j \leftarrow \bar{M}_j - r_{ij}\bar{Q}_i$$

$$\begin{bmatrix} \bar{Q}_1 & \bar{Q}_2 & \bar{Q}_3 & \bar{Q}_4 \\ \uparrow & \uparrow & \uparrow & \uparrow \\ r_{11} & r_{12} & r_{13} & r_{14} \\ \downarrow & \downarrow & \downarrow & \downarrow \\ & r_{22} & r_{23} & r_{24} \\ & & r_{33} & r_{34} \\ & & & r_{44} \end{bmatrix}$$

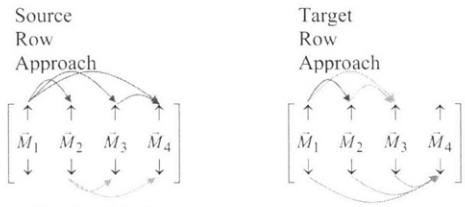
$$\bar{Q}_i = \frac{1}{r_{ii}} \bar{M}_i \quad r_{ii} = \sqrt{\bar{M}_i \cdot \bar{M}_i}$$

$$r_{ij} \leftarrow \bar{M}_j \cdot \bar{Q}_i$$



- ❖ Basic Algorithm
 - Source Row Approach
 - For $i = 1$ to N { "For each Source Column"
 - $r_{ii} = \sqrt{\bar{M}_i \cdot \bar{M}_i}$
 $\bar{Q}_i = \frac{1}{r_{ii}} \bar{M}_i$
 - } Normalize $\sum_{i=1}^N 2N = 2N^2$ operations
 - For $j = i + 1$ to N { "For each target column right of source"
 - $r_{ij} \leftarrow \bar{M}_j \cdot \bar{Q}_i$
 $\bar{M}_j \leftarrow \bar{M}_j - r_{ij}\bar{Q}_i$
 - } $\sum_{i=1}^N (N-i)2N = N^3$ operations

- o Target Row Approach
 - For $i = 1$ to N { "For each Target Column"
 - $\tilde{M}_i \leftarrow \mathbf{M}\tilde{e}_i$ "matrix-vector product"
 - For $j = 1$ to $i-1$ { "For each source column left of target"
 - $r_{ji} \leftarrow \tilde{Q}_j \cdot \tilde{M}_i$
 - $\tilde{M}_i \leftarrow \tilde{M}_i - r_{ji}\tilde{Q}_j$ } $\sum_{j=1}^{i-1} (N-i)2N \approx N^3$ operations
 - $r_{ii} = \sqrt{\tilde{M}_i \cdot \tilde{M}_i}$
 - $\tilde{Q}_i = \frac{1}{r_{ii}}\tilde{M}_i$ } Normalize $\sum_{i=1}^N 2N \approx 2N^2$ operations



❖ Three Step Solve Procedure

		Computational Complexity
Step 1)	Factor $\mathbf{M} = \mathbf{QR}$ (\mathbf{Q} orthonormal, \mathbf{R} upper triangular)	$O(N^3)$
Step 2)	Solve $\mathbf{Q}\tilde{y} = \tilde{b}$. (Very easy : $\tilde{y} = \mathbf{Q}^T \tilde{b}$)	$O(N^2)$
Step 3)	Backsolve the triangular system $\mathbf{R}\tilde{x} = \tilde{y}$	$O(N^2)$

- $\mathbf{M}\tilde{x} = \tilde{b}$
- Step 1) Factor $\mathbf{M} = \mathbf{QR}$
 $\mathbf{M}\tilde{x} = \tilde{b} \Leftrightarrow \mathbf{QR}\tilde{x} = \tilde{b}$
- Step 2) Solve for \tilde{y} : $\mathbf{Q}\tilde{y} = \tilde{b}$

$$y_i = \tilde{Q}_i \cdot \tilde{b} = \tilde{Q}_i^T \tilde{b} = \left[\leftarrow \tilde{Q}_i^T \rightarrow \right] \begin{bmatrix} \uparrow \\ \tilde{b} \\ \downarrow \end{bmatrix}$$

(Very easy : $\tilde{y} = \mathbf{Q}^T \tilde{b}$)

$$\tilde{y} = \mathbf{Q}^T \tilde{b} = \begin{bmatrix} \leftarrow \tilde{Q}_1 \rightarrow \\ \leftarrow \tilde{Q}_2 \rightarrow \\ \leftarrow \tilde{Q}_3 \rightarrow \\ \vdots \\ \leftarrow \tilde{Q}_n \rightarrow \end{bmatrix} \begin{bmatrix} \uparrow \\ \tilde{b} \\ \downarrow \end{bmatrix}$$

$\tilde{b} = y_1\tilde{Q}_1 + y_2\tilde{Q}_2 + y_3\tilde{Q}_3 + \dots + y_n\tilde{Q}_n$
 $\tilde{b} = x_1\tilde{M}_1 + x_2\tilde{M}_2 + x_3\tilde{M}_3 + \dots + x_n\tilde{M}_n$

Step 3)

Backsolve the triangular system $\mathbf{R}\tilde{x} = \tilde{y}$

$$\begin{bmatrix} \mathbf{R} & \begin{bmatrix} \rightarrow \\ \tilde{x} \\ \rightarrow \end{bmatrix} \end{bmatrix} = \begin{bmatrix} \rightarrow \\ \tilde{y} \\ \rightarrow \end{bmatrix}$$

INTRODUCTION TO NUMERICAL SIMULATION

LECTURES 7 & 8.

QR and Krylov-Subspace Matrix Solution Methods

TODAY'S OUTLINE:

- ❖ Minimization View of QR
 - Singular Matrix
 - Basic Minimization Approach
 - Orthogonalized Search Directions
 - QR and Length Minimization Produce Identical Results
- ❖ Arbitrary Subspace Algorithm
 - Orthogonalization of Search Directions
- ❖ Generalized Conjugate Residual Algorithm
 - Krylov-subspace
 - Simplification in the symmetric case
 - Leaky and insulating examples

QR FACTORIZATION

Example.

$$\mathbf{A} = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ -\frac{1}{2} & 1 & 0 \\ - & -\frac{1}{3} & 1 \end{bmatrix} \begin{bmatrix} 2 & -1 & 0 \\ 0 & \frac{3}{2} & -1 \\ 0 & 0 & \frac{4}{3} \end{bmatrix}$$

$\mathbf{L} \qquad \mathbf{U}$

$$\begin{bmatrix} -0.89 & -0.35 & 0.27 \\ 0.44 & -0.72 & 0.53 \\ 0 & 0.60 & 0.80 \end{bmatrix} \begin{bmatrix} -2.2 & 1.8 & -0.45 \\ 0 & -1.7 & 1.9 \\ 0 & 0 & 1.1 \end{bmatrix}$$

$\mathbf{Q} \qquad \mathbf{R}$

Matrix is Singular, column of Q is zero

- ❖ Zero Column
- If a column is zero

$$\vec{M}_i = \sum_{j=1}^{i+1} w_j \vec{M}_j$$

$\{\vec{M}_1, \dots, \vec{M}_i\}$ not linearly independent

What if a column becomes zero?

$$\begin{bmatrix} \uparrow & \uparrow & \uparrow & \uparrow \\ \vec{Q}_1 & 0 & \vec{M}_3 & \dots & \vec{M}_N \\ \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & \dots & r_{1N} \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}$$

$$\begin{bmatrix} \uparrow & \uparrow & \uparrow & \uparrow \\ \vec{Q}_1 & 0 & \vec{Q}_3 & \dots & \vec{Q}_N \\ \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & \dots & r_{1N} \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & r_{33} & \dots & r_{3N} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & r_{NN} \end{bmatrix}$$

Matrix MUST be Singular!!

1. Do not try to normalize the column.
2. Do not use the column as a source for orthogonalization.
3. Perform backward substitution as well as possible.

$$\mathbf{QR}\vec{x} = \vec{b} \Rightarrow \mathbf{R}\vec{x} = \mathbf{Q}^T \vec{b}$$

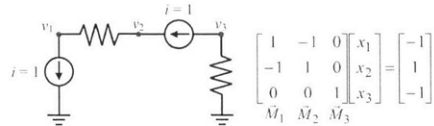
$$\Rightarrow \begin{bmatrix} r_{11} & r_{12} & r_{13} & \dots & r_{1N} \\ 0 & 0 & 0 & \dots & 0 \\ & & r_{33} & \dots & r_{3N} \\ & & & \ddots & \vdots \\ & & & & r_{NN} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} \leftarrow \vec{Q}_1 \rightarrow \\ 0 & \dots & 0 \\ \leftarrow \vec{Q}_3 \rightarrow \\ \vdots \\ \leftarrow \vec{Q}_N \rightarrow \end{bmatrix} \begin{bmatrix} \uparrow \\ \vec{b} \\ \downarrow \end{bmatrix}$$

QR Factorization

Problem $\mathbf{M}\vec{x} = \vec{b}$:

- (1) $\mathbf{Q} \cdot \mathbf{R}\vec{x} = \vec{b}$
- (2) Solve $\mathbf{Q}\vec{y} = \vec{b} \Rightarrow \vec{y} = \mathbf{Q}^T \vec{b}$
- (3) Solve $\mathbf{R}\vec{x} = \vec{y} \rightarrow \vec{x}$

❖ Singular Example

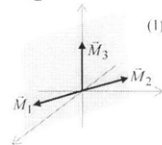


(1) Column \vec{M}_1

$$\eta_{11} = \sqrt{\vec{M}_1 \cdot \vec{M}_1} = \sqrt{2} \quad \vec{Q}_1 = \frac{1}{\eta_{11}} \vec{M}_1 = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \\ 0 \end{bmatrix}$$

$$\eta_{12} = \vec{Q}_1 \cdot \vec{M}_2 = -\frac{1}{\sqrt{2}} - \frac{1}{\sqrt{2}} = -\sqrt{2}$$

$$\eta_{13} = \vec{Q}_1 \cdot \vec{M}_3 = 0$$



$$\vec{M}_2 \leftarrow \vec{M}_2 - \eta_{12} \vec{Q}_1 = \begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix} - \begin{bmatrix} -1 \\ -1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\vec{M}_3 \leftarrow \vec{M}_3 - \eta_{13} \vec{Q}_1 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} \frac{1}{\sqrt{2}} & 0 & 0 \\ -\frac{1}{\sqrt{2}} & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \sqrt{2} & -\sqrt{2} & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$\vec{Q}_1 \quad \vec{M}_2 \quad \vec{M}_3$

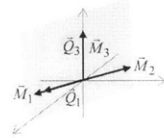
Column \vec{M}_2
 Since \vec{M}_2 is zero, don't need to do anything.
 Move on to the next step

$$\begin{bmatrix} \frac{1}{\sqrt{2}} & 0 & 0 \\ -\frac{1}{\sqrt{2}} & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \sqrt{2} & -\sqrt{2} & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$\vec{Q}_1 \quad \vec{Q}_2 \quad \vec{M}_3$

Column \vec{M}_3
 $r_{33} = \sqrt{\vec{M}_3 \cdot \vec{M}_3} = 1 \quad \vec{Q}_3 = \frac{1}{r_{33}} \vec{M}_3$

$$\underbrace{\begin{bmatrix} \frac{1}{\sqrt{2}} & 0 & 0 \\ -\frac{1}{\sqrt{2}} & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}}_{\mathbf{Q}} \underbrace{\begin{bmatrix} \sqrt{2} & -\sqrt{2} & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}}_{\mathbf{R}}$$

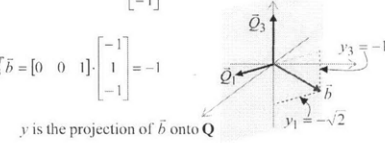


(2) $\vec{y} = \mathbf{Q}^T \vec{b}$

$$y_1 = \vec{Q}_1^T \vec{b} = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \end{bmatrix} \cdot \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} = -\frac{1}{\sqrt{2}} - \frac{1}{\sqrt{2}} = -\sqrt{2}$$

$$y_2 = 0$$

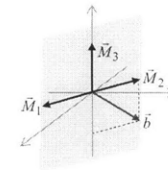
$$y_3 = \vec{Q}_3^T \vec{b} = \begin{bmatrix} 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} = -1$$



(3) $\mathbf{R}\vec{x} = \vec{y}$

$$\begin{bmatrix} \sqrt{2} & -\sqrt{2} & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -\sqrt{2} \\ 0 \\ -1 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -1+a \\ a \\ -1 \end{bmatrix}$$



$\vec{b} \in \text{span}\{\mathbf{M}\} \rightarrow$

The resulting nodal matrix is SINGULAR, but a solution exists!
 \rightarrow infinite number of solutions.

If $\tilde{b} \notin \text{span}\{\mathbf{M}\} \rightarrow$ there is no solution

$$\begin{bmatrix} 1 & -1 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 2 \\ -1 \end{bmatrix}$$

(1) Step (1) stays the same

$$\tilde{y} = \mathbf{Q}^T \tilde{b}$$

$$y_1 = \mathbf{Q}_1^T \tilde{b} = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 2 \\ -1 \end{bmatrix} = -\frac{2}{\sqrt{2}} = -\sqrt{2}$$

$$y_2 = 0$$

$$y_3 = \mathbf{Q}_3^T \tilde{b} = \begin{bmatrix} 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 2 \\ -1 \end{bmatrix} = -1$$

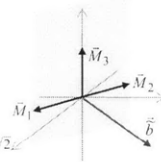
\tilde{y} is the same as the above example

(3) $\mathbf{R}\tilde{x} = \tilde{y}$

Since y is the same as the example above, the solution is the same.

This solution will give us the solution \tilde{b} , but not \tilde{b} (we know this because we choose this \tilde{b} outside of $\text{span}\{\mathbf{M}\}$ knowing that we cannot get a solution). What is going on?

We find that \tilde{b} is the projection of \tilde{b} onto the subspace of the columns of \mathbf{M} .



Minimization View of QR

❖ Alternative Formulations

Definition of the residual, R :

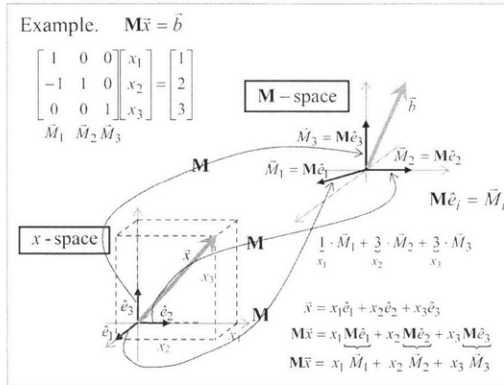
$$R(\tilde{x}) = \tilde{b} - \mathbf{M}\tilde{x}$$

Find \tilde{x} which satisfies	Minimize over all \tilde{x}
$\mathbf{M}\tilde{x} = \tilde{b}$	$\ R(\tilde{x})\ ^2$

If $\tilde{b} \in \text{range}\{\mathbf{M}\}$

$$\text{Then } \mathbf{M}\tilde{x} = \tilde{b} \Leftrightarrow \min_{\tilde{x}} \|R(\tilde{x})\|^2 = 0$$

Minimization extends to singular or non-square case!



❖ One-dimensional Minimization

Suppose $\tilde{x} = x_1 \tilde{e}_1$ and therefore $\mathbf{M}\tilde{x} = x_1 \mathbf{M}\tilde{e}_1 = x_1 \tilde{M}_1$

One-dimensional minimization:

$$\|R(\tilde{x})\|^2 = R(\tilde{x})^T R(\tilde{x}) = (\tilde{b} - x_1 \mathbf{M}\tilde{e}_1)^T (\tilde{b} - x_1 \mathbf{M}\tilde{e}_1)$$

$$= \tilde{b}^T \tilde{b} - 2x_1 \tilde{b}^T \mathbf{M}\tilde{e}_1 + x_1^2 (\mathbf{M}\tilde{e}_1)^T (\mathbf{M}\tilde{e}_1)$$

Note:
$\tilde{b}^T x_1 \mathbf{M}\tilde{e}_1 = (x_1 \mathbf{M}\tilde{e}_1)^T \tilde{b}$
scalar scalar

$$\frac{d}{dx} R(\tilde{x})^T R(\tilde{x}) = -2\tilde{b}^T \mathbf{M}\tilde{e}_1 + 2x_1 (\mathbf{M}\tilde{e}_1)^T (\mathbf{M}\tilde{e}_1) = 0$$

$$x_1 = \frac{\tilde{b}^T \mathbf{M}\tilde{e}_1}{\tilde{e}_1^T \mathbf{M}^T \mathbf{M}\tilde{e}_1}$$

Normalization

$$(\mathbf{M}\tilde{e}_1)^T (\mathbf{M}\tilde{e}_1) = \tilde{M}_1^T \tilde{M}_1$$

Example. $M\vec{x} = \vec{b}$

$$\begin{bmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

Let's look for the best solution using only $\tilde{M}_1 = M\hat{e}_1$

$\begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} x^1 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$

M-space

$\tilde{M}_1 = M\hat{e}_1$

$\hat{Q}_1 = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \\ 0 \end{bmatrix}$

$y_1 = \vec{b} \cdot \hat{Q}_1 = -\frac{1}{\sqrt{2}}$

$M\vec{x}^1 = y_1 \hat{Q}_1 = \begin{bmatrix} -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ 0 \end{bmatrix}$

$\vec{x}^1 = y_1 \hat{e}_1 = \begin{bmatrix} -\frac{1}{\sqrt{2}} \\ 0 \\ 0 \end{bmatrix}$

- Normalize $M\hat{e}_1 \rightarrow \hat{Q}_1$
 $\hat{Q}_1 = \frac{M\hat{e}_1}{\|M\hat{e}_1\|} = \frac{\begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix}}{\sqrt{2}} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \\ 0 \end{bmatrix}$
- Projection
 The Closest is $\vec{y}_1 = \vec{b} \cdot (M\hat{e}_1) = \vec{b} \cdot \hat{Q}_1 = -\frac{1}{\sqrt{2}}$
 $M\vec{x}^1 = (\vec{b} \cdot \hat{Q}_1) \hat{Q}_1 = y_1 M\hat{e}_1 = \begin{bmatrix} -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ 0 \end{bmatrix}$
- Best Approximation
 $\vec{x}^1 = y_1 \hat{e}_1 = -\frac{1}{\sqrt{2}} \cdot \hat{e}_1 = -\frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} -\frac{1}{\sqrt{2}} \\ 0 \\ 0 \end{bmatrix}$

x-space

M-space

$\vec{p}_1 = \frac{\hat{e}_1}{\|M\hat{e}_1\|}$

$x_1 M\hat{e}_1 = y_1 M\hat{p}_1$

$y_1 = \vec{b} \cdot M\hat{p}_1$

$M\hat{p}_1 = \hat{Q}_1 = \frac{M\hat{e}_1}{\|M\hat{e}_1\|}$

Compare to the result from minimization:
 $x_1 = \frac{\vec{b}^T M\hat{e}_1}{\hat{e}_1^T M^T M\hat{e}_1}$

The orthogonal projection minimizes the residual!

Search & minimize in one direction first

Orthogonal Projection:

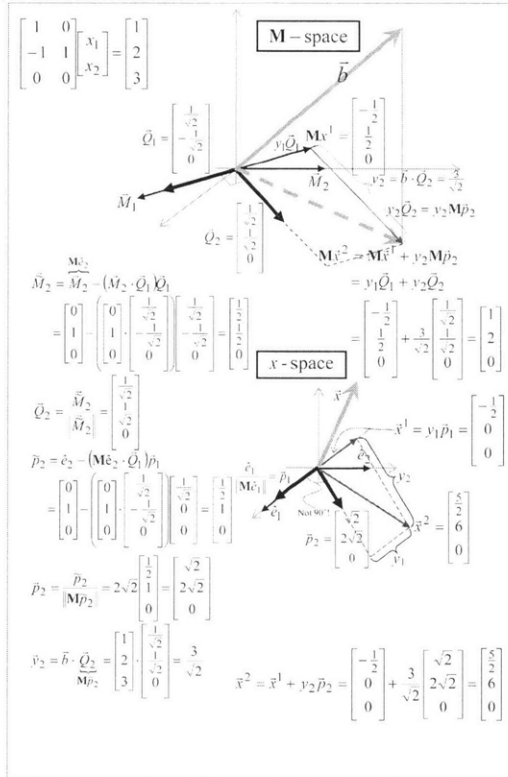
$x_1 \tilde{M}_1 = y_1 \hat{Q}_1 = \left(\vec{b} \cdot \frac{\tilde{M}_1}{\|\tilde{M}_1\|} \right) \frac{\tilde{M}_1}{\|\tilde{M}_1\|}$

$x_1 \tilde{M}_1 = \frac{\vec{b}^T M\hat{e}_1}{\hat{e}_1^T M^T M\hat{e}_1} \tilde{M}_1$

$x_1 = \frac{\vec{b}^T M\hat{e}_1}{\hat{e}_1^T M^T M\hat{e}_1}$

One-dimensional minimization yields same result as orthogonal projection on the column!

❖ Two-dimensional Minimization



Now $\tilde{x} = x_1\tilde{e}_1 + x_2\tilde{e}_2$ and $\tilde{M}\tilde{x} = x_1\tilde{M}\tilde{e}_1 + x_2\tilde{M}\tilde{e}_2$
Residual Minimization

$$\begin{aligned} \|R(\tilde{x})\|^2 &= R(\tilde{x})^T R(\tilde{x}) = (\tilde{b} - x_1\tilde{M}\tilde{e}_1 - x_2\tilde{M}\tilde{e}_2)^T (\tilde{b} - x_1\tilde{M}\tilde{e}_1 - x_2\tilde{M}\tilde{e}_2) \\ &= \tilde{b}^T \tilde{b} - 2x_1\tilde{b}^T \tilde{M}\tilde{e}_1 + x_1^2 (\tilde{M}\tilde{e}_1)^T (\tilde{M}\tilde{e}_1) - 2x_2\tilde{b}^T \tilde{M}\tilde{e}_2 + x_2^2 (\tilde{M}\tilde{e}_2)^T (\tilde{M}\tilde{e}_2) \\ &\quad + \underbrace{2x_1x_2 (\tilde{M}\tilde{e}_1)^T (\tilde{M}\tilde{e}_2)}_{\text{coupling term}} \end{aligned}$$

Minimization is difficult because of the coupling term.

More general search directions

$x = y_1\tilde{p}_1 + y_2\tilde{p}_2$ such that:

$$\tilde{M}x = y_1\tilde{M}\tilde{p}_1 + y_2\tilde{M}\tilde{p}_2 \quad \text{span}\{\tilde{p}_1, \tilde{p}_2\} = \text{span}\{\tilde{e}_1, \tilde{e}_2\}$$

$$\begin{aligned} R(\tilde{x})^T R(\tilde{x}) &= \tilde{b}^T \tilde{b} - 2y_1\tilde{b}^T \tilde{M}\tilde{p}_1 + y_1^2 (\tilde{M}\tilde{p}_1)^T (\tilde{M}\tilde{p}_1) \\ &\quad - 2y_2\tilde{b}^T \tilde{M}\tilde{p}_2 + y_2^2 (\tilde{M}\tilde{p}_2)^T (\tilde{M}\tilde{p}_2) + \underbrace{2y_1y_2 (\tilde{M}\tilde{p}_1)^T (\tilde{M}\tilde{p}_2)}_{\text{coupling term}} \end{aligned}$$

If $\tilde{p}_i^T \tilde{M}^T \tilde{M} \tilde{p}_j = 0$ Minimizations Decouple!!!

If search directions are M-orthonormal

$$\frac{\tilde{M}\tilde{p}_i \cdot \tilde{M}\tilde{p}_j}{\tilde{Q}_i \cdot \tilde{Q}_j} = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases}$$

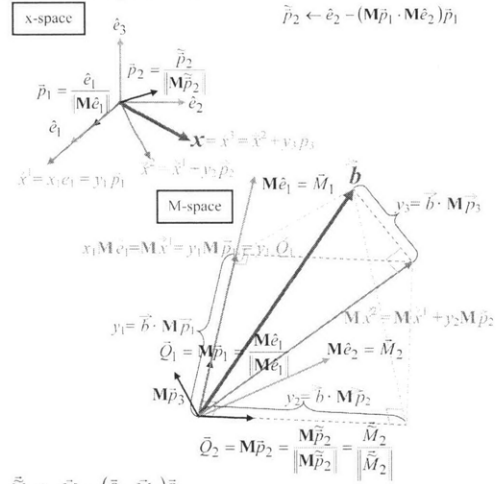
Decoupled minimizations can be done individually

$$\begin{aligned} \text{Minimize: } & -2y_1(\tilde{b}^T \tilde{M}\tilde{p}_1) + y_1^2 (\tilde{M}\tilde{p}_1)^T (\tilde{M}\tilde{p}_1) \\ & -2y_1(\tilde{b} \cdot \tilde{Q}_1) + y_1^2 (\tilde{Q}_1 \cdot \tilde{Q}_1) \end{aligned}$$

$$\text{Differentiating: } -2(\tilde{b} \cdot \tilde{Q}_1) + 2y_1 = 0 \Rightarrow y_1 = \tilde{b} \cdot \tilde{Q}_1$$

Minimization yields same result as orthogonal projection!

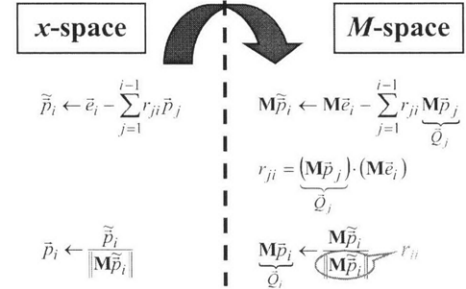
❖ Picture of arbitrary subspace method



$\tilde{M}_2 \leftarrow \tilde{M}_2 - (\tilde{Q}_1 \cdot \tilde{M}_2)\tilde{Q}_1$
 $\mathbf{M}\tilde{p}_2 \leftarrow \mathbf{M}\tilde{e}_2 - (\mathbf{M}\tilde{p}_1 \cdot \mathbf{M}\tilde{e}_2)\mathbf{M}\tilde{p}_1$

$\{\tilde{e}_1, \tilde{e}_2, \dots, \tilde{e}_i\}$ $\{\tilde{p}_1, \tilde{p}_2, \dots, \tilde{p}_i\}$	$\left\{ \frac{\tilde{Q}_1}{\ \mathbf{M}\tilde{p}_1\ }, \frac{\tilde{Q}_2}{\ \mathbf{M}\tilde{p}_2\ }, \dots, \frac{\tilde{Q}_i}{\ \mathbf{M}\tilde{p}_i\ } \right\}$ <i>Orthonormal</i>
$\tilde{x}^i = \underbrace{(\tilde{b} \cdot \mathbf{M}\tilde{p}_1)}_{y_1}\tilde{p}_1 + \dots + \underbrace{(\tilde{b} \cdot \mathbf{M}\tilde{p}_i)}_{y_i}\tilde{p}_i$	
$\tilde{x}^{i+1} = \underbrace{(\tilde{b} \cdot \mathbf{M}\tilde{p}_1)}_{y_1}\tilde{p}_1 + \dots + \underbrace{(\tilde{b} \cdot \mathbf{M}\tilde{p}_i)}_{y_i}\tilde{p}_i + \underbrace{(\tilde{b} \cdot \mathbf{M}\tilde{p}_{i+1})}_{y_{i+1}}\tilde{p}_{i+1}$	

Forming M-orthonormal minimization directions
 The i^{th} search direction equals M-orthonormalized unit vector



- ❖ Minimization Algorithm
- For $i = 1$ to N
 - $\tilde{p}_i \leftarrow \tilde{e}_i$ "For each Target Column"
 - For $j = 1$ to $i - 1$
 - $r_{ji} \leftarrow \frac{\tilde{Q}_j \cdot \tilde{M}_i}{\tilde{M}_j \cdot \tilde{M}_j}$ "For each Source Column left of target"
 - $\tilde{p}_i \leftarrow \tilde{p}_i - r_{ji}\tilde{p}_j$ $\left. \begin{matrix} \mathbf{M} \text{- Orthogonalize} \\ \text{Search Direction} \end{matrix} \right\}$
 - $r_{ii} \leftarrow \frac{\tilde{Q}_i}{\|\mathbf{M}\tilde{p}_i\|}$ $\left. \begin{matrix} \mathbf{M} \text{- Normalize Search Direction} \\ \text{Search Direction} \end{matrix} \right\}$
 - $\tilde{p}_i \leftarrow \frac{\tilde{p}_i}{r_{ii}}$
 - $y_i \leftarrow \tilde{b} \cdot \frac{\tilde{Q}_i}{\|\mathbf{M}\tilde{p}_i\|}$ Calculate projection on new direction
 - $x \leftarrow x + y_i\tilde{p}_i$ Update the solution

i directions

$$\begin{aligned} & \{\hat{e}_1, \hat{e}_2, \dots, \hat{e}_i\} \\ & \{\bar{p}_1, \bar{p}_2, \dots, \bar{p}_i\} \\ & \bar{x}^i = \underbrace{(\bar{b} \cdot \mathbf{M}\bar{p}_1)}_{y_1} \bar{p}_1 + \dots + \underbrace{(\bar{b} \cdot \mathbf{M}\bar{p}_i)}_{y_i} \bar{p}_i \end{aligned}$$

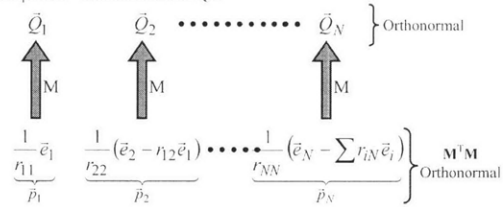
$\{\mathbf{M}\bar{p}_1, \mathbf{M}\bar{p}_2, \dots, \mathbf{M}\bar{p}_i\}$
orthonormal

i+1 directions

$$\begin{aligned} & \{\hat{e}_1, \hat{e}_2, \dots, \hat{e}_i, \hat{e}_{i+1}\} \\ & \{\bar{p}_1, \bar{p}_2, \dots, \bar{p}_i, \bar{p}_{i+1}\} \\ & \bar{x}^{i+1} = \underbrace{(\bar{b} \cdot \mathbf{M}\bar{p}_1)}_{y_1} \bar{p}_1 + \dots + \underbrace{(\bar{b} \cdot \mathbf{M}\bar{p}_i)}_{y_i} \bar{p}_i + \underbrace{(\bar{b} \cdot \mathbf{M}\bar{p}_{i+1})}_{y_{i+1}} \bar{p}_{i+1} \end{aligned}$$

$\{\mathbf{M}\bar{p}_1, \mathbf{M}\bar{p}_2, \dots, \mathbf{M}\bar{p}_i, \mathbf{M}\bar{p}_{i+1}\}$
orthonormal

❖ Comparison – Minimization and QR



$$\begin{bmatrix} I & \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \end{bmatrix} \begin{matrix} \uparrow \\ \\ \downarrow \end{matrix} = \bar{b} \quad \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} \text{ exact solution}$$

Solve this – choose the first search direction to be:

- \hat{e}_1 $\bar{x}^1 = x_1 \hat{e}_1 = b_{11} \hat{e}_1$ (first step)
- \hat{e}_{12} $\bar{x}^1 = x_1 \hat{e}_{12} = b_{12} \hat{e}_{12}$ (first step)
- \bar{b} $\bar{x}^1 = v_1 \bar{b}$ (search direction, first step)

❖ Search Direction
 Orthogonalized unit vectors \rightarrow search directions
 $\left\{ \bar{e}_1, \bar{e}_2, \dots, \bar{e}_N \right\}$ unit vectors $\xrightarrow{\mathbf{M}^T \mathbf{M}}$ Orthogonalization $\left\{ \bar{p}_1, \bar{p}_2, \dots, \bar{p}_N \right\}$ search directions

Example. Simple Problem.

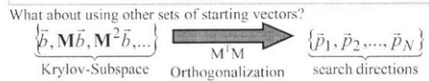
$$\begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_N \end{bmatrix}$$

Step 1. $\bar{x}^1 = \begin{bmatrix} b_1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$ Step 2. $\bar{x}^2 = \begin{bmatrix} b_1 \\ b_2 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$... Step N. $\bar{x}^N = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_N \end{bmatrix}$

If I am searching along $M\bar{p}_1 = \bar{e}_1$ the best I can do is $\bar{b} \cdot M\bar{p}_1 = \bar{x}^1$

$\bar{e}_1, \bar{e}_2, \dots, \bar{e}_N \xrightarrow{M} \bar{e}_1, \bar{e}_2, \dots, \bar{e}_N$
 $\downarrow \qquad \qquad \downarrow$
 $M\bar{p}_1 \qquad \qquad M\bar{p}_N$
 Already M-orthonormal

OR needs N steps!



Example. $\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$

$\bar{e}_1 \leftarrow M\bar{e}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ $\bar{e}_2 \leftarrow M\bar{e}_2 = \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix}$ $\bar{x}^1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ $\bar{x}^2 = \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix}$

What if this is the problem:
 $\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 300 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$

In general, for $L\bar{x} = \bar{b}$ it will take N iterations to get to the solution using this method

If \bar{e}_3 is chosen as the first search direction, what is the result?

Example. $\begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ & & & & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} = \begin{bmatrix} 7 \\ 23 \\ 3 \\ -2 \\ 1 \end{bmatrix}$

$\bar{x}^1 = \begin{bmatrix} 7 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\bar{x}^2 = \begin{bmatrix} 7 \\ 23 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\bar{x}^3 = \begin{bmatrix} 7 \\ 23 \\ 3 \\ 0 \\ 0 \end{bmatrix}$ $\bar{x}^4 = \begin{bmatrix} 7 \\ 23 \\ 3 \\ -2 \\ 0 \end{bmatrix}$ $\bar{x}^5 = \begin{bmatrix} 7 \\ 23 \\ 3 \\ -2 \\ 1 \end{bmatrix}$

$M\bar{v} = \bar{Q}_i = \bar{p}_i$ OR takes n steps

What if $\bar{b} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 23 \end{bmatrix}$ $\bar{x}^1 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$... $\bar{x}^4 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\bar{x}^5 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 23 \end{bmatrix}$

Solution: take $\bar{x}^1 = \bar{b}$ huge error

Example.

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

What is the best initial search direction?

$\vec{b} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$

$\vec{p}_1 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$ $\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$ \mathbf{M} $\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$ x -space

\mathbf{M} is the identity, so $\mathbf{M}\vec{b} = \vec{b}$

$\Rightarrow \vec{p}_1 = \frac{\vec{b}}{\|\vec{b}\|}$

$y_1 = \vec{b} \cdot \mathbf{M}\vec{p}_1 = \|\vec{b}\|$

$y_1 \cdot \mathbf{M}\vec{p}_1 = \|\vec{b}\| \cdot \frac{\vec{b}}{\|\vec{b}\|} = \vec{b}$

Try $\{\vec{b}, \dots\}$

$$\vec{b} \xrightarrow{\mathbf{M}} \mathbf{M}\vec{b} = \vec{b}$$

Normalize

$$\vec{p}_1 = \frac{\vec{b}}{\|\mathbf{M}\vec{b}\|} = \frac{\vec{b}}{\|\vec{b}\|} \xrightarrow{\mathbf{M}} \frac{\mathbf{M}\vec{b}}{\|\mathbf{M}\vec{b}\|} = \frac{\vec{b}}{\|\vec{b}\|}$$

$y_1 = \vec{b} \cdot \mathbf{M}\vec{p}_1 = \|\vec{b}\|$ $\mathbf{M}\vec{p}_1 = \frac{\vec{b}}{\|\vec{b}\|}$

$\vec{p}_1 = y_1 \mathbf{M}\vec{p}_1$

$= \|\vec{b}\| \frac{\mathbf{I} \cdot \vec{b}}{\|\vec{b}\|} = \vec{b}$ One Step

Use QR to solve $\mathbf{M}\mathbf{x} = \mathbf{b}$ \mathbf{M} is singular

\mathbf{x} minimizes $(\mathbf{b} - \mathbf{M}\mathbf{x})^T (\mathbf{b} - \mathbf{M}\mathbf{x})$
 $\mathbf{M}\mathbf{z} = \mathbf{0}$
 \mathbf{z} nonzero exists if \mathbf{M} is singular

ARBITRARY SUBSPACE ALGORITHM

$$\begin{bmatrix} 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots \end{bmatrix}$$

QR Algorithm Search Directions

Approach to Approximately Solving $\mathbf{M}\mathbf{x} = \mathbf{b}$

Pick a k -dimensional subspace $\Rightarrow \left\{ \begin{bmatrix} w_{0,1} \\ \vdots \\ w_{0,k} \end{bmatrix}, \dots, \begin{bmatrix} w_{k-1,1} \\ \vdots \\ w_{k-1,k} \end{bmatrix} \right\} \equiv \{\vec{w}_0, \dots, \vec{w}_{k-1}\}$

Approximate \vec{x}^k as a weighted sum of $\{\vec{w}_0, \dots, \vec{w}_{k-1}\} \Rightarrow \vec{x}^k = \sum_{i=0}^{k-1} \alpha_i \vec{w}_i$

The residual is defined as $\vec{r}^k \equiv \vec{b} - \mathbf{M}\vec{x}^k$

If $\vec{x}^k = \sum_{i=0}^{k-1} \alpha_i \vec{w}_i \Rightarrow \vec{r}^k = \vec{b} - \mathbf{M}\vec{x}^k = \vec{b} - \sum_{i=0}^{k-1} \alpha_i \mathbf{M}\vec{w}_i$

Residual Minimizing Idea:

Pick α_i 's to minimize

$$\|\vec{r}^k\|_2^2 = (\vec{r}^k)^T (\vec{r}^k) = \left(\vec{b} - \sum_{i=0}^{k-1} \alpha_i \mathbf{M}\vec{w}_i \right)^T \left(\vec{b} - \sum_{i=0}^{k-1} \alpha_i \mathbf{M}\vec{w}_i \right)$$

Computational Approach

Minimizing $\|\vec{r}^k\|_2^2 = \left\| \vec{b} - \sum_{i=0}^{k-1} \alpha_i \mathbf{M}\vec{w}_i \right\|_2^2$ is easy if $\mathbf{M}\vec{w}_i$ were orthonormal!

Create a set of vectors $\{\vec{p}_0, \vec{p}_1, \dots, \vec{p}_{k-1}\}$ such that

$span\{\vec{p}_0, \dots, \vec{p}_{k-1}\} = span\{\vec{w}_0, \dots, \vec{w}_{k-1}\}$

$$\mathbf{M}\vec{p}_i \cdot \mathbf{M}\vec{p}_j = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases}$$

Gram-Schmidt on $M\tilde{w}_i$'s

$$\begin{aligned} \tilde{w}_0 &\rightarrow \tilde{p}_0 \\ \tilde{w}_1 &\rightarrow \tilde{p}_1 \quad \text{want } (M\tilde{p}_1) \cdot (M\tilde{p}_0) = 0 \\ &\quad \tilde{p}_1 = \tilde{w}_1 - \alpha\tilde{w}_0 = \tilde{w}_1 - \alpha\tilde{p}_0 \\ &\quad \Rightarrow (M\tilde{w}_1 - M\alpha\tilde{p}_0)^T (M\tilde{p}_0) = 0 \\ &\quad \Rightarrow \underbrace{(M\tilde{w}_1)^T (M\tilde{p}_0)}_{\text{if } \tilde{w}_1 \perp \tilde{p}_0 \Rightarrow 0} = \alpha \underbrace{(M\tilde{p}_0)^T (M\tilde{p}_0)}_1 \end{aligned}$$

Tie this back to the QR algorithm.

$$\begin{aligned} [\tilde{M}_1 \quad \dots \quad \tilde{M}_N] &\Rightarrow [\underbrace{M\hat{e}_1 \quad \dots \quad M\hat{e}_N}_{\text{columns orthogonal}}] \\ &= [\tilde{Q}_1 \quad \dots \quad \tilde{Q}_N] [R] \end{aligned}$$

M-space

$$M\tilde{p}_i \leftarrow M\tilde{w}_i - \sum_{j=0}^{i-1} \underbrace{\left(\frac{\tilde{Q}_j^T}{\tilde{Q}_j} \cdot M\tilde{w}_i \right)}_{\tilde{r}_{ij}} \underbrace{M\tilde{p}_j}_{\tilde{Q}_j}$$

x-space

$$\tilde{p}_i \leftarrow \tilde{w}_i - \sum_{j=0}^{i-1} (M\tilde{p}_j \cdot M\tilde{w}_i) \tilde{p}_j$$

Arbitrary Subspace Solution Algorithm

Given M, b and a set of search directions $\{\tilde{w}_0, \tilde{w}_1, \dots, \tilde{w}_k\}$

1) orthogonalize the $M\tilde{w}_i$'s

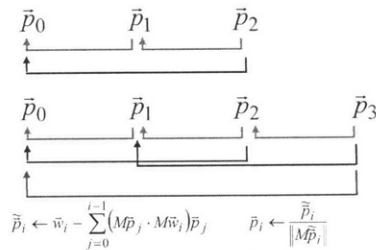
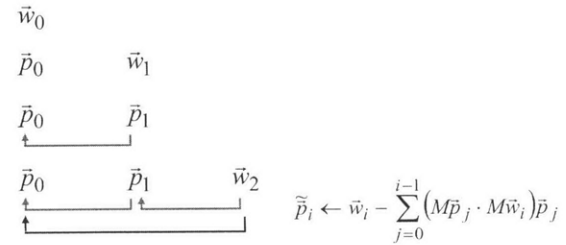
$$\text{for } i = 1 \text{ to } k \quad \tilde{p}_i = \tilde{w}_i - \sum_{j=0}^{i-1} \frac{(M\tilde{w}_i)^T (M\tilde{p}_j)}{(M\tilde{p}_j)^T (M\tilde{p}_j)} \tilde{p}_j$$

2) compute the \tilde{x} minimizing solution \tilde{x}^{k+1}

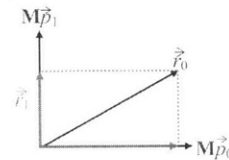
$$\tilde{x}^{k+1} = \sum_{i=0}^k \frac{(e^0)^T (M\tilde{p}_i)}{(M\tilde{p}_i)^T (M\tilde{p}_i)} \tilde{p}_i$$

Orthogonalization of Search Directions

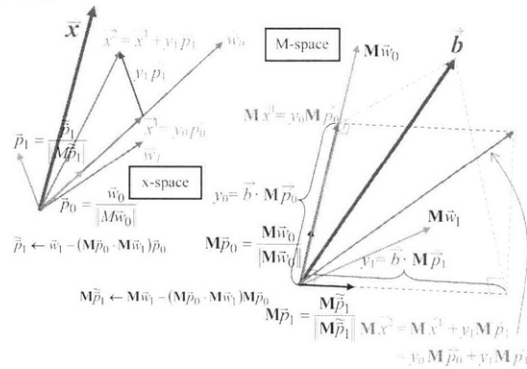
M-orthonormalizing \tilde{w}_i



Compute the \tilde{x} minimizing solution \tilde{x}^{k+1}



The Picture



Arbitrary Subspace Generation

If $span\{\tilde{w}_0, \tilde{w}_1, \dots, \tilde{w}_k\} = span\{r^0, Mr^0, \dots, M^k r^0\}$

$$\tilde{x}^{k+1} = \sum_{i=0}^k \alpha_i M^i r^0 = \underbrace{z_k(M)}_{k^{th} \text{ order polynomial}} r^0$$

$$r^{k+1} = r^0 - \sum_{i=0}^k \alpha_i M^i r^0 = (I - Mz_k(M))r^0$$

Note: for any $\alpha_0 \neq 0$

$$span\{r^0, r^1\} = r^0 - \alpha_0 M r^0 = span\{r^0, M r^0\}$$

$$\tilde{x}^3 = \alpha_1 \tilde{b} + \alpha_2 M \tilde{b} + \alpha_3 M^2 \tilde{b} = \underbrace{(\alpha_1 I + \alpha_2 M + \alpha_3 M^2)}_{\text{polynomial}} \tilde{b}$$

Krylov Subspace Steps:

- (1) Pick a Search Direction
- (2) Make search direction $M^i M$ orthogonal
- (3) Update x (minimize residual)
- (4) Update residual

$$\begin{aligned} r^0 &= \tilde{b} \\ p^0 &= r^0 \\ x^0 &= \alpha_0 p^0 \end{aligned} \quad \text{First Pass}$$

Algorithm

For $k = 0$ to $N - 1$
 $\tilde{p}_k \leftarrow \tilde{w}_k$ The algorithm is the same even if I use random search directions.
 For $j = 1$ to $k - 1$
 $\tilde{p}_k \leftarrow \tilde{p}_k - (M\tilde{p}_j \cdot M\tilde{p}_k)\tilde{p}_j$ M - Orthogonalize Search Direction
 $\tilde{p}_k \leftarrow \frac{\tilde{p}_k}{\|M\tilde{p}_k\|}$ M - Normalize search direction
 $y_k \leftarrow \tilde{b} \cdot (M\tilde{p}_k)$ Calculate projection on new direction
 $\tilde{x}^{k+1} \leftarrow \tilde{x}^k + y_k \tilde{p}_k$ Update the solution

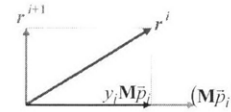
After i^{th} iteration:

$$\tilde{b} = y_0 M\tilde{p}_0 + \dots + y_{i-1} M\tilde{p}_{i-1} + \underbrace{r^i}_{\text{residual}}$$

Since they are all orthogonal to $(M\tilde{p}_i) \Rightarrow y_i = \tilde{b} \cdot (M\tilde{p}_i) = r^i \cdot (M\tilde{p}_i)$

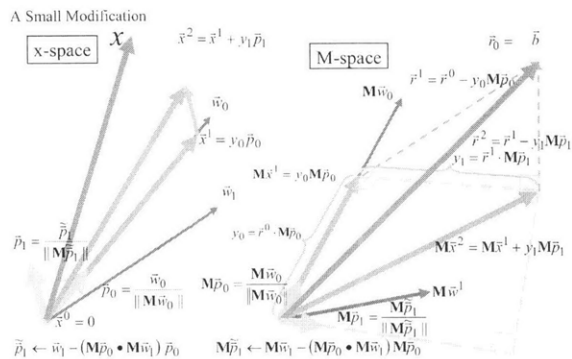
After $i + 1$ iterations:

$$\tilde{b} = y_0 M\tilde{p}_0 + \dots + y_{i-1} M\tilde{p}_{i-1} + \underbrace{y_i M\tilde{p}_i + r^{i+1}}_{r_i}$$



Residual Formula.

$$\begin{aligned} \tilde{x}^{j+1} &= \tilde{x}^j + \alpha_j \tilde{p}_j \\ r^{j+1} &= \tilde{b} - Mx^{j+1} = \tilde{b} - M(x^j + \alpha_j \tilde{p}_j) = \underbrace{\tilde{b} - Mx^j}_{r^j} - M\alpha_j \tilde{p}_j \end{aligned}$$



Algorithm

For $i = 1$ to N

$\tilde{p}_i \leftarrow \tilde{w}_i$

For $j = 1$ to $i - 1$

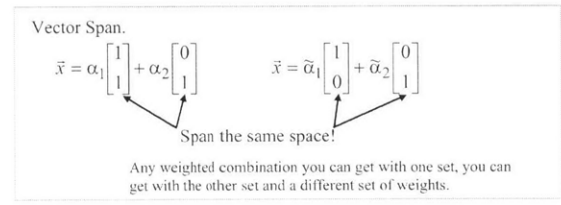
$\tilde{p}_i \leftarrow \tilde{p}_i - (M \tilde{p}_j \cdot M \tilde{p}_i) \tilde{p}_j$ M - Orthogonalize Search Direction

$\tilde{p}_i \leftarrow \frac{\tilde{p}_i}{\|M \tilde{p}_i\|}$ M - Normalize search direction

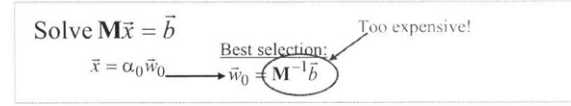
$y_i \leftarrow \tilde{r}^i \cdot (M \tilde{p}_i)$ Calculate projection on new direction

$\tilde{x}^{i+1} \leftarrow \tilde{x}^i + y_i \tilde{p}_i$ Update the solution

$\tilde{r}^{i+1} \leftarrow \tilde{r}^i - y_i M \tilde{p}_i$ Update the residual



GENERALIZED CONJUGATE RESIDUAL ALGORITHM



Selection of the Search Directions

Criteria for selecting $\tilde{w}_0, \dots, \tilde{w}_{k-1}$

All that matters is the $span\{\tilde{w}_0, \dots, \tilde{w}_{k-1}\}$

$\exists \alpha_i$'s such that $\tilde{b} - M \tilde{x}^k = \tilde{b} - \sum_{i=0}^{k-1} \alpha_i M \tilde{w}_i$ is small

$M^{-1} \tilde{b} \approx$ in the $span\{\tilde{w}_0, \dots, \tilde{w}_{k-1}\}$ for $k < N$

One choice, unit vectors, $\tilde{x}^k \in span\{\tilde{e}_1, \dots, \tilde{e}_k\}$

Generates the QR algorithm if $k = N$

Can be terrible if $k < N$

Example.

$$\begin{bmatrix} 1 & & & 0 \\ & \ddots & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} 1 \\ \vdots \\ 0 \end{bmatrix}$$

Only need (1) \bar{w} to solve this problem \rightarrow the first one

$$\bar{w}'s = \left\{ \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots \right\}$$

$$\begin{bmatrix} 1 & & & 0 \\ & \ddots & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 1 \end{bmatrix}$$

If you start from the first search direction, this will take all N search directions to find the correct answer.
Only need (1) \bar{w} to solve this problem \rightarrow the last one

$$\begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

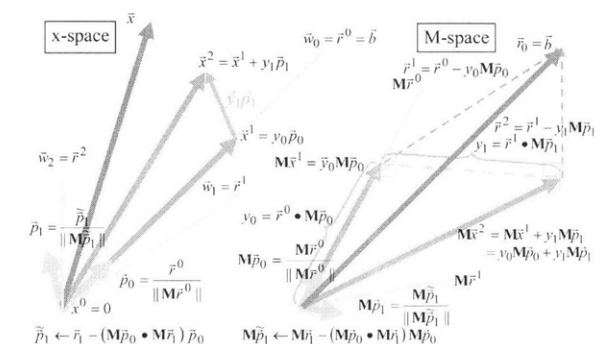
Need all N \bar{w} 's to solve this problem
Could also solve this problem with ONE search direction if that search direction is all ones.

Historical Development
 Assume $\mathbf{M} = \mathbf{M}^T$ (symmetric) and $\bar{x}^T \mathbf{M} \bar{x} > 0$ (positive definite)
 Then the solution of $\mathbf{M}\bar{x} = \bar{b}$ corresponds to the location of the minimum of $f(\bar{x}) = \frac{1}{2} \bar{x}^T \mathbf{M} \bar{x} - \bar{x}^T \bar{b}$
 $\nabla_x f(\bar{x}) = \mathbf{M}\bar{x} - \bar{b} \Rightarrow \bar{x} = \mathbf{M}^{-1} \bar{b}$ minimizes f
 Gradient of $f(x)$ = Residual at x
 Idea: search along the gradient i.e. the steepest descent directions for $f(x)$ i.e. along the current residual direction.
 Pick $\text{span}\{\bar{w}_0, \dots, \bar{w}_{k-1}\} = \text{span}\{\nabla_x f(\bar{x}^0), \dots, \nabla_x f(\bar{x}^{k-1})\} = \text{span}\{\bar{r}^0, \bar{r}^1, \dots, \bar{r}^{k-1}\}$
 Does not extend to non-symmetric, non positive definite case.

Krylov Subspace
 Note: $\text{span}\{\nabla_x f(\bar{x}^0), \dots, \nabla_x f(\bar{x}^{k-1})\} = \text{span}\{\bar{r}^0, \dots, \bar{r}^{k-1}\}$

If $\text{span}\{\bar{w}_0, \dots, \bar{w}_{k-1}\} = \text{span}\{\bar{r}^0, \dots, \bar{r}^{k-1}\}$
 Then $\bar{r}^k = \bar{r}^0 - \sum_{i=0}^{k-1} \alpha_i \mathbf{M} \bar{r}^i$
 and $\text{span}\{\bar{r}^0, \dots, \bar{r}^{k-1}\} = \text{span}\{\bar{r}^0, \mathbf{M} \bar{r}^0, \dots, \mathbf{M}^{k-1} \bar{r}^0\}$
 Krylov Subspace

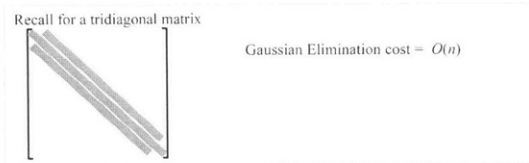
The Picture



Computational Approach
 $\bar{r}^0 = \bar{b} - \mathbf{M}\bar{x}^0 = \bar{b}$
 REPEAT
 Pick next Search Direction
 $\bar{p}_k \leftarrow \bar{r}^k$
 For $j = 1$ to $k-1$
 M-Orthogonalize Search Direction
 $\bar{p}_k \leftarrow \bar{p}_k - (\mathbf{M}\bar{p}_j \cdot \mathbf{M}\bar{p}_k) \bar{p}_j$
 M - Normalize search direction
 $\bar{p}_k \leftarrow \frac{\bar{p}_k}{\|\mathbf{M}\bar{p}_k\|}$
 Optimal step size in new direction
 $y_k \leftarrow \bar{r}^k \cdot (\mathbf{M}\bar{p}_k)$
 Update the solution
 $\bar{x}^{k+1} \leftarrow \bar{x}^k + y_k \bar{p}_k$
 Update the residual
 $\bar{r}^{k+1} \leftarrow \bar{r}^k - y_k \mathbf{M}\bar{p}_k$

Computational Cost
 Algorithm Cost for the k^{th} iteration of GCR:
 $\bar{r}^0 = \bar{b} - \mathbf{M}\bar{x}^0 = \bar{b}$

REPEAT
 $\tilde{p}_k \leftarrow \tilde{r}^k$
 For $j = 1$ to $k-1$ $O(k)$ inner products: $O(kn)$ mult.
 $\tilde{p}_k \leftarrow \tilde{p}_k - (\mathbf{M}\tilde{p}_j) \cdot (\mathbf{M}\tilde{p}_j) \tilde{p}_j$ Matrix vector product: $O(n)$ if sparse
 $\tilde{p}_k \leftarrow \frac{\tilde{p}_k}{\|\mathbf{M}\tilde{p}_k\|}$
 $y_k \leftarrow \tilde{r}^k \cdot (\mathbf{M}\tilde{p}_k)$ Vector inner products: $O(n)$ mult.
 $\tilde{x}^{k+1} \leftarrow \tilde{x}^k + y_k \tilde{p}_k$ $O(n)$ multiplications
 $\tilde{r}^{k+1} \leftarrow \tilde{r}^k + y_k \mathbf{M}\tilde{p}_k$ $O(n)$ multiplications
 total cost = $O(n) + O(2n) + \dots + O(kn) = O(k^2 n)$
 If \mathbf{M} is sparse, as k (number of iterations) approaches n : $O(n^3)$
 Better Converge Fast!!



Algorithm Cost for iteration k
 For $j = 1$ to $k-1$ $O(k)$ inner products: $O(kn)$ mult.
 $\tilde{p}_k \leftarrow \tilde{p}_k - (\mathbf{M}\tilde{p}_j) \cdot (\mathbf{M}\tilde{p}_j) \tilde{p}_j$
 If \mathbf{M} is sparse the inner products are the dominant cost:
 total cost = $O(n) + O(2n) + \dots + O(kn) = O(k^2 n)$
 As k (# of iters) approaches N , total cost = $O(n^3)$
 Better Converge Fast!!
 Symmetric Case – An amazing fact that will not be derived
 If $\mathbf{M} = \mathbf{M}^T$ then $\mathbf{M}\tilde{r}^k \perp \mathbf{M}\tilde{p}_j$ for $j < k-2$
 For $j = 1$ to $k-1$ Only ONE inner products: $O(n)$ mult.
 $\tilde{p}_k \leftarrow \tilde{r}_k - (\mathbf{M}\tilde{p}_{k-1}) \cdot (\mathbf{M}\tilde{p}_{k-1}) \tilde{p}_{k-1}$
 If k (# of iters) $\rightarrow n$, then symmetric, sparse, GCR is worst case $O(n^3)$
 Better Converge Fast!!

		SPARSE	DENSE	Per Iteration Cost
$\mathbf{M}\tilde{r}^k$	Matrix-Vector Product	$O(n)$	$O(n^2)$	
$\sum (\mathbf{M}\tilde{r}^k)^T (\mathbf{M}\tilde{p}_j)$	M-Orthogonalize	$O(kn) \cdot O(n)$	$O(kn) \cdot O(n)$	
$(\tilde{r}^k)^T (\mathbf{M}\tilde{p}_k)$	Optimal Step Size Update	$O(n)$	$O(n)$	
	Total	Symmetric $O(kn) + O(k^2 n)$ matrix-vector product cost inner products cost	$O(kn^2) + O(k^2 n)$ matrix-vector product cost inner products cost	
	Worst Case As $K \rightarrow n$	$O(n^2) + O(n^3)$ $O(n^2) \cdot O(n^2)$ Need to reduce k !!	$O(n^3) + O(n^4)$ $O(n^3) \cdot O(n^2)$ Need to reduce k AND matrix-vector cost	
	Best Case For small K (5-10 iterations)	$O(n) + O(n)$ $O(n)$	$O(n^2) + O(n)$ $O(n^2)$	

Note: If "fast" $\{-O(n)\}$ matrix-vector product then $O(n)$ total also for dense case.

Comparison – GCR & LU

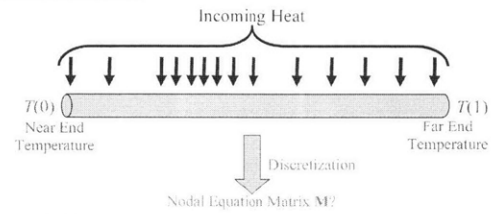
	Sparse Symmetric	Sparse Asymmetric	Dense
GCR	$O(Kn)$	$O(K^2 n)$	$O(Kn^2)$
LU	$O(n^{1.2-1.8})$	$O(n^{1.2-1.8})$	$O(n^3)$

- Making GCR converge fast can be very problem specific
- If the number of iterations k is about constant and small GCR is much faster than LU for dense problems.

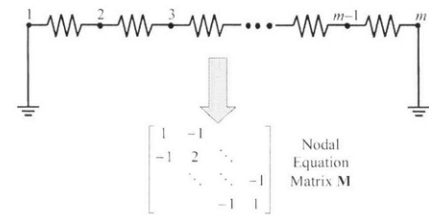
	Sparse	Dense
GCR	$O(n)$	$O(n^2)$
LU	$O(n^{1.2-1.8})$	$O(n^3)$

Convergence Examples

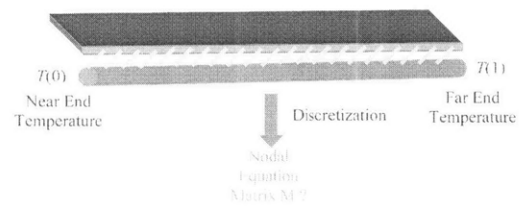
- ❖ "No-leak" Examples
 - Insulated Bar and Matrix



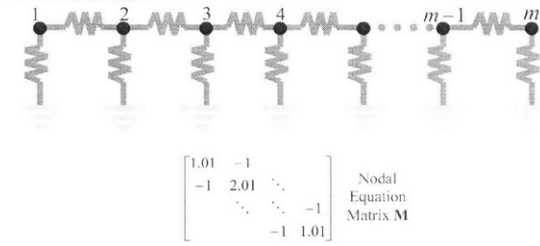
- Circuit and Matrix



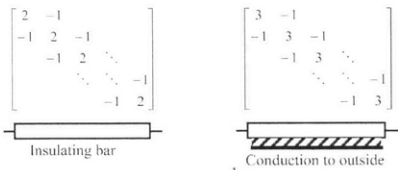
- ❖ "Leaky" Examples
 - Conducting Bar and Matrix



- Circuit and Matrix



Example. Heat conducting bar.



$$u_{ii} = \frac{d}{2 \text{ or } 3} - \frac{1}{u_{i-1,j-1}}$$

Suppose $d=2$

$$u_{22} = 2 - \frac{1}{2} = \frac{3}{2}$$

$$u_{33} = 2 - \frac{1}{\frac{3}{2}} = \frac{4}{3}$$

$$\vdots$$

$$u_{ii} = \frac{i+1}{i}$$

Suppose $d=3$

$$u_{22} = 3 - \frac{1}{3} = \frac{8}{3}$$

$$u_{33} = 3 - \frac{1}{\frac{8}{3}} = \frac{21}{8}$$

$$\vdots$$

pattern is not so obvious

Recurrence formula: $u_{ii} = d - \frac{1}{u_{i-1,j-1}}$

$$\approx d - \frac{1}{d - \frac{1}{d - \frac{1}{d - \frac{1}{d - \dots}}}}$$

Continued Fraction

Let the asymptotic value be x : $u_{ij} = x$ & $u_{i-1,j-1} = x$ for $\lim_{i \rightarrow \infty}$

$$\Rightarrow x = d - \frac{1}{x} \Rightarrow x^2 - dx + 1 = 0$$

For $d=2$: $x \rightarrow 1$

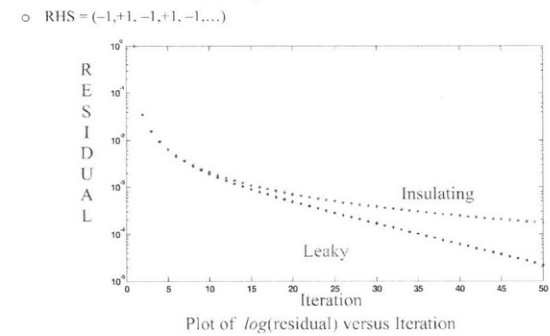
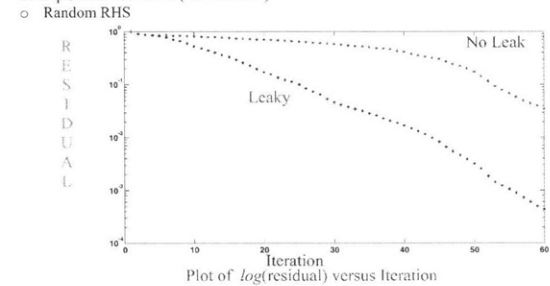
For $d=3$: $x \rightarrow \frac{3-\sqrt{5}}{2}$

Continued Fraction...

$$u_{ii} = \begin{bmatrix} x[i] \leftarrow d \cdot x[i-1] - x[i-2] \\ y[i] \leftarrow d \cdot y[i-1] - y[i-2] \end{bmatrix}$$

simple linear difference equations

❖ GCR Performance
Example with 100 nodes (102 resistors)



INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 9.

GCR Convergence

TODAY'S OUTLINE:

- ❖ Review - Eigenvalues
 - Norms and Spectral Radius
 - Spectral Mapping Theorem
- ❖ GCR Convergence Analysis
 - Polynomial Error Bound
 - Chebychev Bounds
 - Gershgorin Circle Theorem

EIGENVALUES

$$\mathbf{M}\vec{u}_i - \lambda_i \vec{u}_i = 0$$

$$\Rightarrow (\mathbf{M} - \lambda_i \mathbf{I})\vec{u}_i = 0$$

Example.

Circuit is symmetric - expect the node voltages might be the same

$$\begin{bmatrix} 1.1 & -1 \\ -1 & 1.1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \xrightarrow{\text{expected eigenvalue}} (0.1) \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

corresponds to 10Ω resistor

Suppose the voltages are different at the nodes...

$$\begin{bmatrix} 1.1 & -1 \\ -1 & 1.1 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} \xrightarrow{\text{expected eigenvalue}} (2.1) \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

Example. Block Diagonal Matrix

$$\begin{bmatrix} 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & -1 & 2 \end{bmatrix}$$

Example. Lower Triangular Matrix

$$\begin{bmatrix} M_{11} & 0 & 0 & \dots & 0 \\ M_{21} & M_{22} & 0 & \dots & 0 \\ M_{31} & \dots & \dots & \dots & \vdots \\ \vdots & \dots & M_{N-1,N-2} & M_{N-1,N-1} & 0 \\ M_{N,1} & \dots & M_{N,N-2} & M_{N,N-1} & M_{N,N} \end{bmatrix}$$

Eigenvalues are the diagonal elements.

Eigenvalues and eigenvectors of a matrix \mathbf{M} satisfy

$$\mathbf{M}\vec{u}_i = \lambda_i \vec{u}_i$$

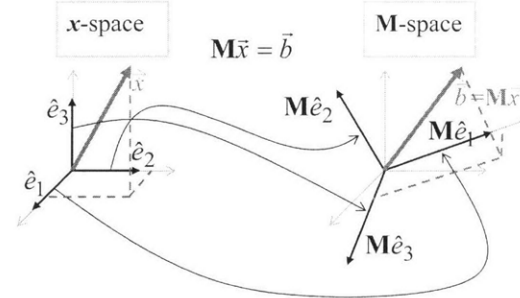
eigenvalue
eigenvector

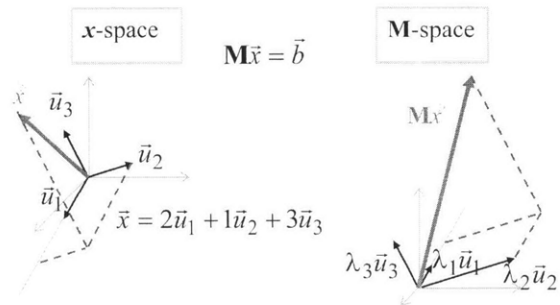
Or, λ_i is an eigenvalue of \mathbf{M} if

$\mathbf{M} - \lambda_i \mathbf{I}$ is singular

\vec{u}_i is an eigenvector of \mathbf{M} if

$$(\mathbf{M} - \lambda_i \mathbf{I})\vec{u}_i = 0$$





$$\mathbf{M}\bar{x} = 2\underbrace{\lambda_1}_{-\frac{1}{2}}\bar{u}_1 + 1\underbrace{\lambda_2}_{2}\bar{u}_2 + 3\underbrace{\lambda_3}_{1}\bar{u}_3$$

$$\mathbf{M}\bar{x} = 2\left(-\frac{1}{2}\right)\bar{u}_1 + 1(2)\bar{u}_2 + 3(1)\bar{u}_3$$

Easy to do – the **M** space uses the same vectors as the **x**-space, they are just stretched by multiplier λ_i .

A Simplifying Assumption

$$\mathbf{M} \begin{bmatrix} \uparrow & \uparrow & \uparrow & \dots & \uparrow \\ \bar{u}_1 & \bar{u}_2 & \bar{u}_3 & \dots & \bar{u}_N \\ \downarrow & \downarrow & \downarrow & \dots & \downarrow \end{bmatrix} = \begin{bmatrix} \uparrow & \uparrow & \uparrow & \dots & \uparrow \\ \lambda_1\bar{u}_1 & \lambda_2\bar{u}_2 & \lambda_3\bar{u}_3 & \dots & \lambda_N\bar{u}_N \\ \downarrow & \downarrow & \downarrow & \dots & \downarrow \end{bmatrix}$$

$$\Rightarrow \left. \begin{matrix} \mathbf{M}\bar{u}_1 = \lambda_1\bar{u}_1 \\ \mathbf{M}\bar{u}_2 = \lambda_2\bar{u}_2 \\ \vdots \end{matrix} \right\} \Rightarrow \mathbf{M}\mathbf{U} = \mathbf{U}\mathbf{\Lambda}$$

The set of all eigenvalues of **M** is known as the Spectrum of **M**.

$$\mathbf{M}\mathbf{U} = \mathbf{U} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_N \end{bmatrix} \Rightarrow \mathbf{U}^{-1}\mathbf{M}\mathbf{U} = \mathbf{\Lambda} \text{ or } \mathbf{M} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1}$$

Almost all $N \times N$ matrices have N linearly independent eigenvectors (i.e. **U** nonsingular)
 Does NOT imply distinct eigenvalues, i.e. λ_i can equal λ_j
 Does NOT imply **M** is nonsingular

Example.) **U** nonsingular \nleftrightarrow Distinct Eigenvalues

$$\mathbf{U} = \mathbf{I} = \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix}$$

- All eigenvalues are equal, $\lambda_i = 1 \forall i$
- Eigenvectors are distinct, $\bar{u}_i = \delta_i$

Example.) **U** nonsingular \nleftrightarrow **M** nonsingular

$$\mathbf{M} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$

zero eigenvalue
M singular

has L.I. eigenvectors: $\bar{v}_1 = \hat{\rho}_1 \quad \bar{v}_2 = \hat{\rho}_2$

$$\left. \begin{matrix} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = 1 \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = 0 \begin{bmatrix} 0 \\ 1 \end{bmatrix} \end{matrix} \right\} \mathbf{U} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

U nonsingular

Example.) Real Matrices \rightarrow Real or Complex Conjugate Eigenvalues

$$\mathbf{M} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1} \quad \text{Real Matrix}$$

$$\lambda \text{ eigenvalue} \Leftrightarrow (\mathbf{M} - \lambda\mathbf{I}) \text{ is singular}$$

$$\mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1} - \lambda\mathbf{U}\mathbf{U}^{-1} \rightarrow \text{singular}$$

$$\mathbf{U} \begin{bmatrix} \lambda_1 - \lambda & & \\ & \ddots & \\ & & \lambda_N - \lambda \end{bmatrix} \mathbf{U}^{-1} \rightarrow \text{singular}$$

$\phi(\lambda) = \pi(\lambda_i - \lambda) = 0 \rightarrow$ characteristic polynomial eigenvalues \rightarrow real or complex conjugate pairs

from the *Fundamental Theorem of Algebra*

Symmetric Matrices $\rightarrow \mathbf{U}^T \mathbf{U} = \mathbf{I}$ Orthonormal Eigenvectors

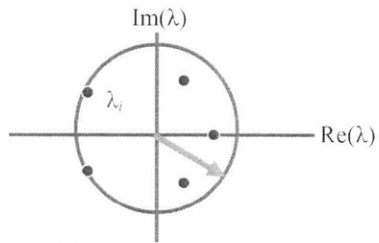
$$\mathbf{M} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1} \quad \mathbf{M}^T = (\mathbf{U}^{-1})^T \mathbf{\Lambda} \mathbf{U}^T$$

$$\text{Symmetric} \rightarrow \mathbf{M} = \mathbf{M}^T$$

$$\mathbf{U}\mathbf{U}^{-1} = (\mathbf{U}^{-1})^T \mathbf{U}^T \Leftrightarrow \mathbf{U}^{-1} = \mathbf{U}^T$$

\rightarrow **U** orthonormal

Spectral Radius



The spectral Radius of M is the radius of the smallest circle, centered at the origin, which encloses all of M 's eigenvalues.

Suppose $\bar{y} = M\bar{x}$

How much larger is \bar{y} than \bar{x} ? OR How much does M magnify \bar{x} ?

$$\begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Rightarrow \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

★ This vector will get the most magnification from matrix M

Look at the direction corresponding to the largest eigenvalue

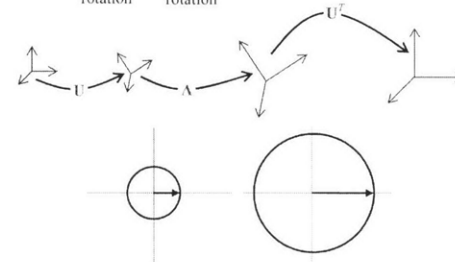
Any induced norm is a bound on the spectral radius
If M is symmetric then $\|M\|_2 = \text{spectral radius}$

Example. $\begin{bmatrix} 1 & 0 \\ 100 & 1 \end{bmatrix}$ $eig\left(\begin{bmatrix} 1 & 0 \\ 100 & 1 \end{bmatrix}\right) = \{1, 1\}$
 $\left\| \begin{bmatrix} 1 & 0 \\ 100 & 1 \end{bmatrix} \right\|_{\infty, 1} = 101$

$$M = UAU^{-1} = U\Lambda U^T \quad M = M^T$$

$$\|M\|_2 = \|UAU^{-1}\|_2 = \|U\|_2 \|\Lambda\|_2 \|U^{-1}\|_2 = \|\Lambda\|_2 = \max_i |\lambda_i|$$

rotation stretch rotation



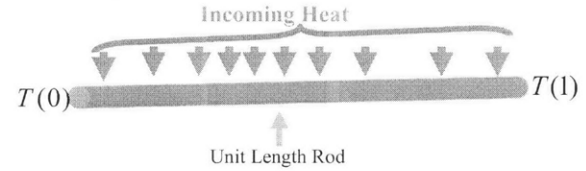
Induced Matrix Norms – standard induced l -norms

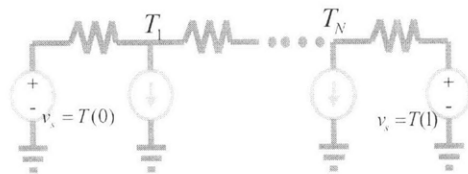
Definition:

$$\|M\|_l = \max_x \frac{\|Mx\|_l}{\|x\|_l} = \max_{\|x\|_l=1} \|Mx\|_l$$

$$\|M\|_2 = \max_i \sqrt{|\lambda_i(M^H M)|} = \max_i |\lambda_i(M)| = \begin{cases} \text{for symmetric real} \\ \text{matrices} = \text{max eigenvalue} \end{cases}$$

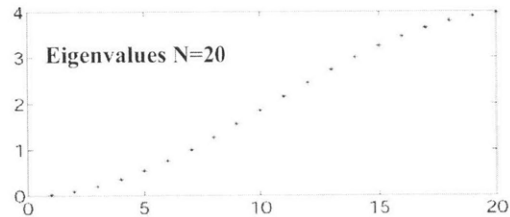
Heat Flow Example





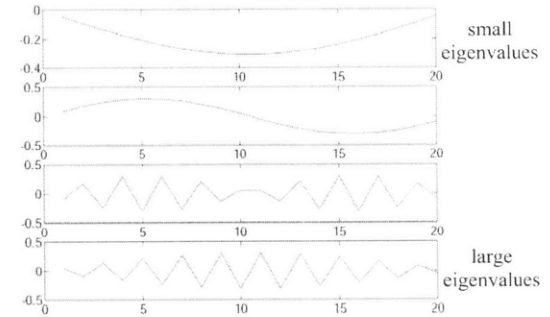
- \bar{x} = Distribution of temperatures
- $M\bar{x}$ = Distribution of incoming/outgoing heat
- $M\bar{x} = \lambda\bar{x}$ The "shape" of the temperature distribution is the same as the incoming/outgoing heat
- λ small means that for that eigenvector distribution of temperatures the resulting heat flow is small.

$$\begin{bmatrix} 1 & -1 & & & \\ -1 & 2 & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & & -1 & 1 \end{bmatrix}$$



heat flow $(h_i) \propto T_i - T_{i+1}$

Four Eigenvectors – Which ones?



GCR CONVERGENCE ANALYSIS

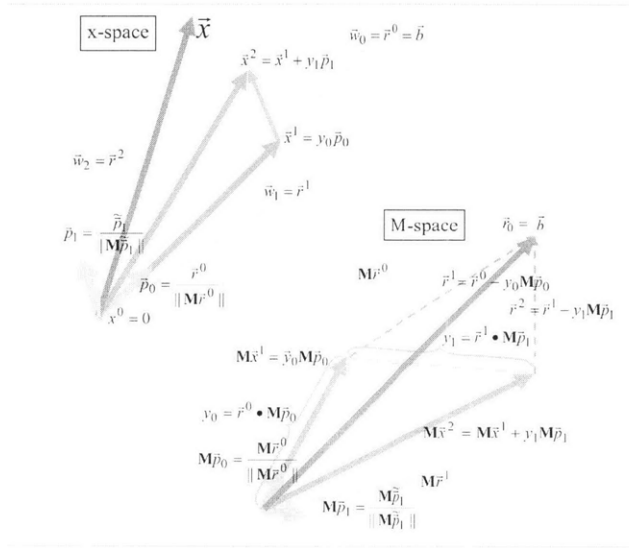
The Iterative Algorithm

- $x^0 = 0$
- $r^0 = b$
- REPEAT
 - Compute $M r^k$
 - $M^T M$ Orthonormalize search directions
 - Update solution x^{k+1}
 - Update residue r^{k+1}
- UNTIL residue small "enough"

Problem Statement

Let's estimate an upper bound for how much we have reduced the residual after k iterations

$$\frac{\|r^k\|}{\|r^0\|} \leq ?$$



Krylov Subspace

Note : $span\{\nabla_x f(x^0), \dots, \nabla_x f(x^{k-1})\} = span\{\vec{r}^0, \dots, \vec{r}^{k-1}\}$
 If $span\{\vec{w}_0, \dots, \vec{w}_{k-1}\} = span\{\vec{r}^0, \dots, \vec{r}^{k-1}\}$
 then $r^k = r^0 - \sum_{i=0}^{k-1} \alpha_i M r^i$
 and $span\{\vec{r}^0, \dots, \vec{r}^{k-1}\} = span\{\vec{r}^0, M\vec{r}^0, \dots, M^{k-1}\vec{r}^0\}$
 Krylov Subspace

$k = 1$
 Step 1.
 $\vec{p}_0 \leftarrow \frac{\vec{r}^0}{\|\mathbf{M}\vec{r}^0\|} \in span\{\vec{r}^0\}$
 $y_0 \leftarrow \vec{r}^{0T} \mathbf{M}\vec{p}_0$ scalar
 $\vec{x}^1 \leftarrow 0 + y_0 \vec{p}_0 \in span\{\vec{p}_0\} = span\{\vec{r}^0\}$
 $\vec{r}^1 \leftarrow \vec{r}^0 - y_0 \mathbf{M}\vec{p}_0 \in span\{\vec{r}^0, \mathbf{M}\vec{p}_0\} = span\{\vec{r}^0, \mathbf{M}\vec{r}^0\}$

$k = 2$
 Step 2.
 $\vec{p}_1 \leftarrow \frac{\vec{r}^1 - \overbrace{(\mathbf{M}\vec{p}_0)^T (\mathbf{M}\vec{r}^1)}^{\text{scalar}} \vec{p}_0}{\|\mathbf{M}\vec{p}_1\|} \in span\{\vec{r}^1, \vec{p}_0\} = span\{\vec{r}^0, \mathbf{M}\vec{r}^0\}$
 $\vec{p}_1 \leftarrow \frac{\vec{r}^1}{\|\mathbf{M}\vec{p}_1\|} \in span\{\vec{r}^0, \mathbf{M}\vec{r}^0\}$
 $y_1 \leftarrow \vec{r}^{1T} \mathbf{M}\vec{p}_0$ scalar
 $\vec{x}^2 \leftarrow \vec{x}^1 + y_1 \vec{p}_1 \in span\{\vec{x}^1, \mathbf{M}\vec{p}_1\} = span\{\vec{r}^0, \mathbf{M}\vec{r}^0\}$
 $\vec{r}^2 \leftarrow \vec{r}^1 - y_1 \mathbf{M}\vec{p}_1 \in span\{\vec{r}^1, \mathbf{M}\vec{p}_1\} = span\{\vec{r}^0, \mathbf{M}\vec{r}^0, \mathbf{M}^2\vec{r}^0\}$

$\vec{x}^k \leftarrow \alpha_0 \vec{r}^0 + \alpha_1 \mathbf{M}\vec{r}^0 + \dots + \alpha_{k-1} \mathbf{M}^{k-1} \vec{r}^0 = \xi_{k-1}(\mathbf{M})\vec{r}^0$
 $\vec{r}^k \leftarrow \beta_0 \vec{r}^0 + \beta_1 \mathbf{M}\vec{r}^0 + \dots + \beta_k \mathbf{M}^k \vec{r}^0 = \varphi_k(\mathbf{M})\vec{r}^0$

Polynomial View

$\vec{p}_{k-1} \leftarrow \vec{r}^{k-1} - (\mathbf{M}\vec{p}_{k-2})^T (\vec{r}^{k-1}) \vec{p}_{k-2}$
 $\vec{p}_{k-1} \in span\{\vec{r}^0, \mathbf{M}\vec{r}^0, \dots, \mathbf{M}^{k-1}\vec{r}^0\}$
 $\vec{x}^k \leftarrow \vec{x}^{k-1} + y_{k-1} \vec{p}_{k-1}$
 $\vec{x}^k \in span\{\vec{r}^0, \mathbf{M}\vec{r}^0, \dots, \mathbf{M}^{k-1}\vec{r}^0\}$
 $\vec{x}^k = \xi_{k-1}(\mathbf{M})\vec{r}^0 \rightarrow (k-1)^{\text{th}}$ order polynomial in \mathbf{M}
 $\vec{r}^k \leftarrow \vec{r}^{k-1} - y_{k-1} \mathbf{M}\vec{p}_{k-1}$
 $\vec{r}^k \in span\{\vec{r}^0, \mathbf{M}\vec{r}^0, \dots, \mathbf{M}^{k-1}\vec{r}^0, \mathbf{M}^k \vec{r}^0\}$
 $\vec{r}^k = \varphi_k(\mathbf{M})\vec{r}^0 \rightarrow k^{\text{th}}$ order polynomial in \mathbf{M}

$$\begin{aligned} \varphi_k(0) \quad \mathbf{M} &= 0 & \mathbf{M}\tilde{x} &= \tilde{b} & 0\tilde{x} &= \tilde{b} & \tilde{r}^k &= \tilde{b} \\ \varphi_k(0) &= \beta_0 \mathbf{I} + \beta_1 \mathbf{M} + \beta_2 \mathbf{M}^2 + \dots + \beta_k \mathbf{M}^k & &= \beta_0 \mathbf{I} & & & & \\ \tilde{b} = \tilde{r}^k &= \varphi_k(0)\tilde{r}^0 & \Rightarrow \varphi_k(0) &= \mathbf{I} & & & & \end{aligned}$$

Residual Minimization

If $\tilde{x}^k \in \text{span}\{\tilde{r}^0, \mathbf{M}\tilde{r}^0, \dots, \mathbf{M}^{k-1}\tilde{r}^0\}$ minimizing $\|\tilde{r}^k\|_2^2$

- 1.) $\tilde{x}^k = \tilde{z}_{k-1}(\mathbf{M})\tilde{r}^0$
 $\tilde{z}_{k-1}(\mathbf{M})$ is the k^{th} order poly minimizing $\|\tilde{r}^k\|_2^2$
- 2.) $\tilde{r}^k = \varphi_k(\mathbf{M})\tilde{r}^0$
 $= \tilde{b} - \mathbf{M}\tilde{x}^k = (\mathbf{I} - \mathbf{M}\tilde{z}_k(\mathbf{M}))\tilde{r}^0$

where $\varphi_k(\mathbf{M})\tilde{r}^0$ is the k^{th} order poly minimizing $\|\tilde{r}^k\|_2^2$ subject to $\varphi_k(0) = \mathbf{I}$

Polynomial Property only a function of solution space and residual minimization

Explicitly construct Krylov subspace:

$$\begin{aligned} \tilde{r}^0 &= \tilde{b} - \mathbf{M}\tilde{x}_0^0 \\ \mathbf{M}\tilde{r}^0 &= \mathbf{M}\tilde{r}^0 \\ \mathbf{M}^2\tilde{r}^0 &= \mathbf{M}\cdot\mathbf{M}\tilde{r}^0 \\ &\vdots \\ \mathbf{M}^k\tilde{r}^0 &= \mathbf{M}\cdot\mathbf{M}^{k-1}\tilde{r}^0 \end{aligned}$$

Orthogonalize & Project

Residual Minimizing Optimality Property:

$$\|\tilde{r}^k\| = \|\varphi_k(\mathbf{M})\tilde{r}^0\| \leq \|\tilde{\varphi}_k(\mathbf{M})\tilde{r}^0\| \leq \|\tilde{\varphi}_k(\mathbf{M})\| \|\tilde{r}^0\|$$

$\tilde{\varphi}_k$ is any k^{th} order poly such that $\tilde{\varphi}_k(0) = \mathbf{I}$

Therefore

Any polynomial which satisfies the constraint can be used to get

an upper bound on $\frac{\|\tilde{r}^k\|}{\|\tilde{r}^0\|} \leq \|\tilde{\varphi}_k(\mathbf{M})\|$

Bound on quantity of residual reduction

$$\|\tilde{r}^{k+1}\| < \epsilon \longrightarrow \mathbf{M}\tilde{x}^{k+1} \cong \tilde{b}$$

$$\|\tilde{x}^{k+1} - \tilde{x}^* = \mathbf{M}^{-1}\tilde{b}\| < \epsilon$$

$$\left. \begin{aligned} \tilde{b} - \mathbf{M}\tilde{x}^{k+1} &= \tilde{r}^{k+1} \\ \tilde{b} - \mathbf{M}\tilde{x}^{k+1} &= 0 \end{aligned} \right\} \mathbf{M}(\tilde{x}^{k+1} - \tilde{x}^*) = -\tilde{r}^{k+1}$$

$$\Rightarrow \|\tilde{x}^{k+1} - \tilde{x}^*\| = \|\mathbf{M}^{-1}\tilde{r}^{k+1}\| \leq \|\mathbf{M}^{-1}\| \|\tilde{r}^{k+1}\| = \|\mathbf{M}^{-1}\| \epsilon$$

$$\tilde{\varphi}_2(x) = 3 + 2x - 4x^2$$

$$\begin{aligned} \tilde{\varphi}_2(\mathbf{M}) &= 3\mathbf{I} + 2\mathbf{M} - 4\mathbf{M}^2 = 3\mathbf{U}\mathbf{U}^{-1} + 2\mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1} - 4\mathbf{U}\mathbf{\Lambda}^2\mathbf{U}^{-1} \\ &= \mathbf{U}(3 + 2\mathbf{\Lambda} - 4\mathbf{\Lambda}^2)\mathbf{U}^{-1} \end{aligned}$$

$$\mathbf{M} = \mathbf{U} \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_3 \end{bmatrix} \mathbf{U}^{-1}$$

$\lambda_j \rightarrow \varphi_2(\lambda_j)$

$$\tilde{\varphi}_2(\mathbf{M}) = \mathbf{U} \begin{bmatrix} 3 + 2\lambda_1 - 4\lambda_1^2 & & \\ & 3 + 2\lambda_2 - 4\lambda_2^2 & \\ & & 3 + 2\lambda_3 - 4\lambda_3^2 \end{bmatrix} \mathbf{U}^T$$

Spectral Mapping Theorem.

Given a polynomial $\tilde{\varphi}_p(x) = a_0 + a_1x + \dots + a_px^p$

Apply the polynomial to a matrix $\tilde{\varphi}_p(\mathbf{M}) = a_0\mathbf{I} + a_1\mathbf{M} + \dots + a_p\mathbf{M}^p$

Then $\text{spectrum}(\tilde{\varphi}_p(\mathbf{M})) = \tilde{\varphi}_p(\text{spectrum}(\mathbf{M}))$

Proof.

Note a property of matrix powers

$$\mathbf{M}\mathbf{M} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1}\mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1} = \mathbf{U}\mathbf{\Lambda}^2\mathbf{U}^{-1} \Rightarrow \mathbf{M}^p = \mathbf{U}\mathbf{\Lambda}^p\mathbf{U}^{-1}$$

Apply to the polynomial of the matrix

$$\tilde{\varphi}_p(\mathbf{M}) = a_0\mathbf{U}\mathbf{U}^{-1} + a_1\mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1} + \dots + a_p\mathbf{U}\mathbf{\Lambda}^p\mathbf{U}^{-1}$$

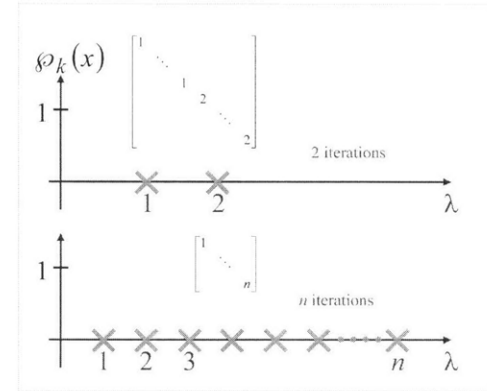
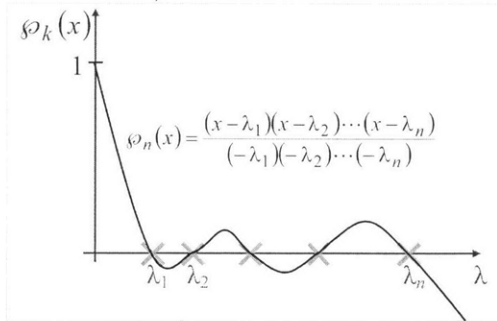
Factoring $\tilde{\varphi}_p(\mathbf{M}) = \mathbf{U} \underbrace{(a_0\mathbf{I} + a_1\mathbf{\Lambda} + \dots + a_p\mathbf{\Lambda}^p)}_{\text{diagonal}} \mathbf{U}^{-1}$

$$\text{spectrum}(\tilde{\varphi}_p(\mathbf{M})) = \text{spectrum}(a_0 I + a_1 \lambda + \dots + a_p \lambda^p) = \tilde{\varphi}_p(\text{spectrum}(\mathbf{M}))$$

$$\begin{bmatrix} a_0 + a_1 \lambda_1 + \dots + a_p \lambda_1^p & & \\ & \ddots & \\ & & a_0 + a_1 \lambda_N + \dots + a_p \lambda_N^p \end{bmatrix}$$

Norm of matrix polynomials.

$$\begin{aligned} \|\tilde{\varphi}_k(\mathbf{M})\| &= \left\| \begin{bmatrix} \uparrow & \uparrow \\ \tilde{u}_1 & \dots & \tilde{u}_N \\ \downarrow & & \downarrow \end{bmatrix} \begin{bmatrix} \tilde{\varphi}_k(\lambda_1) \\ \vdots \\ \tilde{\varphi}_k(\lambda_N) \end{bmatrix} \right\| \\ &\leq \underbrace{\left\| \begin{bmatrix} \uparrow & \uparrow \\ \tilde{u}_1 & \dots & \tilde{u}_N \\ \downarrow & & \downarrow \end{bmatrix} \right\|}_{\text{condition number of } \mathbf{M}'\text{'s eigenspace}} \left\| \begin{bmatrix} \uparrow & \uparrow \\ \tilde{u}_1 & \dots & \tilde{u}_N \\ \downarrow & & \downarrow \end{bmatrix} \begin{bmatrix} \tilde{\varphi}_k(\lambda_1) \\ \vdots \\ \tilde{\varphi}_k(\lambda_N) \end{bmatrix} \right\| \\ &\leq \text{cond}(\mathbf{U}) \max_i \|\tilde{\varphi}_k(\lambda_i)\| \end{aligned}$$



Symmetric positive semidefinite case

If $\mathbf{M} = \mathbf{M}^T$ then

1. \mathbf{M} has orthonormal eigenvectors

$$\text{cond}(\mathbf{U}) = \left\| \begin{bmatrix} \uparrow & \uparrow \\ \tilde{u}_1 & \dots & \tilde{u}_N \\ \downarrow & & \downarrow \end{bmatrix} \right\| = 1$$

$$\Rightarrow \|\tilde{\varphi}_k(\mathbf{M})\| = \max_i \|\tilde{\varphi}_k(\lambda_i)\|$$

2. \mathbf{M} has real eigenvalues

If \mathbf{M} is positive definite, then $\lambda(\mathbf{M}) > 0$

Bound on quantity of residual reduction

$$|r^{k+1}| < \varepsilon \longrightarrow Mx^{k+1} \cong \bar{b}$$

$$\left\{ \begin{array}{l} |x^{k+1}| < \varepsilon \\ x^* = M^{-1}\bar{b} \end{array} \right.$$

$$\left. \begin{array}{l} \bar{b} - Mx^{k+1} = r^{k+1} \\ \bar{b} - Mx^{k+1} = 0 \end{array} \right\} M(x^{k+1} - x^*) = -r^{k+1}$$

$$\Rightarrow |x^{k+1} - x^*| = \|M^{-1}r^{k+1}\| \leq \|M^{-1}\| |r^{k+1}| = \|M^{-1}\| \varepsilon$$

Polynomial Error Bound

Upper bound on residual after k iterations

$$\frac{\|r^k\|}{\|r^0\|} = |\varphi_k(\mathbf{M})| \quad \varphi_k(\mathbf{M}) \text{ is the GCR polynomial of order } k$$

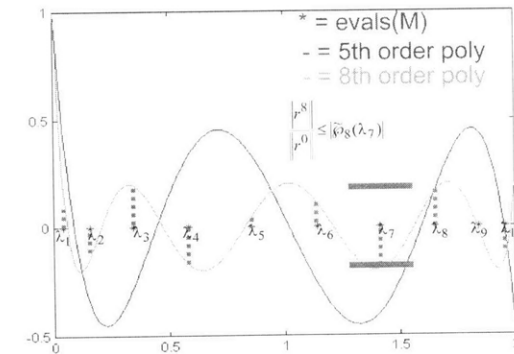
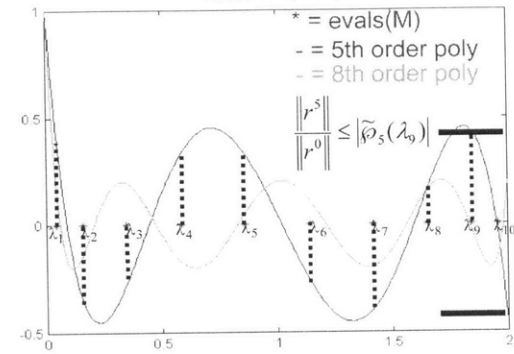
$$\leq |\tilde{\varphi}_k(\mathbf{M})| \quad \tilde{\varphi}_k(\mathbf{M}) \text{ is any } k^{\text{th}} \text{ order polynomial such that } \tilde{\varphi}_k(0) = \mathbf{I}$$

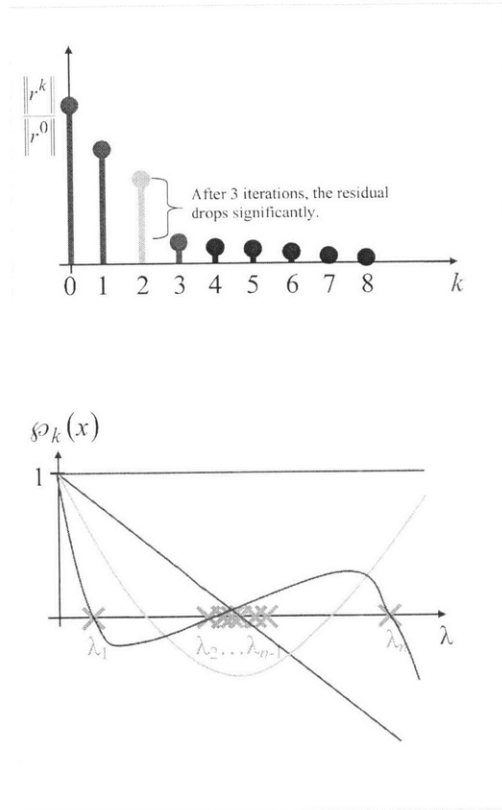
$$\leq \text{cond}(\mathbf{U}) \cdot \max_i |\tilde{\varphi}_k(\lambda_i)| \quad \mathbf{U} \text{ are the eigenvectors of } \mathbf{M}$$

$$\leq \max_i |\tilde{\varphi}_k(\lambda_i)| \quad \lambda_i \text{ are the eigenvalues of } \mathbf{M}$$

if \mathbf{M} is symmetric then $\text{cond}(\mathbf{U}) = 1$

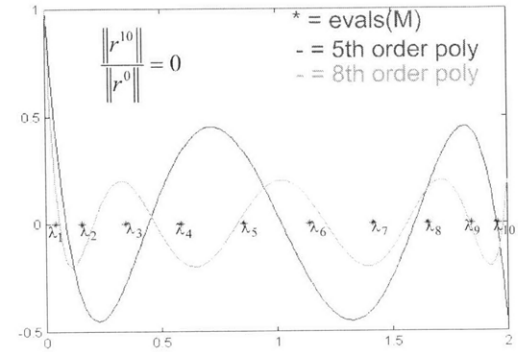
Residual Poly Picture for Heat Conducting Bar Matrix
No loss to air ($n = 10$)



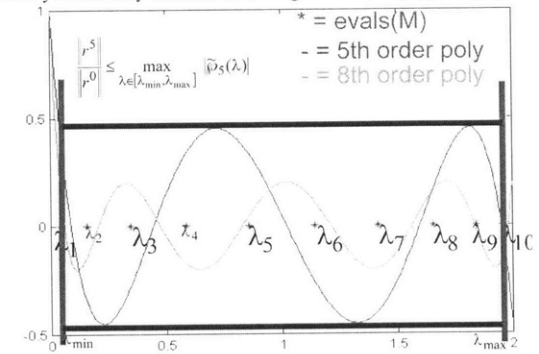


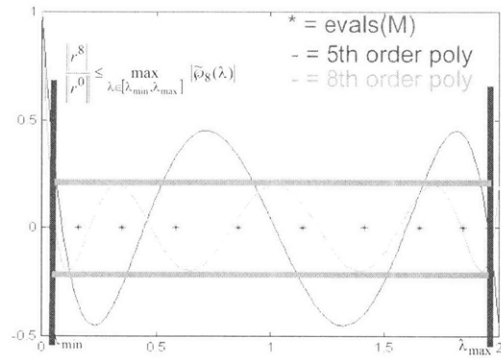
Chebyshev Bounds

How many iterations to converge?



What if you don't know how many and where the eigenvalues are?
 What if you can only estimate their range $\lambda_i \in [\lambda_{\min}, \lambda_{\max}]$





Polynomial Min-Max Problem
 [Convergence for $\mathbf{M} = \mathbf{M}^{-1}$]

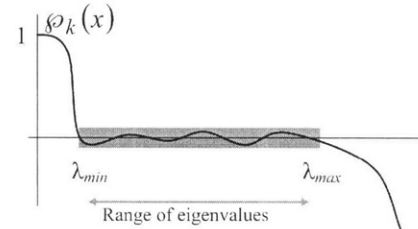
Consider $\lambda(M) \in [\lambda_{\min}, \lambda_{\max}]$, $\lambda_{\min} > 0$

Then a good polynomial (i.e. $\|\tilde{\varphi}_k(\mathbf{M})\|$ is small) can be found by solving the min - max problem

$$\min_{\substack{k^{\text{th}} \text{ order } x \in [\lambda_{\min}, \lambda_{\max}] \\ \text{polyss.t.} \\ \tilde{\varphi}_k(0)=1}} \max_{x \in [\lambda_{\min}, \lambda_{\max}]} |\tilde{\varphi}_k(x)|$$

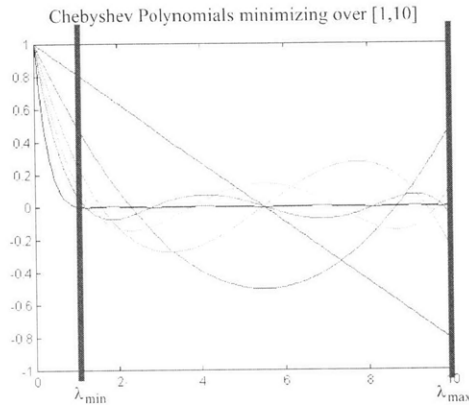
The min-max problem is exactly solved by Chebyshev Polynomials

Find the polynomial that minimizes the value within the range $[\lambda_{\min}, \lambda_{\max}]$ and is equal to 1 at $\lambda=0$



$$\min_{\substack{k^{\text{th}} \text{ order } x \in [\lambda_{\min}, \lambda_{\max}] \\ \text{polyss.t.} \\ \tilde{\varphi}_k(0)=1}} \max_{x \in [\lambda_{\min}, \lambda_{\max}]} |\tilde{\varphi}_k(x)| = \max_{x \in [\lambda_{\min}, \lambda_{\max}]} \left| \frac{C_k \left(1 + 2 \frac{\lambda_{\min} - x}{\lambda_{\max} - \lambda_{\min}} \right)}{C_k \left(1 + 2 \frac{\lambda_{\min}}{\lambda_{\max} - \lambda_{\min}} \right)} \right|$$

The Chebyshev Polynomial $C_k(x) \equiv \cos(k \cos^{-1}(x))$ $x \in [-1, 1]$



Chebyshev Bounds

$$\begin{aligned} \min_{\substack{k^{\text{th}} \text{ order} \\ \text{polys.s.t.} \\ \tilde{\varphi}_k(0)=1}} \max_{x \in [\lambda_{\min}, \lambda_{\max}]} \tilde{\varphi}_k(x) &= \max_{x \in [\lambda_{\min}, \lambda_{\max}]} \frac{C_k \left(1 + 2 \frac{\lambda_{\min} - x}{\lambda_{\max} - \lambda_{\min}} \right)}{C_k \left(1 + 2 \frac{\lambda_{\min}}{\lambda_{\max} - \lambda_{\min}} \right)} \\ &= \frac{1}{\left| C_k \left(1 - 2 \frac{\lambda_{\max}}{\lambda_{\max} - \lambda_{\min}} \right) \right|} \\ &\leq 2 \left(\frac{\sqrt{\lambda_{\max}} - 1}{\sqrt{\lambda_{\min}}} \right)^k \end{aligned}$$

Example.

Given the ratio $\frac{\lambda_{\max}}{\lambda_{\min}} = 4$, what is the residual after k iterations?

$$\|r^k\| \leq 2 \left(\frac{\sqrt{4} - 1}{\sqrt{4} + 1} \right)^k \|r^0\| \Rightarrow \|r^k\| \leq \frac{2}{3^k} \|r^0\|$$

Given the ratio $\frac{\lambda_{\max}}{\lambda_{\min}} = 100$, what is the residual after k iterations?

$$\|r^k\| \leq 2 \left(\frac{\sqrt{100} - 1}{\sqrt{100} + 1} \right)^k \|r^0\| \Rightarrow \|r^k\| \leq 2 \left(\frac{9}{11} \right)^k \|r^0\|$$

Given the ratio $\frac{\lambda_{\max}}{\lambda_{\min}} = 10,000$, what is the residual after k iterations?

$$\|r^k\| \leq 2 \left(\frac{\sqrt{10,000} - 1}{\sqrt{10,000} + 1} \right)^k \|r^0\| \Rightarrow \|r^k\| \leq 2 \left(\frac{99}{101} \right)^k \|r^0\|$$

Which case converges faster?
If the ratio of eigenvalues is large, then the bound on the convergence will be slower.

Chebyshev Result

If $\lambda(\mathbf{M}) \in [\lambda_{\min}, \lambda_{\max}]$, $\lambda_{\min} > 0$

$$\|r^k\| \leq 2 \left(\frac{\sqrt{\lambda_{\max}} - 1}{\sqrt{\lambda_{\min}}} \right)^k \|r^0\|$$

$$\text{cond}(\mathbf{M}) = \|\mathbf{M}\| \|\mathbf{M}^{-1}\| = \lambda_{\max} \cdot \frac{1}{\lambda_{\min}}$$

$$\begin{aligned} \mathbf{M} &= \mathbf{U} \mathbf{\Lambda} \mathbf{U}^{-1} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T & \lambda_{\max} &= \|\mathbf{M}\|_2 \\ \mathbf{M}^{-1} &= \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^{-1} = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^T & \lambda_{\min} &= ? \end{aligned}$$

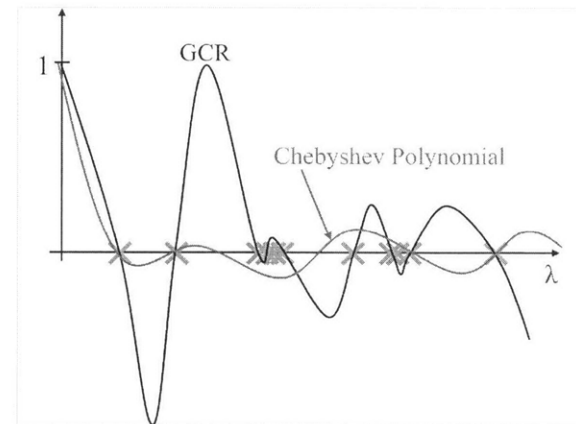
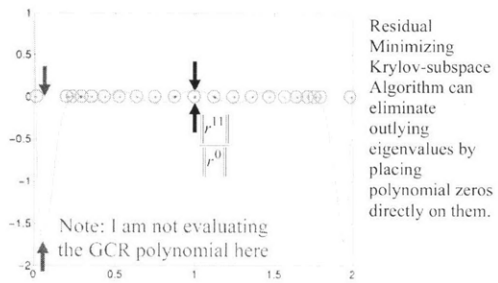
$$\|\mathbf{M}^{-1}\|_2 = \max_i \frac{1}{\lambda_i} = \frac{1}{\lambda_{\min}} \Rightarrow \lambda_{\min} = \left[\|\mathbf{M}^{-1}\|_2 \right]^{-1}$$

$$\frac{\lambda_{\max}}{\lambda_{\min}} = \|\mathbf{M}\|_2 \|\mathbf{M}^{-1}\|_2 = \text{cond}(\mathbf{M})$$

Examples.

$cond(\mathbf{M}) = 100$	$\left \frac{r^k}{r^0} \right \leq 2 \left(\frac{10-1}{10+1} \right)^k = 2 \left(\frac{9}{11} \right)^k$	At each iteration, improves by ~10%
$cond(\mathbf{M}) = 10,000$	$\left \frac{r^k}{r^0} \right \leq 2 \left(\frac{100-1}{100+1} \right)^k$	Worst Case (bound)
$cond(\mathbf{M}) = 10^{12}$	$\left \frac{r^k}{r^0} \right \leq 2 \left(\frac{10^6-1}{10^6+1} \right)^k$	At each iteration, improves by a very, very small amount.

Heat Conducting Bar example
 Sometimes GCR can do much better than Chebyshev bound



Gershgorin Circle Theorem

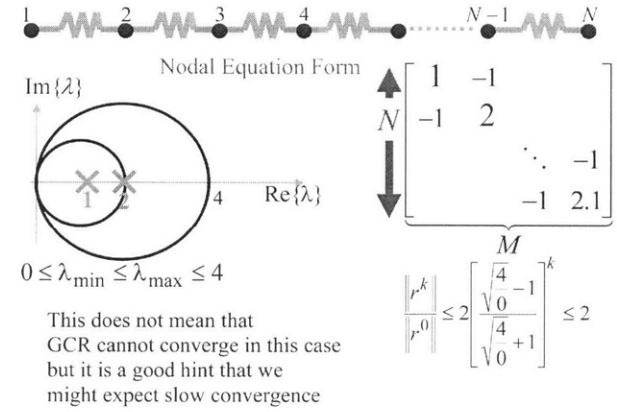
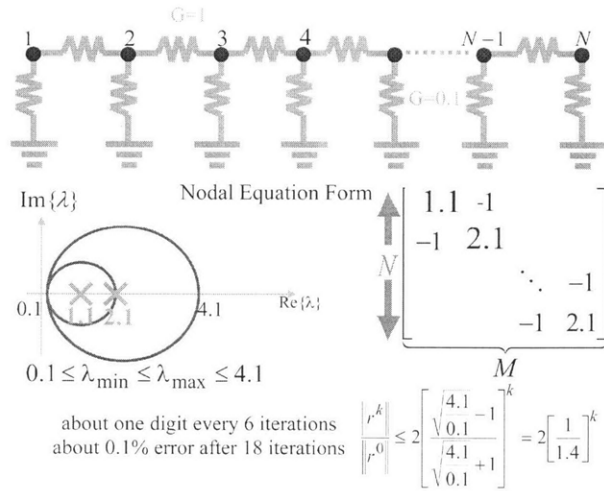
❖ Theorem Statement

Given a matrix $\mathbf{M} = \begin{bmatrix} m_{1,1} & \dots & m_{1,N} \\ \vdots & \ddots & \vdots \\ m_{N,1} & \dots & m_{N,N} \end{bmatrix}$

For each eigenvalue of \mathbf{M} there exists an $i, 1 < i < N$ such that

$$|\lambda - m_{i,i}| \leq \sum_{j \neq i} |m_{i,j}|$$

We say that the eigenvalues are contained in the union of the Gershgorin circles



INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 10.
Preconditioners

TODAY'S OUTLINE:

- ❖ Preconditioners
 - Diagonal Preconditioner
 - Blockdiagonal Preconditioner
 - Incomplete Factorization Preconditioner
- ❖ GCR for different Right Hand Sides
 - Recycling the Krylov subspace

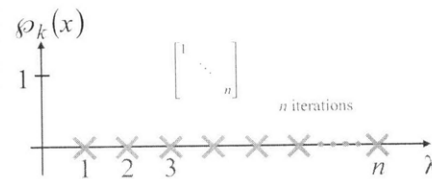
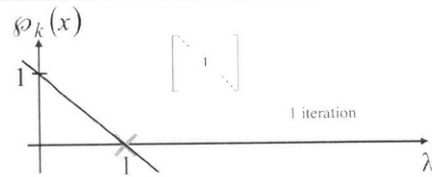
PRECONDITIONERS

Diagonal Preconditioner

❖ Diagonal Example

$$\begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 1 & 0 \\ 0 & \dots & 0 & 0 & 1 \end{bmatrix} \quad \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 2 & 0 & \dots & 0 \\ 0 & 0 & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & N-1 & 0 \\ 0 & \dots & 0 & 0 & N \end{bmatrix}$$

For which problem will GCR converge faster?



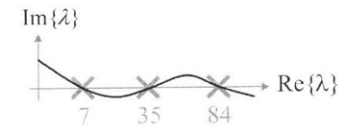
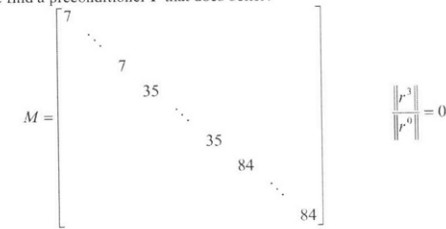
Suppose $M\bar{x} = \bar{b}$ converges slowly

Try $PM\bar{x} = P\bar{b}$ for some P (left pre-conditioner matrix)

e.g. How many iterations for this diagonal M to converge?

$$M = D = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 2 & 0 & \dots & 0 \\ 0 & 0 & 3 & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & 0 & N \end{bmatrix} \quad P = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & \frac{1}{2} & 0 & \dots & 0 \\ 0 & 0 & \frac{1}{3} & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & 0 & \frac{1}{N} \end{bmatrix}$$

Can we find a preconditioner P that does better?



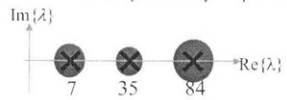
❖ Diagonally Dominant Example

How many iterations for this diagonally dominant M to "get very close" to the exact answer?

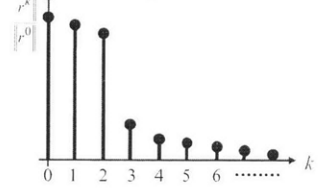
$$M = D + \delta = \begin{bmatrix} 7 & * & * & * & * & * & * & * \\ * & \ddots & * & * & * & * & * & * \\ * & * & 7 & * & * & * & * & * \\ * & * & * & 35 & * & * & * & * \\ * & * & * & * & \ddots & * & * & * \\ * & * & * & * & * & 35 & * & * \\ * & * & * & * & * & * & 84 & * \\ * & * & * & * & * & * & * & \ddots \\ * & * & * & * & * & * & * & * & 84 \end{bmatrix}$$

* is a very small entry compared to diagonals

$$\left| \frac{r^3}{r^0} \right| \approx \text{small}$$



Can we find a preconditioner that does even better?

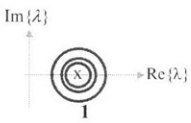


❖ Diagonal Preconditioner

Try as preconditioner the inverse of the diagonal

$$D^{-1}(D + \delta) = \begin{bmatrix} 1 & * & * & * & * & * & * & * \\ * & \ddots & * & * & * & * & * & * \\ * & * & 1 & * & * & * & * & * \\ * & * & * & 1 & * & * & * & * \\ * & * & * & * & \ddots & * & * & * \\ * & * & * & * & * & 1 & * & * \\ * & * & * & * & * & * & 1 & * \\ * & * & * & * & * & * & * & \ddots \\ * & * & * & * & * & * & * & * & 1 \end{bmatrix}$$

$$\left| \frac{r^1}{r^0} \right| \approx \text{small}$$



❖ General Idea

Suppose $M\bar{x} = \bar{b}$ converges slowly

Try $PM\bar{x} = P\bar{b}$ for some P (left pre-conditioner matrix)

If $PM = I$ then convergence happens in one step

however $P = M^{-1}$ is VERY hard to compute

If $PM \approx I$ then we hope convergence happens in "few" steps

Any general idea for picking P ?

Pick \tilde{M} such that :

e.g. if M is diagonally dominant

a) $\tilde{M} \approx M$

its diagonal $\tilde{M} = \text{diag}(M)$ is :

b) \tilde{M} is easy to invert or factor

- a good approximation of M

Preconditioner : $P = \tilde{M}^{-1}$

- and easy to invert

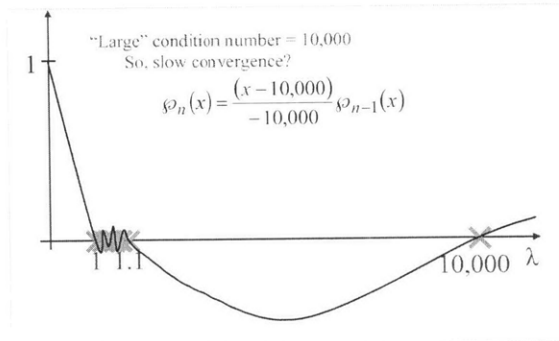
Let $A = D + A_{nd}$

$$\begin{bmatrix} \diagdown & \\ & \diagdown \end{bmatrix}$$

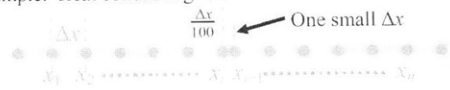
Apply GCR to

$$(D^{-1}A)\bar{x} = (I + D^{-1}A_{nd})\bar{x} = D^{-1}\bar{b}$$

- The inverse of a diagonal is cheap to compute
- Usually improves convergence



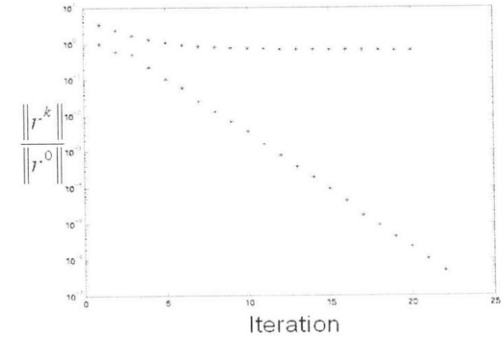
Example. Heat conducting bar



$$\begin{bmatrix}
 2+\gamma & -1 & & & & & \\
 -1 & 2+\gamma & & & & & \\
 & & \ddots & & & & \\
 & & & -1 & 1+\gamma+100 & -100 & \\
 & & & -100 & 1+\gamma+100 & -1 & \\
 & & & & -1 & \ddots & -1 \\
 & & & & & -1 & 2+\gamma
 \end{bmatrix}
 \begin{bmatrix}
 \hat{u}_1 \\
 \vdots \\
 \hat{u}_n
 \end{bmatrix}
 =
 \begin{bmatrix}
 f(x_1) \\
 \vdots \\
 f(x_n)
 \end{bmatrix}$$

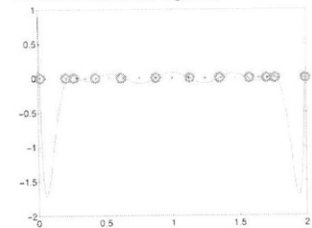
$\lambda_{\max} > 100$
 λ_{\min}

For the heat conducting bar, which convergence curve is GCR?



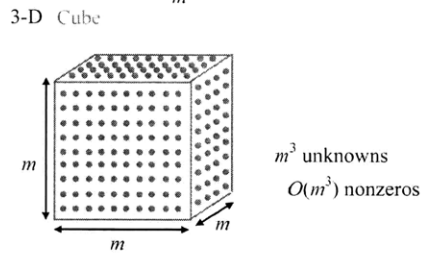
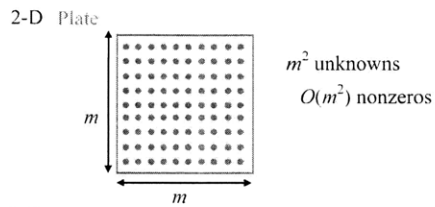
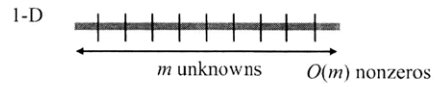
Heat Conducting Bar example continued...

Preconditioned Matrix Eigenvalues



Residual Minimizing Krylov-subspace Algorithm can eliminate outlying eigenvalues by placing polynomial zeros directly on them.

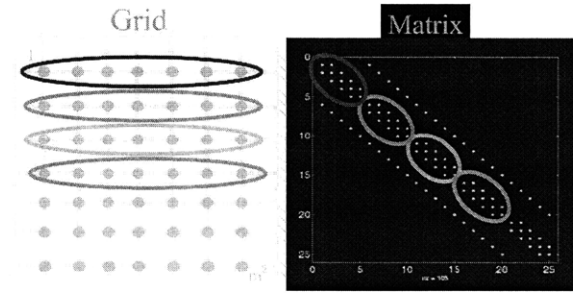
Heat Flow Comparison



Dimension	Dense GE	Sparse GE	GCR
1	$O(m^3)$	$O(m)$	$O(m^2)$
2	$O(m^6)$	$O(m^3)$	$O(m^3)$
3	$O(m^9)$	$O(m^6)$	$O(m^4)$

Blockdiagonal Preconditioner

❖ Line Schemes

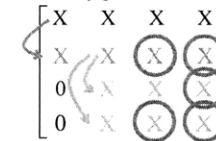


Tridiagonal Matrices factor quickly

Incomplete Factorization Preconditioner

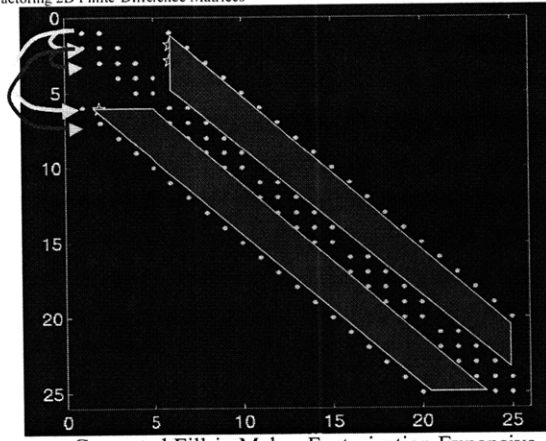
❖ Fill-Ins Example

Fill-ins Propagate



Fill-ins from Step 1 results in Fill-ins in step 2.

❖ Factoring 2D Finite-Difference Matrices



Generated Fill-in Makes Factorization Expensive

❖ Key Idea

- Throw away fill-ins
 - Throw away all fill-ins
 - Throw away only fill-ins with small values
 - Throw away fill-ins produced by other fill-ins
 - Throw away fill-ins produced by fill-ins of other fill-ins, etc.

○ Key Idea

Pick \tilde{M} such that :

a.) Pick $\tilde{M} \approx M$

b.) \tilde{M} is easy to factor : $\tilde{M} = \tilde{L}\tilde{U}$

Use as preconditioner : $P = \tilde{M}^{-1}$

$$M\bar{x} = \bar{b}$$

$$\tilde{M}^{-1}M\bar{x} = \tilde{M}^{-1}\bar{b}$$

But do NOT calculate $P = \tilde{M}^{-1}$!!!!!!!!!!!!!!!!!!!!

OR do NOT calculate $\tilde{M}^{-1}M$!!!!!!!!!!!!!!!!!!!!

$$\tilde{M}^{-1}M\bar{x} = \tilde{M}^{-1}\bar{b} \quad \tilde{M} = \tilde{L}\tilde{U}$$

$$\bar{b} = \tilde{M}^{-1}\bar{b}$$

solve for \bar{b} : $\tilde{L}\tilde{U}\bar{b} = \bar{b}$

At each iteration of GCR we need : $\tilde{M}^{-1}M r^k = \tilde{r}^k$

Then solve for \tilde{r}^k : $\tilde{L}\tilde{U}\tilde{r}^k = M r^k$

$\tilde{M}^{-1}\bar{b} = \bar{\tilde{b}}$	$\tilde{M} = \tilde{L}\tilde{U}$	$\tilde{M}\bar{\tilde{b}} = \bar{b}$	Solve $\tilde{L}\bar{y} = \bar{\tilde{b}} \rightarrow \bar{y}$
		$\tilde{L}\tilde{U}\bar{\tilde{b}} = \bar{b}$	Solve $\tilde{U}\bar{\tilde{b}} = \bar{y} \rightarrow \bar{\tilde{b}}$

$\tilde{M}^{-1}M \rightarrow ?$ using GCR

$\tilde{M}^{-1}M\bar{x} = \bar{\tilde{b}}$

x-space

$w^0 = \tilde{b} = \tilde{r}^0$

M-space

$\tilde{M}^{-1}M r^k \rightarrow \tilde{r}^k$

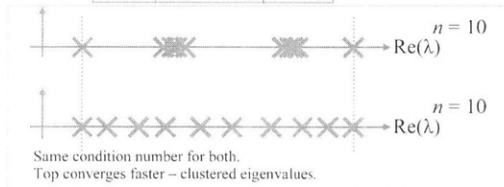
$$\tilde{M} = \tilde{L}\tilde{U}$$

$$\tilde{M}^{-1}M r^k = \tilde{r}^k \Rightarrow \tilde{M} \tilde{r}^k = M r^k$$

$$\tilde{L}\tilde{U} \tilde{r}^k = M r^k$$

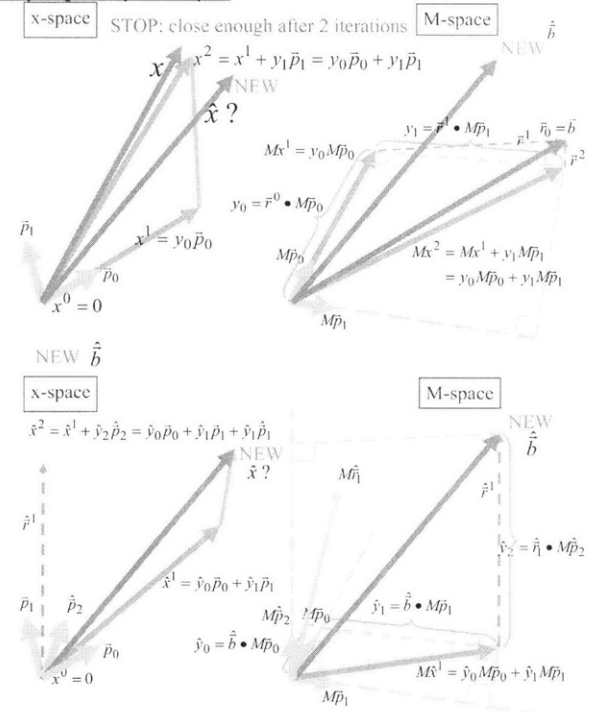
- ❖ Complete Algorithm
 - Pick $\tilde{\mathbf{M}}$ such that :
 - a.) Pick $\tilde{\mathbf{M}} \approx \mathbf{M}$
 - b.) $\tilde{\mathbf{M}}$ is easy to factor : $\tilde{\mathbf{M}} = \tilde{\mathbf{L}}\tilde{\mathbf{U}}$
 - Factor $\tilde{\mathbf{M}} = \tilde{\mathbf{L}}\tilde{\mathbf{U}}$
 - Calculate new RHS solving for \tilde{b} : $\tilde{\mathbf{L}}\tilde{\mathbf{U}}\tilde{b} = \tilde{b}$
 - At the k^{th} step of GCR:
 - Calculate $\mathbf{M}r^k$
 - solve for \tilde{r}^k : $\tilde{\mathbf{L}}\tilde{\mathbf{U}}\tilde{r}^k = \mathbf{M}r^k$
- ❖ GCR Summary
 - Comparison CGR & GE
 - Making GCR converge fast can be very problem specific: pre-conditioning
 - Assume we have a good pre-conditioner (e.g. # iterations $k < 10,20$)

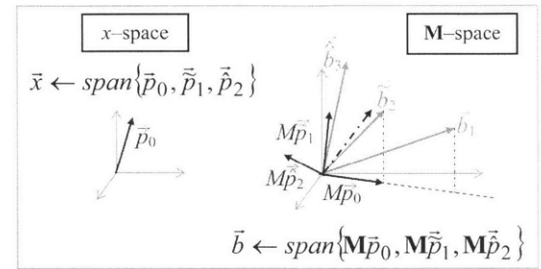
	Sparse	Dense
GCR	$O(n)$	$O(n^2)$
GE	$O(n^{1.2-1.8})$	$O(n^3)$



- When would one use GCR?
 - Need to have fast convergence rate:
 1. $cond(\mathbf{M})$ is not large
 2. or eigenvalues of \mathbf{M} are clustered in few groups
 3. or we have a good "pre-conditioner"
 - And one of the following:
 1. need fewer than 16 digits of precision
 2. or \mathbf{M} is dense
 3. or have a fast matrix-vector product algorithm
- And what about if you need to re-solve a system with different RHS?

GCR FOR DIFFERENT RIGHT HAND SIDES
 Recycling the Krylov Subspace





INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 11.

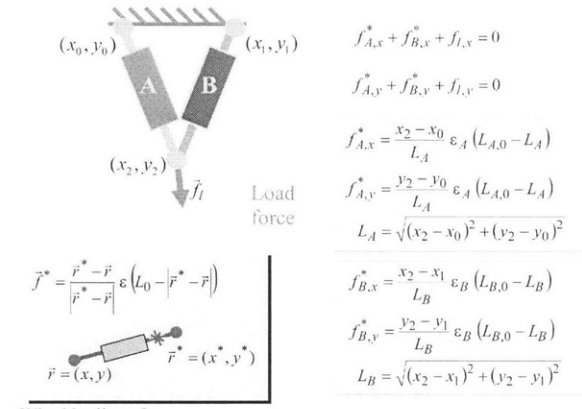
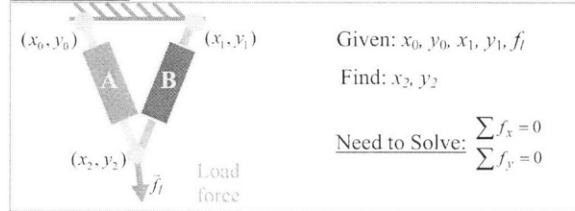
1-D Nonlinear Solution Methods

TODAY'S OUTLINE:

- ❖ Nonlinear Problems
 - Struts and Circuit Example
- ❖ Richardson and Linear Convergence
- ❖ Newton's Method
 - Derivation of the basic algorithm
 - Convergence Analysis
 - Multidimensional Newton's Method

NONLINEAR PROBLEMS

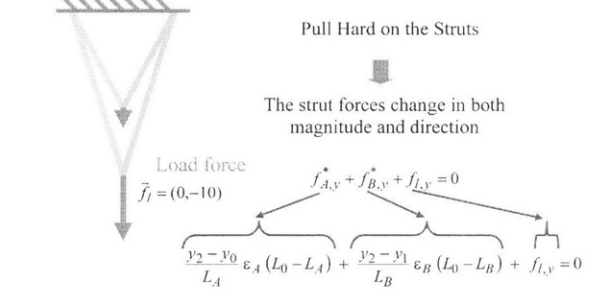
Strut Example

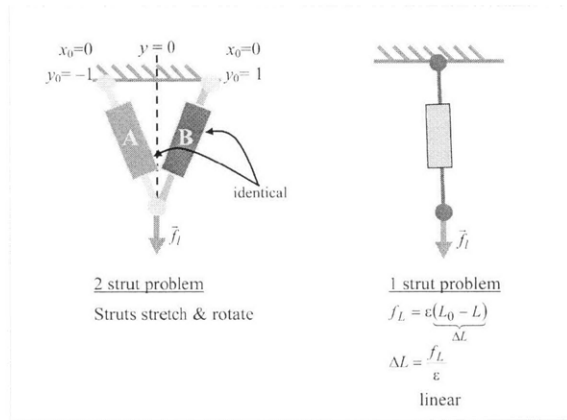


$$\vec{f}^* = \frac{\vec{f}^* - \vec{f}}{|\vec{f}^* - \vec{f}|} \epsilon (L_0 - |\vec{f}^* - \vec{f}|)$$

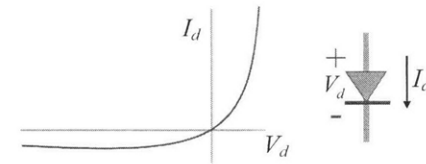
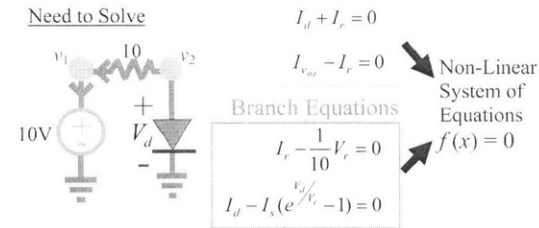
$\vec{f} = (x, y)$ $\vec{f}^* = (x^*, y^*)$

Why Nonlinear?





Circuit Example

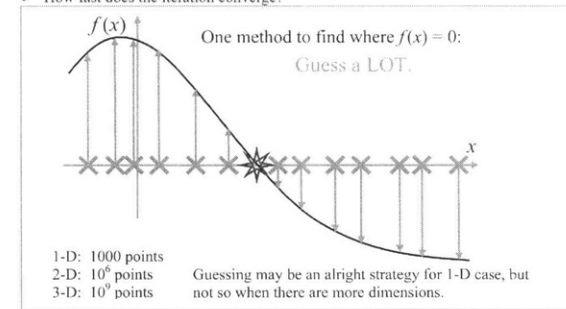


Solve Iteratively

Hard to find analytical solution for $f(x) = 0$
Solve iteratively
guess a solution $x^0 = x_0$
repeat for $k = 0, 1, 2, \dots$
 $x^{k+1} = W(x^k)$
until $f(x^{k+1}) \approx 0$

Ask

- Does the iteration converge to correct solution?
- How fast does the iteration converge?



RICHARDSON AND LINEAR CONVERGENCE

Richardson Iteration

❖ Definition

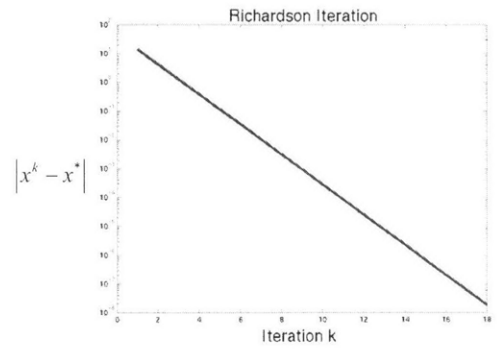
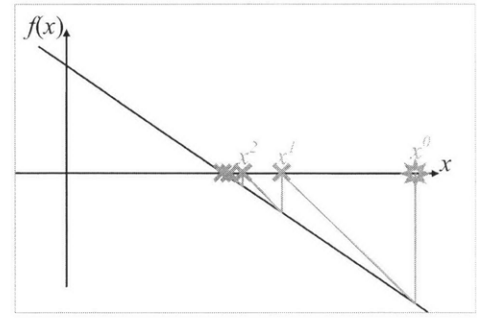
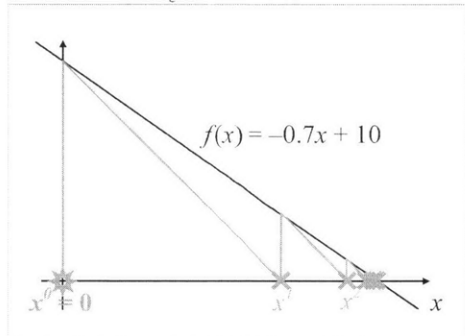
Richardson Iteration Definition

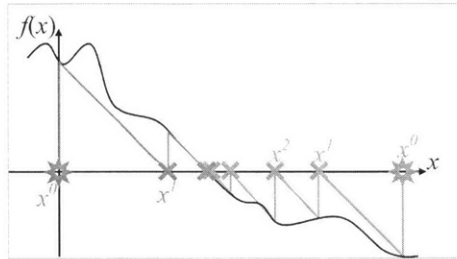
$x^{k+1} = x^k + f(x^k)$

An iteration stationary point is a solution

$x^{k+1} = x^k$
 $\Rightarrow f(x^k) = 0$
 $\Rightarrow x^k = x^*$ (Solution)

❖ Example 1
 $f(x) = -0.7x + 10$
 Start with $x^0 = 0$
 $x^1 = x^0 + f(x^0) = 0 + 10 = 10$
 $x^2 = x^1 + f(x^1) = 10 + (-7 + 10) = 13$
 $x^3 = x^2 + f(x^2) = 13 + (-0.7 \cdot 13 + 10) = 13.9$
 $x^4 = x^3 + f(x^3) = 13.9 + (-0.7 \cdot 13.9 + 10) = 14.17$
 $x^5 = 14.25$
 $x^6 = 14.27$
 $x^7 = 14.28$
 $x^8 = 14.28 \leftarrow$ Converged!





❖ Example 2

$$f(x) = 2x + 10$$

Start with $x^0 = 0$

$$x^1 = x^0 + f(x^0) = 0 + 10 = 10$$

$$x^2 = x^1 + f(x^1) = 10 + (2 \cdot 10 + 10) = 40$$

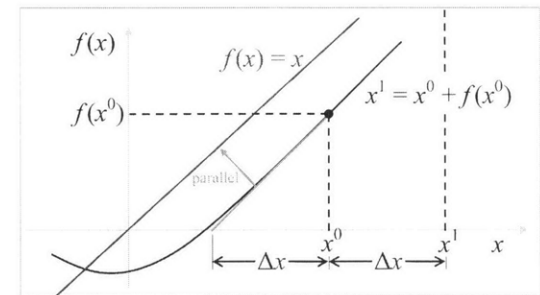
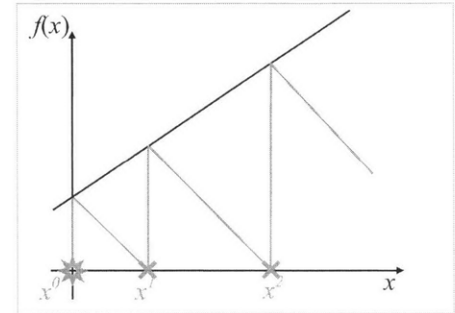
$$x^3 = x^2 + f(x^2) = 40 + (2 \cdot 40 + 10) = 130$$

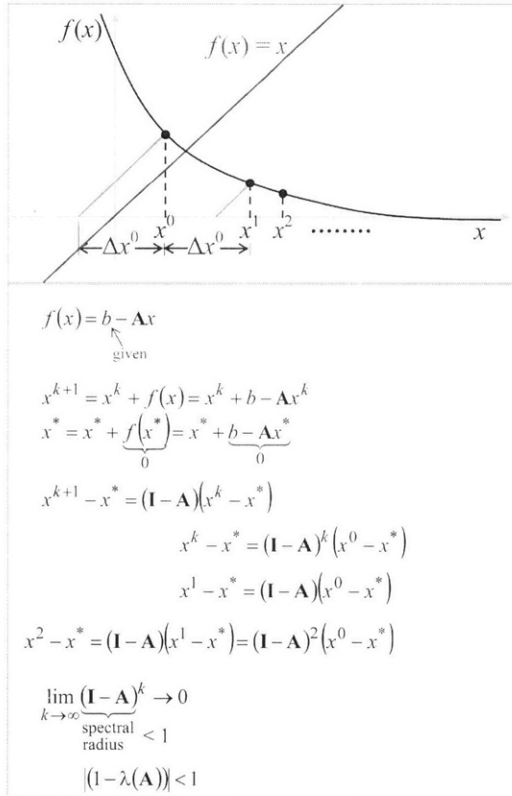
$$x^4 = x^3 + f(x^3) = 130 + (2 \cdot 130 + 10) = 400$$

⋮

⋮

No Convergence!!





❖ Convergence

$x^{k+1} = x^k + f(x^k)$
 $x^* = x^* + \underbrace{f(x^*)}_0$

$f(x) = f(x^*) + \frac{df}{dx}(x^*)(x - x^*) + \frac{1}{2} \frac{d^2 f}{dx^2}(x^*)(x - x^*)^2 + \dots$
 $x^{k+1} - x^* = x^k - x^* + f(x^*) + \frac{df}{dx}(x^*)(x^k - x^*) + \frac{1}{2} \frac{d^2 f}{dx^2}(x^*)(x^k - x^*)^2 + \dots$

$|x^{k+1} - x^*| = \left| \underbrace{1 + \frac{df}{dx}(x^*) + \frac{1}{2} \frac{d^2 f}{dx^2}(x^*)(x^k - x^*) + \dots}_{K^k} \right| |x^k - x^*|$

Find a bound: $K^k \leq K^* \quad x^k \rightarrow x^*$

$K^k \leq \left| 1 + \frac{df}{dx}(x^*) + \epsilon \right|$?

$|x^{k+1} - x^*| = K^* |x^k - x^*|$ for large k

If $K^* < 1$ then Richardson converges linearly.

$ x^{k+1} - x^* \leq \frac{1}{2} x^k - x^* $	$ x^k - x^* = 10$
	$ x^{k+1} - x^* \leq \frac{1}{2} \cdot 10 = 5$
	$ x^{k+10} - x^* \leq \left(\frac{1}{2}\right)^{10} \cdot 10$

○ Setup

Iteration Equation $x^{k+1} = x^k + f(x^k)$

Exact Solution $x^* = x^* + \underbrace{f(x^*)}_0$

Computing Differences $x^{k+1} - x^* = x^k - x^* + \underbrace{f(x^k) - f(x^*)}_{\text{Need to Estimate}}$

$$|x^{k+1} - x^*| = g(x^k - x^*)$$

$$|x^{k+1} - x^*| = \gamma |x^k - x^*|$$

$$\gamma \leq 1$$

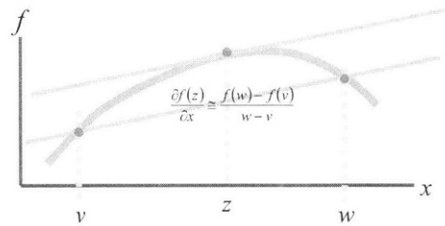
$$K^* = \left| 1 + \frac{\partial f}{\partial x}(x^*) + \frac{1}{2} \frac{\partial^2 f}{\partial x^2}(x - x^*) + \varepsilon \right|$$

Example 1. $f(x) = -0.7x + 10$
 $K^* = |1 + (-0.7) + 0 + \varepsilon| = |0.3 + \varepsilon|$

Example 2. $f(x) = 2x + 10$
 $K^* = |1 + 2 + 0 + \varepsilon| = |3 + \varepsilon|$
 No convergence

o Mean Value Theorem

$$f(w) - f(v) = \frac{\partial f(z)}{\partial x} (w - v) \quad z \in [v, w]$$



Use Mean Value Theorem

Iteration Equation $x^{k+1} = x^k + f(x^k)$
 Exact Solution $x^* = x^* + \underset{=0}{f(x^*)}$

Computing Differences

$$x^{k+1} - x^* = x^k - x^* + f(x^k) - f(x^*) = \left(1 + \frac{\partial f(\bar{x})}{\partial x} \right) (x^k - x^*)$$

$$\bar{x} \in [x^k, x^*]$$

$$|x^{k+1} - x^*| = \left| x^k - x^* + \frac{\partial f}{\partial x}(\bar{x})(x^k - x^*) \right|$$

$$= \left| 1 + \frac{\partial f}{\partial x}(\bar{x}) \right| (x^k - x^*)$$

$$\leq \gamma (x^k - x^*)$$

$$\leq \gamma \gamma (x^{k-1} - x^*)$$

$$\leq \gamma^k (x^0 - x^*)$$

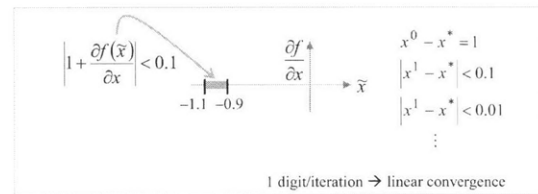
o Richardson Theorem

IF $\left| 1 + \frac{\partial f(\bar{x})}{\partial x} \right| \leq \gamma < 1$ for all \bar{x} s.t. $|\bar{x} - x^*| < \delta$

AND $|x^0 - x^*| < \delta$

THEN $|x^{k+1} - x^*| \leq \gamma |x^k - x^*|$

OR $\lim_{k \rightarrow \infty} |x^{k+1} - x^*| = \lim_{k \rightarrow \infty} |x^0 - x^*| = 0$



Example 1. $f(x) = -0.7x + 10$
 $\gamma = 1 + \frac{\partial f}{\partial x} = 1 + (-0.7) = 0.3 < 1$
 Convergence

Example 2. $f(x) = 2x + 10$
 $\gamma = 1 + \frac{\partial f}{\partial x} = 1 + 2 = 3 > 1$
 No convergence

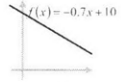
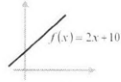
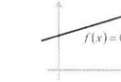
- ❖ Problems
 - Convergence is only linear
 - $x, f(x)$ not in the same units:
 - x is a voltage, $f(x)$ a current in circuits
 - x is a displacement, $f(x)$ a force in struts
 - adding two different physical quantities

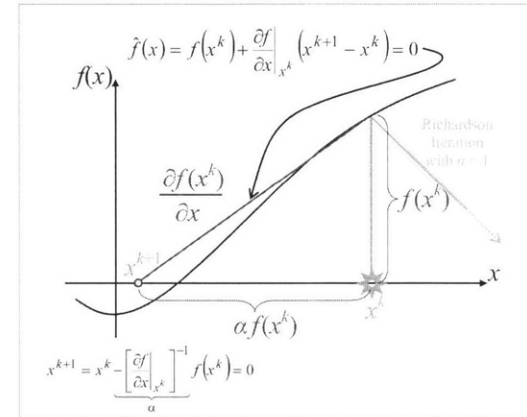
$$x^{k+1} \leftarrow x^k + \alpha f(x^k)$$

V voltage	i current
T temperature	h heat flow
x position	F force

- ❖ Advantage
 - Simple – only need to calculate $f(x)$ and update

$$x^{k+1} \leftarrow x^k + \alpha f(x^k)$$

Example 1.  $\alpha = 1$ $\gamma = 1 - 0.7 = 0.3$	Example 2.  $\alpha = 1$ → No convergence!! $\alpha = -1$ $\gamma = 1 - 2 = 1$ → Still no convergence!! $\alpha = ?$	Example 3.  $\alpha = -1$ $\gamma = 1 - 0.7 = 0.3$
--	---	---



NEWTON'S METHOD

Derivation of the Basic Algorithm

- ❖ 1-D Reminder
 - Newton Idea

Problem: Find x^* such that $f(x^*) = 0$

Approximate $f(x)$ with its Taylor series about x^k :

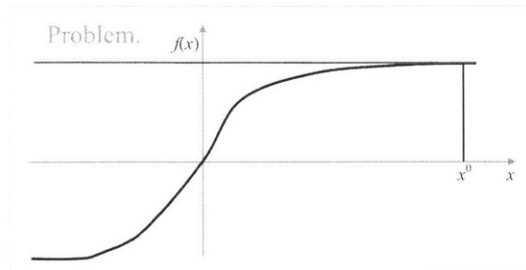
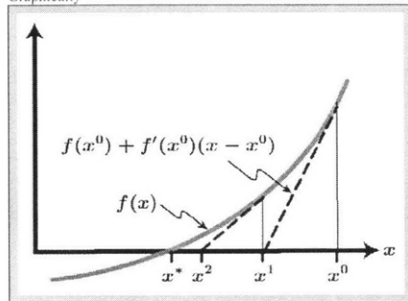
$$\hat{f}(x) \cong f(x^k) + \frac{\delta f}{\delta x}(x^k)(x - x^k) \text{ a straight line}$$

Find the solution of the approximation x^{k+1} :

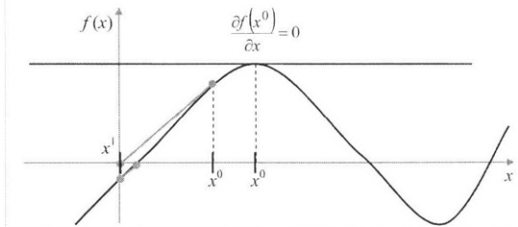
$$f(x^k) + \frac{\delta f}{\delta x}(x^k)(x^{k+1} - x^k) = 0$$

$$\Rightarrow x^{k+1} = x^k - \left[\frac{\partial f}{\partial x}(x^k) \right]^{-1} f(x^k)$$

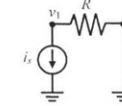
Graphically



Problem.



Example.



$$f(v_1) = \frac{1}{R}v_1 + i_s$$

$$\frac{\partial f}{\partial v} = \frac{1}{R}$$

$$\left(\frac{\partial f}{\partial v}\right)^{-1} (f(v)) = \underbrace{R(f(v))}_{\text{Voltage}}$$

Resistance Current

Newton Algorithm

$x^0 = \text{Initial Guess}, k = 0$

Repeat {

$$x^{k+1} = x^k - \left[\frac{\partial f}{\partial x}(x^k)\right]^{-1} f(x^k)$$

$k = k + 1$

} Until ?

$$\|x^{k+1} - x^k\| < \text{threshold}$$

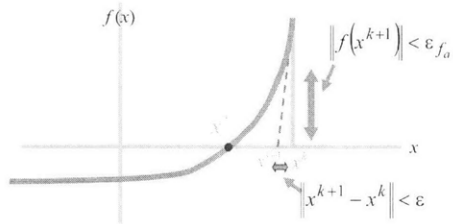
$$\|f(x^{k+1})\| < \text{threshold}$$

Convergence Analysis

❖ Convergence Checks

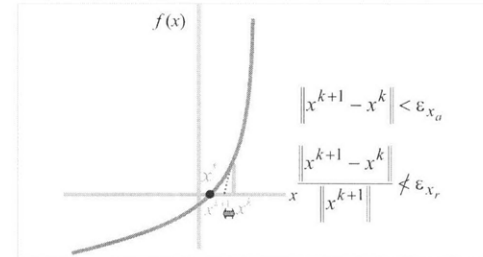
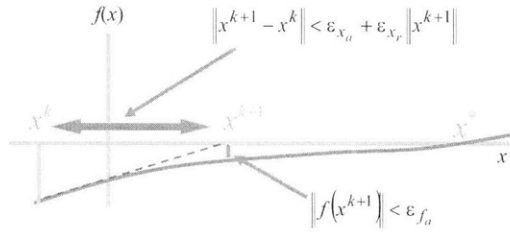
- o $f(x)$

Need an " $f(x)$ " check to avoid false convergence.



- o Δx

Need a " Δx " check to avoid false convergence.

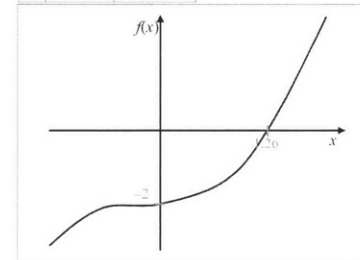


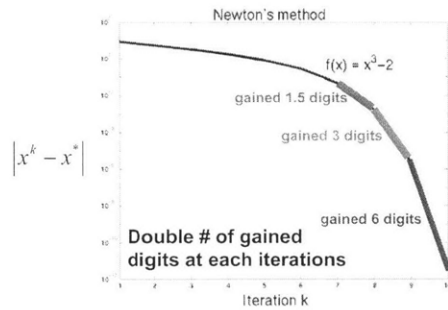
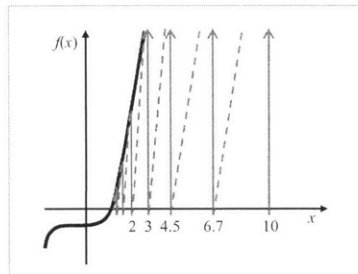
❖ Convergence Example

$f(x) = x^3 - 2 \quad x^* = \sqrt[3]{2} \approx 1.259921$

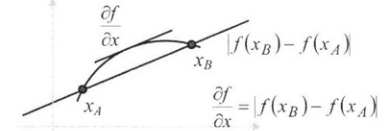
k	x^k	$ x^k - x^* $
0	10.0	8.740
1	6.673333	5.413
•	•	•
8	1.261665	1.744e-03
9	1.259924	2.410e-06
10	1.259921	4.609e-12

Asymptotically,
 $|x^{k+1} - x^*| \approx C |x^k - x^*|^\alpha$
 $C = 0.7951$
 $\alpha = 2.000$ **Quadratic**





❖ Taylor Expansion



Taylor Series Expansion

$$f(x) = f(x^k) + \frac{\partial}{\partial x} f(x^k)(x - x^k) + \frac{1}{2!} \frac{\partial^2}{\partial x^2} f(x^k)(x - x^k)^2 + \dots$$

$$\dots + \frac{1}{n!} \frac{\partial^n}{\partial x^n} f(x^k)(x - x^k)^n + O((x - x^k)^{n+1})$$

$$f(x) = f(x^k) + \frac{\partial}{\partial x} f(x^k)(x - x^k) + \frac{1}{2!} \frac{\partial^2}{\partial x^2} f(x^k)(x - x^k)^2 + \dots$$

$$\dots + \frac{1}{n!} \frac{\partial^n}{\partial x^n} f(x^k)(x - x^k)^n + \frac{1}{(n+1)!} \frac{\partial^{n+1}}{\partial x^{n+1}} f(\tilde{x})(x - x^k)^{n+1}$$

$$\tilde{x} \in [x, x^k]$$

Exact Taylor expansion about the iteration x^k evaluated at the solution x^*

$$0 = f(x^*) = f(x^k) + \frac{df}{dx}(x^k)(x^* - x^k) + \frac{1}{2} \frac{d^2f}{dx^2}(\tilde{x})(x^* - x^k)^2 \quad (1)$$

some $\tilde{x} \in [x^k, x^*]$

Mean Value Theorem
truncates Taylor Series

Approximate Taylor expansion about the solution x^k evaluated at the next iteration x^{k+1}

$$0 = f(x^{k+1}) + \frac{df}{dx}(x^k)(x^{k+1} - x^k) \quad (2)$$

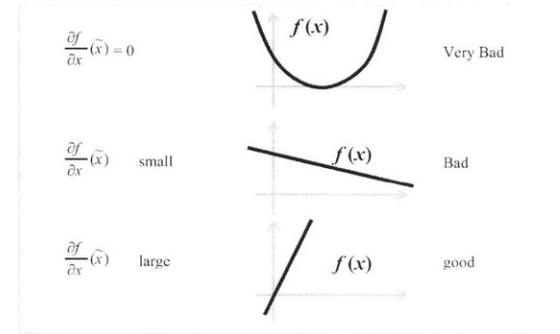
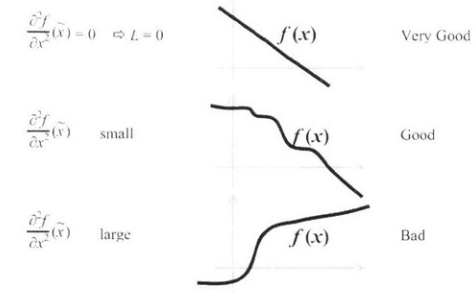
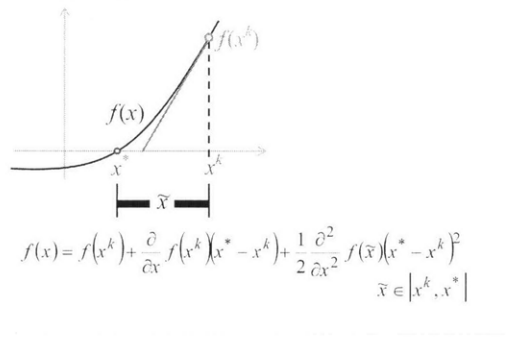
Subtracting (1) from (2)

$$\frac{df}{dx}(x^k)(x^{k+1} - x^*) = \frac{d^2f}{dx^2}(\tilde{x})(x^* - x^k)^2$$

Dividing through

$$(x^{k+1} - x^*) = \left[\frac{df}{dx}(x^k) \right]^{-1} \frac{d^2f}{dx^2}(\tilde{x})(x^* - x^k)^2$$

Suppose $\left| \left[\frac{df}{dx}(x) \right]^{-1} \frac{d^2f}{dx^2}(x) \right| \leq L$ for all x
 then $x^{k+1} - x^* \leq L |x^k - x^*|^2$
 Convergence is quadratic if L is bounded.



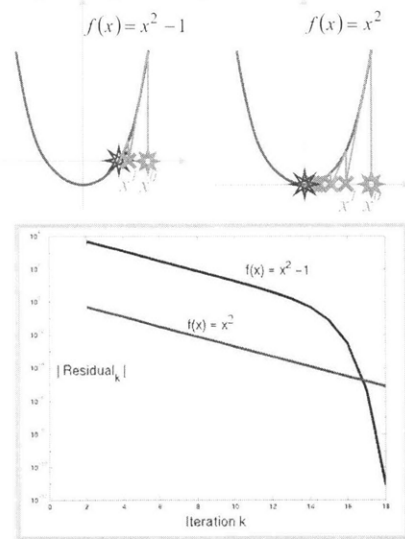
❖ Example 1
 $f(x) = x^2 - 1 = 0$, find x ($x^* = 1$)
 $\frac{df}{dx}(x^k) = 2x^k$
 $2x^k(x^{k+1} - x^k) = -[(x^k)^2 - 1]$
 $2x^k(x^{k+1} - x^*) + \frac{1}{2} 2(x^* - x^k) = -[(x^k)^2 - 1]$ or $(x^{k+1} - x^*) = \frac{1}{2x^k}(x^k - x^*)^2$
 convergence is quadratic
 $x^0 = 1.1$ $x^* = 1$
 $(x^1 - x^*) = \frac{1}{2 \cdot 1.1} (1.1^0 - 1)^2 = \frac{1}{2.2} = 0.0045$

❖ Example 2
 $f(x) = x^2 = 0$, $x^* = 0$
 $\frac{df}{dx}(x^k) = 2x^k$
 $\Rightarrow 2x^k(x^{k+1} - 0) = (x^k - 0)^2$
 $x^{k+1} - 0 = \frac{1}{2}(x^k - 0)$ for $x^k \neq x^* = 0$
 or $(x_{k+1} - x^*) = \frac{1}{2}(x_k - x^*)$

$\frac{df}{dx}(x^k)(x^{k+1} - x^*) = \frac{d^2f}{dx^2}(\tilde{x})(x^k - x^*)^2$

Note: $\left(\frac{df}{dx}\right)^{-1}$ not bounded away from zero

convergence is linear
 ❖ Graph – convergence of examples 1 & 2



❖ Theorem

Suppose $\left[\frac{df}{dx}(x) \right]^{-1} \frac{d^2f}{dx^2}(x) \leq L$ for all x
 if $L|x_0 - x^*| \leq \gamma < 1$ then x_k converges to x^*

Proof: $|x_1 - x^*| \leq L|x_0 - x^*||x_0 - x^*|$
 $\Rightarrow |x_1 - x^*| \leq \gamma|x_0 - x^*|$
 $\Rightarrow |x_2 - x^*| \leq L\gamma|x_0 - x^*||x_1 - x^*|$
 or $|x_2 - x^*| \leq \gamma^2|x_0 - x^*| \leq \gamma^3|x_0 - x^*|$
 $\Rightarrow |x_3 - x^*| \leq \gamma^4|x_2 - x^*| \leq \gamma^7|x_0 - x^*|$

$$|x^{k+1} - x^*| = \left[\frac{\partial f}{\partial x}(x^k) \right]^{-1} \frac{\partial^2 f}{\partial x^2}(x^k) |x^k - x^*|^2$$

$$\leq L|x^k - x^*|^2$$

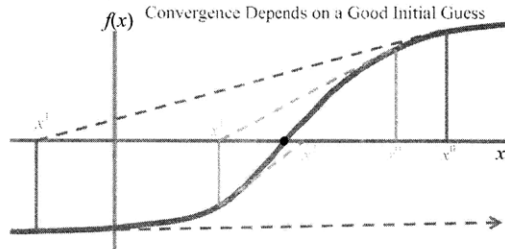
$k = 0: |x^1 - x^*| \leq L|x^0 - x^*||x^0 - x^*|$
 $\leq \gamma|x^0 - x^*|$
 $k = 1: |x^2 - x^*| \leq L|x^1 - x^*||x^1 - x^*|$
 $\leq L\gamma|x^0 - x^*||\gamma|x^0 - x^*||$
 $\leq \gamma L|x^0 - x^*||\gamma|x^0 - x^*||$
 $\leq \gamma^3|x^0 - x^*|$
 \vdots
 $k: |x^{k+1} - x^*| \leq \gamma^{2(k+1)-1}|x^0 - x^*|$

$$\gamma = \left[\frac{\partial f}{\partial x}(x) \right]^{-1} \left[\frac{\partial^2 f}{\partial x^2}(x) \right] |x^0 - x^*| < 1$$

❖ Theorem

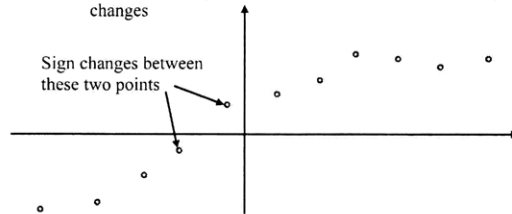
If L is bounded $\left(\frac{df}{dx} \right.$ bounded away from zero; $\left. \frac{d^2f}{dx^2} \right)$ bounded then Newton's method is guaranteed to converge given a "close enough" guess
 Always converges?

❖ 1D Pictorial Example



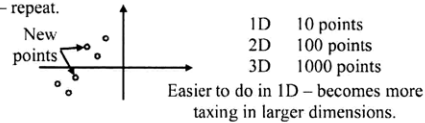
Aside on Global Convergence

What if your initial guess does not converge?
 Take a sampling of guesses and determine where the sign changes



Run Newton's Method again using one of the two points surrounding the sign change as your initial guess.

If it still does not converge? Sample points within the range of those two points and refine the area of the sign change - repeat.



Example: Heat Conducting Bar – Insulating

Pure Neumann Problem:
 Boundary conditions – heat flow only
 leads to a singular matrix
 If $h_0 = h_n$ there are multiple solutions,
 if $h_0 \neq h_n$ then there is NO solution

Example: Heat Conducting Bar – Lossy

Known Temperature

Krylov Subspace Methods converge more rapidly
 (small ratio of eigenvalues)

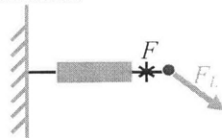
$$\mathbf{M} \begin{bmatrix} T_0 \\ \vdots \\ T_n \end{bmatrix} = \begin{bmatrix} h_0 \\ \vdots \\ h_n \end{bmatrix} \quad \begin{bmatrix} T_0 \\ \vdots \\ T_n \end{bmatrix} = \mathbf{M}^{-1} \begin{bmatrix} h_0 \\ \vdots \\ h_n \end{bmatrix}$$

$\lambda = \text{eigenvalue of } \mathbf{M} \implies \frac{1}{\lambda} \text{ is an eigenvalue of } \mathbf{M}^{-1}$

Multidimensional Newton's Method

❖ Examples

○ Strut and Joint



$$l = \sqrt{x^2 + y^2}$$

$$F = \epsilon A_s \frac{(l - l_0)}{l_0} = E(l - l_0)$$

$$f_x = \frac{x}{l} F = \frac{x}{l} E(l - l_0)$$

$$f_y = \frac{y}{l} F = \frac{y}{l} E(l - l_0)$$

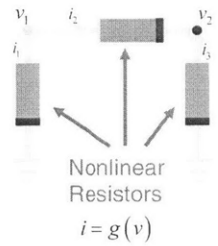
$$F(\vec{x}) = \begin{matrix} f_x + F_{L_x} = 0 \\ f_y + F_{L_y} = 0 \end{matrix}$$

OR

$$\frac{x}{l} E(l_0 - l) + F_{L_x} = 0$$

$$\frac{y}{l} E(l_0 - l) + F_{L_y} = 0$$

○ Nonlinear Resistors



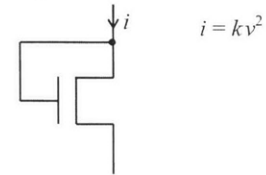
Nodal Analysis

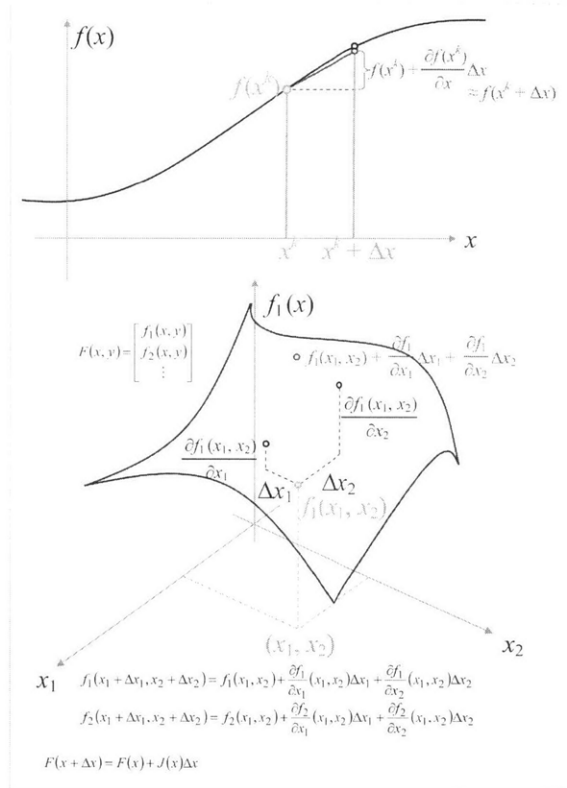
At Node 1: $i_1 + i_2 = 0$
 $\Rightarrow g(v_1) + g(v_1 - v_2) = 0$

At Node 2: $i_3 - i_2 = 0$
 $\Rightarrow g(v_2) - g(v_1 - v_2) = 0$

Two coupled nonlinear equations in two unknowns

Example. Nonlinear circuit element: MOSFET





❖ Jacobian Matrix

$$J_F(x) \Delta x \approx F(x + \Delta x) - F(x)$$

$$J_F(x) \Delta x = \begin{bmatrix} \frac{\partial F_1(x)}{\partial x_1} & \dots & \frac{\partial F_1(x)}{\partial x_N} \\ \vdots & \ddots & \vdots \\ \frac{\partial F_N(x)}{\partial x_1} & \dots & \frac{\partial F_N(x)}{\partial x_N} \end{bmatrix} \begin{bmatrix} \Delta x_1 \\ \vdots \\ \Delta x_N \end{bmatrix}$$

❖ General Setting

Problem: Find x^* such that $F(x^*) = 0$, $x^* \in \mathbb{R}^N$ and $F: \mathbb{R}^N \rightarrow \mathbb{R}^N$
 Approximate $F(x)$ with its Taylor Series about x^k :

$$\tilde{F}(x) \cong F(x^k) + J_F(x^k)(x^{k+1} - x^k)$$

Jacobian Matrix

Find the solution of the approximation: x^{k+1}

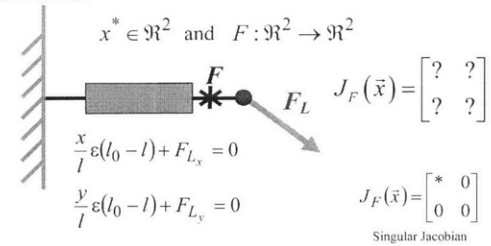
$$F(x^k) + J_F(x^k)(x^{k+1} - x^k) = 0 \Rightarrow x^{k+1} = x^k - [J_F(x^k)]^{-1} F(x^k)$$

❖ Newton Algorithm

- x^0 = initial guess, $k = 0$
- Repeat {
- Compute $F(x^k)$, $J_F(x^k)$
- Solve $J_F(x^k)(x^{k+1} - x^k) = -F(x^k)$ for x^{k+1}
- $k = k + 1$
- } Until $\|x^{k+1} - x^k\|, \|f(x^{k+1})\|$ small enough

❖ Nodal Analysis

- o Strut and Joint



$$J_F(\bar{x}) = \begin{bmatrix} * & 0 \\ 0 & \textcircled{0} \end{bmatrix}$$

$$\frac{\partial}{\partial y} \left[\frac{y}{l} \varepsilon(l_0 - l) + F_{L_y} \right]_{l_0 = \sqrt{x_0^2 + y_0^2}} \Big|_{y=0}$$

$$= \left[\frac{1}{l} \varepsilon(l_0 - l) + \frac{y}{l} \varepsilon \left(\frac{\partial \mathcal{L}}{\partial y} \right) + y \varepsilon(l_0 - l) \left\{ \frac{\partial}{\partial y} \left(\frac{1}{l} \right) \right\} \right]_{l_0 = \sqrt{x_0^2 + y_0^2}} \Big|_{y=0}$$

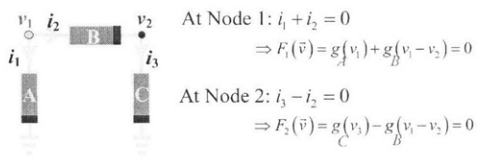
scalar functions

$$\begin{aligned} \rightarrow F_1 \begin{pmatrix} x + \Delta x \\ y + \Delta y \end{pmatrix} &= F_1 \begin{pmatrix} x \\ y \end{pmatrix} + \frac{\partial F_1}{\partial x} \Delta x + \frac{\partial F_1}{\partial y} \Delta y \\ \rightarrow F_2 \begin{pmatrix} x + \Delta x \\ y + \Delta y \end{pmatrix} &= F_2 \begin{pmatrix} x \\ y \end{pmatrix} + \frac{\partial F_2}{\partial x} \Delta x + \frac{\partial F_2}{\partial y} \Delta y \end{aligned}$$

$$\Rightarrow \bar{F}(\bar{x} + \Delta \bar{x}) \approx \bar{F}(\bar{x}) + \begin{bmatrix} \frac{\partial F_1}{\partial x} & \frac{\partial F_1}{\partial y} \\ \frac{\partial F_2}{\partial x} & \frac{\partial F_2}{\partial y} \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix}$$

o Nonlinear Resistor

$$x^* \in \mathbb{R}^2 \text{ and } F: \mathbb{R}^2 \rightarrow \mathbb{R}^2$$



$$J_F(\bar{x}) = \begin{bmatrix} ? & ? \\ ? & ? \end{bmatrix}$$

$$J_F(v_1, v_2) = \begin{bmatrix} \frac{\partial g_A}{\partial v_1} & \frac{\partial g_B}{\partial v_1} \\ -\frac{\partial g_B}{\partial v_1} & \frac{\partial g_C}{\partial v_2} - \frac{\partial g_B}{\partial v_2} \end{bmatrix}$$

❖ Computing the Jacobian and the Function

Consider the contribution of one nonlinear resistor Connected between nodes n_1 and n_2

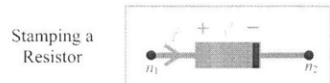


Summing currents at Node n_1 : $F_{n_1}(v) = g(v_{n_1} - v_{n_2}) + \dots$
 Summing currents at Node n_2 : $F_{n_2}(v) = -g(v_{n_1} - v_{n_2}) + \dots$

Differentiating at Node n_1 :

$$\frac{\partial F_{n_1}(v)}{\partial v_{n_1}} = \frac{\partial g(v_{n_1} - v_{n_2})}{\partial v_{n_1}} + \dots = \frac{\partial g}{\partial v}$$

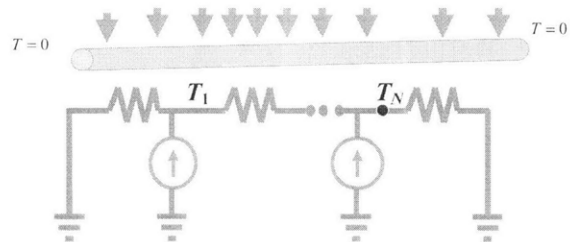
$$\frac{\partial F_{n_2}(v)}{\partial v_{n_2}} = \frac{\partial (-g(v_{n_1} - v_{n_2}))}{\partial v_{n_2}} + \dots = -\frac{\partial g}{\partial v}$$



$$\begin{matrix}
 & n_1 & & n_2 & \\
 & \vdots & & \vdots & \\
 n_1 & \dots & \frac{\partial g(v_{n_1} - v_{n_2})}{\partial v} & \dots & -\frac{\partial g(v_{n_1} - v_{n_2})}{\partial v} & \dots \\
 & \vdots & & \vdots & & \\
 n_2 & \dots & -\frac{\partial g(v_{n_1} - v_{n_2})}{\partial v} & \dots & \frac{\partial g(v_{n_1} - v_{n_2})}{\partial v} & \dots \\
 & \vdots & & \vdots & & \\
 & \underbrace{\hspace{10em}}_{J_F(v)} & & \underbrace{\hspace{10em}}_{F(v)} & &
 \end{matrix}$$

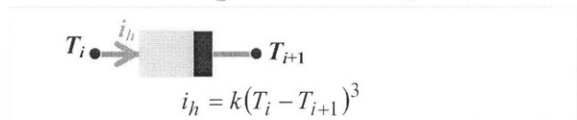
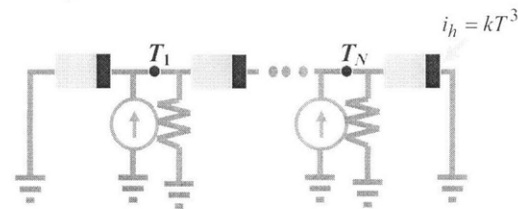
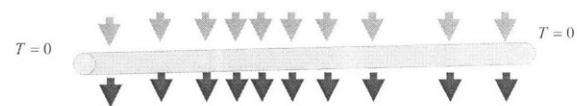
- ❖ More Complete Newton Algorithm
 - x^0 = initial guess, $k = 0$
 - Repeat {
 - Zero J_F and F
 - for each element
 - Compute element currents and derivatives
 - Sum currents to F , sum derivatives to J_F
 - Solve $J_F(x^k)(x^{k+1} - x^k) = -F(x^k)$ for x^{k+1}
 - $k = k + 1$
 - } Until $\|x^{k+1} - x^k\|, \|f(x^{k+1})\|$ small enough

❖ Example: Heat Flow in Immersed Bar



What is the Jacobian?

$$\begin{bmatrix}
 \frac{2}{R} & -\frac{1}{R} & & & \\
 -\frac{1}{R} & \frac{2}{R} & -\frac{1}{R} & & \\
 & -\frac{1}{R} & \frac{2}{R} & \ddots & \\
 & & & \ddots & \ddots \\
 & & & & & \frac{2}{R} & -\frac{1}{R} \\
 & & & & & -\frac{1}{R} & \frac{2}{R}
 \end{bmatrix}
 \begin{bmatrix}
 T_1 \\
 T_2 \\
 T_3 \\
 \vdots
 \end{bmatrix}$$



What is the Jacobian?

$$\begin{bmatrix}
 \frac{1}{R} + 3k(0 - T_1)^2 + 3k(T_1 - T_2)^2 & & -3k(T_1 - T_2)^2 & & \dots \\
 -3k(T_1 - T_2)^2 & \frac{1}{R} + 3k(T_1 - T_2)^2 + 3k(T_2 - T_3)^2 & & & \dots \\
 \vdots & & \vdots & \ddots & \ddots
 \end{bmatrix}$$

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 12.

Multi-Dimensional Newton Method

TODAY'S OUTLINE:

- ❖ Multi-Dimensional Newton Method
 - Convergence Analysis
- ❖ Damped Newton Schemes
 - Global Convergence
 - Difficulty with Singular Jacobians
- ❖ Continuation Schemes (Homotopy Methods)
 - Source/Load Stepping
 - Improving Continuation Efficiency

MULTI-DIMENSIONAL NEWTON METHOD

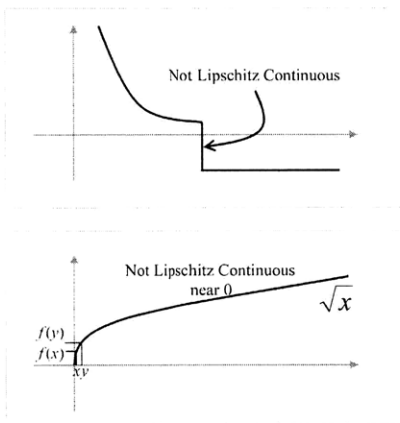
Convergence Analysis

❖ Theorem Statement

If a.) $\|J_F^{-1}(x^k)\| \leq \beta$ (Inverse is bounded)

b.) $\|J_F(x) - J_F(y)\| \leq \ell \|x - y\|$ (Derivative is Lipschitz Cont)

Then Newton's method converges given a sufficiently close initial guess.



❖ Key Lemma

If $\|J_F(x) - J_F(y)\| \leq \ell \|x - y\|$ (Derivative is Lipschitz Continuous)

Then $\|F(x) - F(y) - J_F(y)(x - y)\| \leq \frac{\ell}{2} \|x - y\|^2$

There is no multidimensional mean value theorem

❖ Theorem Proof

By definition of the Newton Iteration and the assumed bound on the inverse of the Jacobian

$$\|x^{k+1} - x^k\| = \|J_F^{-1}(x^k)F(x^k)\| \leq \beta \|F(x^k)\|$$

Again applying the Newton iteration definition

$$\|x^{k+1} - x^k\| \leq \beta \left\| \underbrace{F(x^k) - F(x^{k-1}) - J_F(x^{k-1})(x^k - x^{k-1})}_0 \right\|$$

Finally using the Lemma

$$\|x^{k+1} - x^k\| \leq \frac{\beta \ell}{2} \|x^k - x^{k-1}\|^2$$

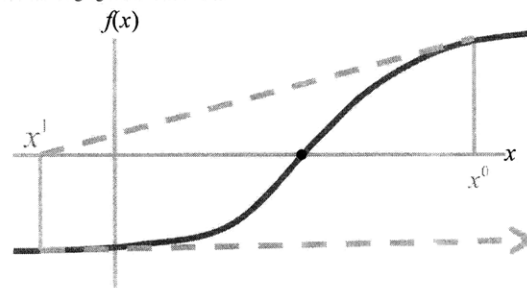
❖ Implications

If the function is not very steep in some direction, or not very smooth...

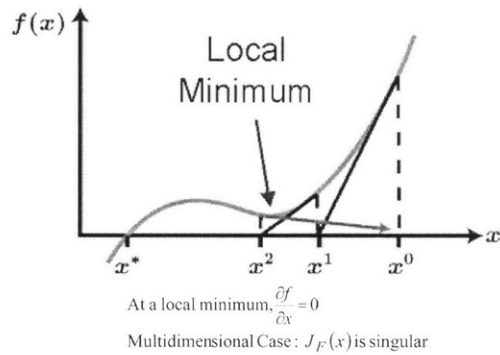
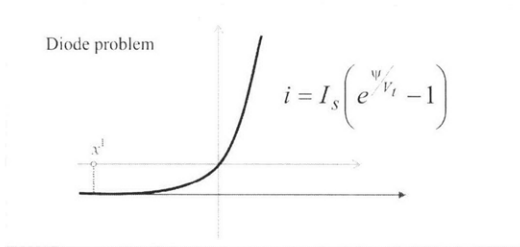
Then Newton's method can be used to find the zero of the function but only if you all ready know the answer...

Need a way to develop Newton methods which converge regardless of initial guess!

❖ Non-converging Case: 1-D Picture



Limiting the changes in X might improve convergence



DAMPED NEWTON SCHEMES

Global Convergence

❖ Newton Algorithm for Solving $F(x) = 0$

\bar{x}^0 = initial guess, $k = 0$

Repeat {

 Compute $F(\bar{x}^k), J_F(\bar{x}^k)$

 Solve $J_F(\bar{x}^k) \Delta \bar{x}^{k+1} = -F(\bar{x}^k)$ for $\Delta \bar{x}^{k+1}$

$\bar{x}^{k+1} = \bar{x}^k + \text{limited}(\Delta \bar{x}^{k+1})$

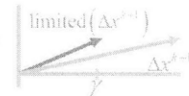
$k = k + 1$

}; Until $\|\Delta \bar{x}^{k+1}\|, \|F(\bar{x}^{k+1})\|$ small enough

❖ Limiting Methods

o Direction Corrupting

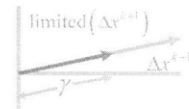
$$\text{limited}(\Delta \bar{x}^{k+1}) = \begin{cases} \Delta \bar{x}_i^{k+1} & \text{if } |\Delta \bar{x}_i^{k+1}| < \gamma \\ \gamma \text{sign}(\Delta \bar{x}_i^{k+1}) & \text{otherwise} \end{cases}$$



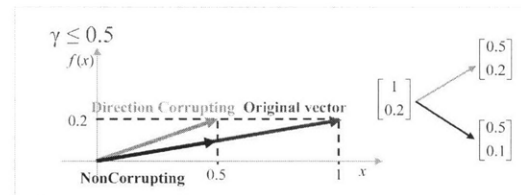
o NonCorrupting

$$\text{limited}(\Delta \bar{x}^{k+1}) = \alpha \Delta \bar{x}^{k+1}$$

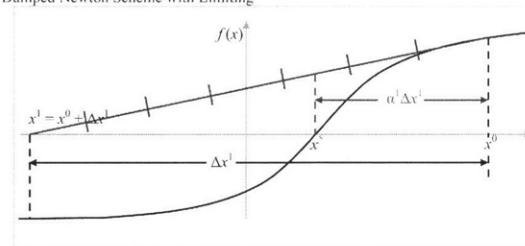
$$\alpha = \min \left\{ 1, \frac{\gamma}{\|\Delta \bar{x}^{k+1}\|} \right\}$$



Heuristics, No Guarantee of Global Convergence



❖ Damped Newton Scheme with Limiting



General Damping Scheme

$$\text{Solve } J_F(\bar{x}^k) \Delta \bar{x}^{k+1} = -F(\bar{x}^k) \text{ for } \Delta \bar{x}^{k+1}$$

$$\bar{x}^{k+1} = \bar{x}^k + \alpha^k \Delta \bar{x}^{k+1}$$

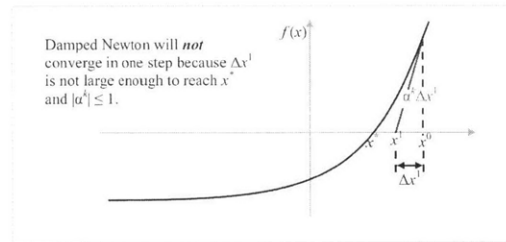
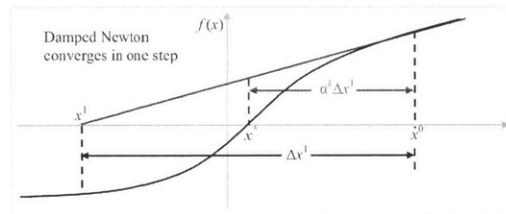
Key Idea: Line Search

Pick α^k to minimize $\|F(x^k + \alpha^k \Delta x^{k+1})\|_2 = [F(x^k + \alpha^k \Delta x^{k+1})]^T [F(x^k + \alpha^k \Delta x^{k+1})]$

Example.

$$\alpha = 1 \quad \|F(x^k + \alpha^k \Delta x^{k+1})\|_2 \gg \|F(x^k)\|_2$$

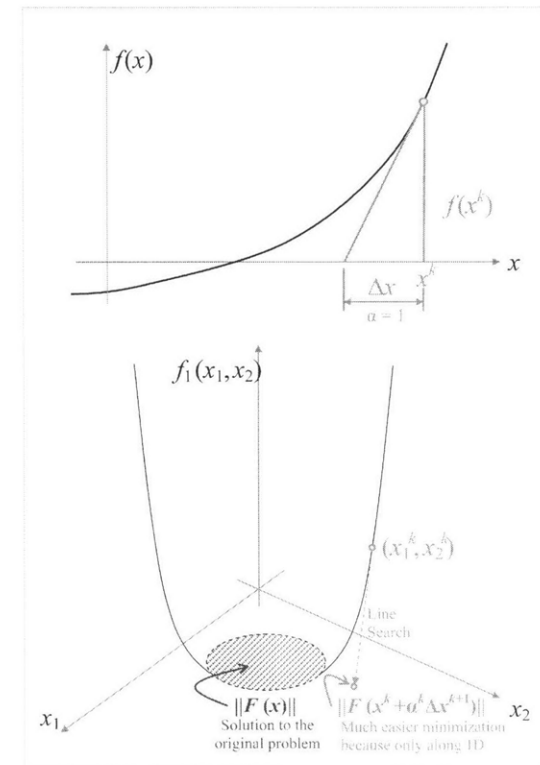
$$\alpha = \frac{1}{2} \quad \|F(x^k + \alpha^k \Delta x^{k+1})\|_2 \ll \|F(x^k)\|_2$$



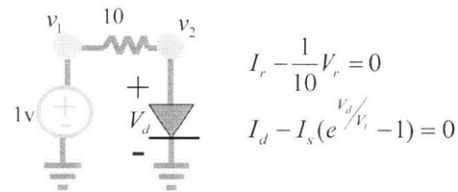
- o Nested Iteration
 - \bar{x}^0 = initial guess, $k = 0$
 - Repeat {
 - Compute $F(\bar{x}^k), J_F(\bar{x}^k)$
 - Solve $J_F(\bar{x}^k) \Delta x^{k+1} = -F(\bar{x}^k)$ for Δx^{k+1}
 - Find $\alpha^k \in (0,1]$ such that $\|F(x^k + \alpha^k \Delta x^{k+1})\|$ is minimized
 - $\bar{x}^{k+1} = \bar{x}^k + \alpha^k \Delta x^{k+1}$
 - $k = k + 1$

} Until $\|\Delta x^{k+1}\|, \|F(\bar{x}^{k+1})\|$ small enough

Method Performs a one-dimensional search in Newton Direction.

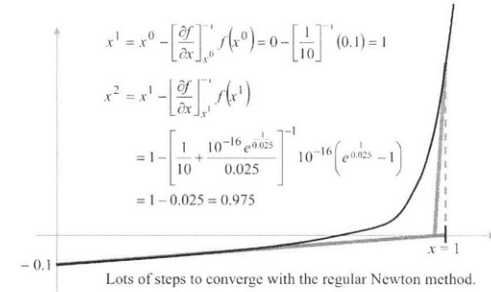
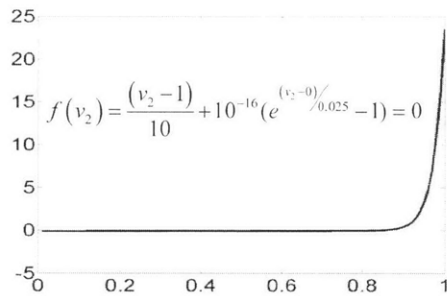


o Example



Nodal Equations with Numerical Values

$$f(v_2) = \frac{(v_2 - 1)}{10} + 10^{-16}(e^{v_2/0.025} - 1) = 0$$



o Nested Iteration

\bar{x}^0 = initial guess, $k = 0$

Repeat {

 Compute $F(\bar{x}^k)$, $J_F(\bar{x}^k)$

 Solve $J_F(\bar{x}^k) \Delta \bar{x}^{k+1} = -F(\bar{x}^k)$ for $\Delta \bar{x}^{k+1}$

 Find $\alpha^k \in (0,1]$ such that $\|F(\bar{x}^k + \alpha^k \Delta \bar{x}^{k+1})\|$ is minimized

$\bar{x}^{k+1} = \bar{x}^k + \alpha^k \Delta \bar{x}^{k+1}$

$k = k + 1$

} Until $\|\Delta \bar{x}^{k+1}\|, \|F(\bar{x}^{k+1})\|$ small enough

How can one find the damping coefficients?

o Convergence Theorem

If a.) $\|J_F^{-1}(x^k)\| \leq \beta$ (Inverse is bounded)

b.) $\|J_F(x) - J_F(y)\| \leq \ell \|x - y\|$ (Derivative is Lipschitz Cont)

Then There exists a set of α^k 's $\in (0,1]$ such that

$$\|F(x^{k+1})\| = \|F(x^k + \alpha^k \Delta x^{k+1})\| < \gamma \|F(x^k)\| \text{ with } \gamma < 1$$

Every Step reduces F—Global Convergence!

Example.

$$\begin{aligned} \left\| F(x^k + \alpha^k \Delta x^{k+1}) \right\|_2^2 &\leq 0.9 \left\| F(x^k) \right\|_2^2 \\ \left\| F(x^{k+1} + \alpha^{k+1} \Delta x^{k+2}) \right\|_2^2 &\leq 0.9 \left\| F(x^{k+1}) \right\|_2^2 \leq 0.9^2 \left\| F(x^k) \right\|_2^2 \\ \left\| F(x^{k+100} + \alpha^{k+100} \Delta x^{k+101}) \right\|_2^2 &\leq 0.9 \left\| F(x^{k+100}) \right\|_2^2 \leq \underset{\rightarrow 0}{0.9^{101}} \left\| F(x^k) \right\|_2^2 \rightarrow 0 \end{aligned}$$

• Proof.

By definition of the Newton Iteration

$$x^{k+1} = x^k - \alpha^k J_F(x^k)^{-1} F(x^k)$$

Newton Direction

Multidimensional Mean Value Lemma

$$\|F(x) - F(y) - J_F(y)(x - y)\| \leq \frac{\ell}{2} \|x - y\|^2$$

Combining

$$\left\| F(x^{k+1}) - F(x^k) - J_F(x^k) \left[\alpha^k J_F(x^k)^{-1} F(x^k) \right] \right\| \leq \frac{\ell}{2} \alpha^k \left\| J_F(x^k)^{-1} F(x^k) \right\|^2$$

Combining terms and moving scalars out of norms

$$\left\| F(x^{k+1}) - (1 - \alpha^k) F(x^k) \right\| \leq (\alpha^k)^2 \frac{\ell}{2} \left\| J_F(x^k)^{-1} F(x^k) \right\|^2$$

Using the Jacobian Bound and splitting the norm

$$\left\| F(x^{k+1}) \right\| \leq \left[(1 - \alpha^k) \left\| F(x^k) \right\| + (\alpha^k)^2 \frac{\beta^2 \ell}{2} \left\| F(x^k) \right\|^2 \right]$$

Yields a quadratic in the damping coefficient.

Simplifying quadratic

$$\left\| F(x^{k+1}) \right\| \leq \left[1 - \alpha^k + (\alpha^k)^2 \frac{\beta^2 \ell}{2} \left\| F(x^k) \right\| \right] \left\| F(x^k) \right\|$$

$$\min_{\alpha^k} \left[1 - \alpha^k + (\alpha^k)^2 C \right] \quad C = \frac{\beta^2 \ell}{2} \left\| F(x^k) \right\|$$

$$-1 + 2\alpha^k C = 0$$

$$\alpha^k = \frac{1}{2C}$$

Two Cases:

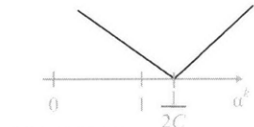
1.) $\left\| F(x^k) \right\| \leq \frac{1}{\beta^2 \ell}$ pick $\alpha^k = 1$ (Standard Newton)

$$\left\| F(x^{k+1}) \right\| \leq \left[1 - 1 + 1^2 \frac{\beta^2 \ell}{2} \left\| F(x^k) \right\| \right] \left\| F(x^k) \right\| \leq \frac{1}{2} \left\| F(x^k) \right\|$$

Case 1. $\frac{1}{2C} \geq 1$

$$\frac{1}{2 \left[\frac{\beta^2 \ell}{2} \left\| F(x^k) \right\| \right]} \geq 1$$

$$\Rightarrow \left\| F(x^k) \right\| \leq \frac{1}{\beta^2 \ell}$$



Take the full Newton step.

pick $\alpha^k = 1$ (Standard Newton)

2.) $\left\| F(x^{k+1}) \right\| > \frac{1}{\beta^2 \ell}$ pick $\alpha^k = \frac{1}{\beta^2 \ell \left\| F(x^k) \right\|}$ (Damped Newton)

$$\left\| F(x^{k+1}) \right\| \leq \left[1 - \frac{1}{\beta^2 \ell \left\| F(x^k) \right\|} \right] \left\| F(x^k) \right\| < \left\| F(x^k) \right\|$$

Case 2. $\frac{1}{2C} < 1$

$\beta^2 \ell \|F(x^k)\| > \frac{1}{2} \Rightarrow \|F(x^k)\| > \frac{1}{\beta^2 \ell}$

$\alpha^k = \frac{1}{2C} = \frac{1}{\beta^2 \ell \|F(x^k)\|}$

$\|F(x^{k+1})\| \leq \left| 1 - \frac{1}{2C} + \frac{1}{\beta^2 \ell} C \right| \|F(x^k)\| = \underbrace{\left| 1 - \frac{1}{2\beta^2 \ell \|F(x^k)\|} \right|}_{\gamma < 1} \|F(x^k)\|$

$0 < \frac{1}{2\beta^2 \ell \|F(x^k)\|} < \frac{1}{2}$

$1 > 1 - \frac{1}{2\beta^2 \ell \|F(x^k)\|} > \frac{1}{2}$

$\frac{1}{2} < 1 - \frac{1}{2\beta^2 \ell \|F(x^k)\|} < 1$

$\gamma^k = \left(1 - \frac{1}{k+1}\right) \quad \|F(x^0)\| = 1$

$\|F(x^k)\| \leq \left(1 - \frac{1}{k+1}\right) \|F(x^{k-1})\| \leq \prod_{i=1}^k \left(1 - \frac{1}{i+1}\right) \cdot 1$

$\left(1 - \frac{1}{2}\right) \left(1 - \frac{1}{3}\right) \left(1 - \frac{1}{4}\right) \left(1 - \frac{1}{5}\right) \left(1 - \frac{1}{6}\right) \dots \neq 0$

If β small ℓ small $\rightarrow \frac{1}{\beta^2 \ell}$ large Always take the full Newton

\downarrow \downarrow
 $\frac{\partial f}{\partial x}$ large $\frac{\partial^2 f}{\partial x^2}$ small

If β large ℓ large $\rightarrow \frac{1}{\beta^2 \ell}$ small

\downarrow \downarrow
 $\frac{\partial f}{\partial x}$ small $\frac{\partial^2 f}{\partial x^2}$ large

$\|F(x^{k+1})\| \leq \gamma^k \|F(x^k)\|$ not good enough, need γ independent from k

The above result does imply $\|F(x^{k+1})\| \leq \|F(x^0)\|$ not yet a convergence theorem

$\Rightarrow 1 - \frac{1}{2\beta^2 \ell \|F(x^k)\|} \leq 1 - \frac{1}{2\beta^2 \ell \|F(x^0)\|} \leq \gamma^0$

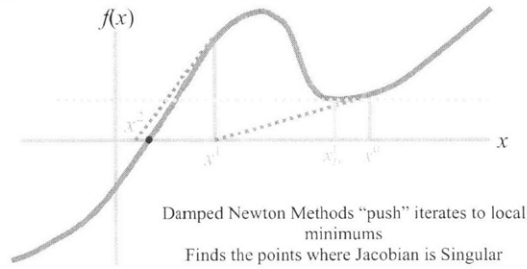
$\|F(x^{k+1})\| \leq 0.9 \|F(x^k)\|$

$\|F(x^{k+2})\| \leq 0.99 \|F(x^{k+1})\|$

$\|F(x^{k+3})\| \leq 0.999 \|F(x^{k+2})\|$

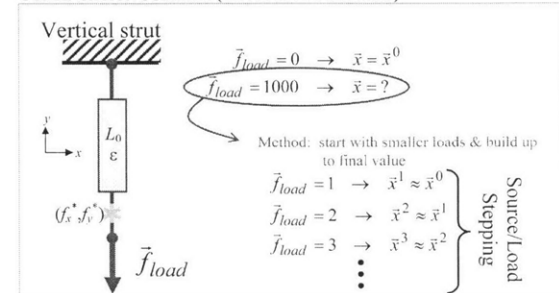
$\lim_{k \rightarrow \infty} \|F(x^k)\| \rightarrow 0?$

Difficulty with Singular Jacobians



	Regular Newton	Damped Newton	Damped Newton with half-step ($\alpha^k = 1/2$)
	X	ⓐ 1 STEP expensive more convenient multidim.	ⓐ fast
	ⓐ	ⓐ SAME # STEPS more expensive	ⓐ more steps than Newton
	ⓐ SEVERAL STEPS	ⓐ 1 STEP minimization costs	ⓐ Fast Cheap

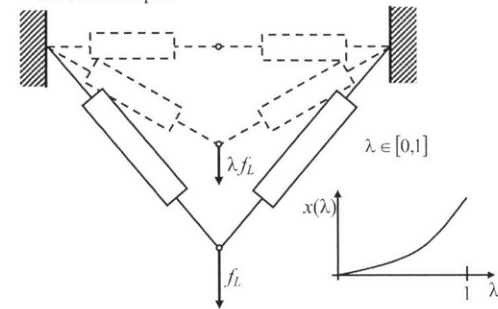
CONTINUATION SCHEMES (HOMOTOPY METHODS)



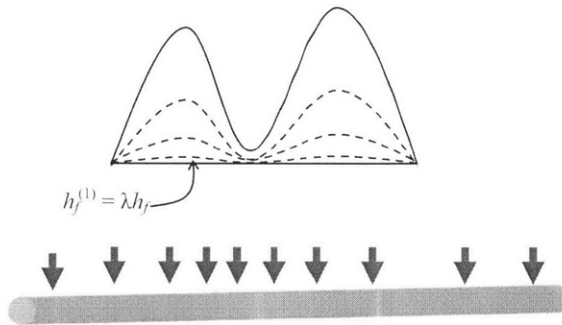
Source/Load Stepping

- ❖ Basic Concepts
 - Newton converges given a close initial guess
 - Generate a sequence of problems
 - Make sure previous problem generates guess for next problem

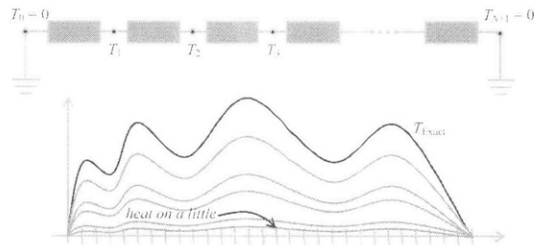
Struts Example.



o Heat-conducting bar example



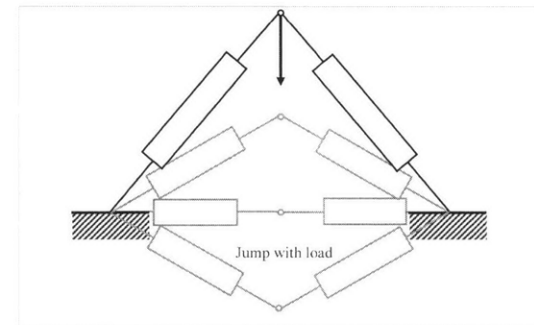
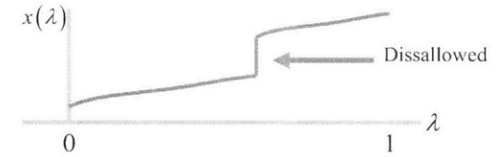
1. Start with heat off, $T=0$ is a very close initial guess
2. Increase the heat slightly, $T=0$ is a good initial guess.
3. Increase heat again...



❖ General Setting

Solve $\bar{F}(x(\lambda), \lambda) = 0$ where :

- a) $\bar{F}(x(0), 0) = 0$ is easy to solve Starts the continuation
- b) $\bar{F}(x(1), 1) = F(x)$ Ends the continuation
- c) $x(\lambda)$ is sufficiently smooth Hard to ensure!



❖ Template Algorithm

Solve $\bar{F}(x(\lambda), \lambda) = 0, x(\lambda_{prev}) = x(0)$

$\delta\lambda = 0.01, \lambda = \delta\lambda$

While $\lambda < 1$ {

$x^0(\lambda) = x(\lambda_{prev})$

Try to solve $\bar{F}(x(\lambda), \lambda) = 0$ with Newton

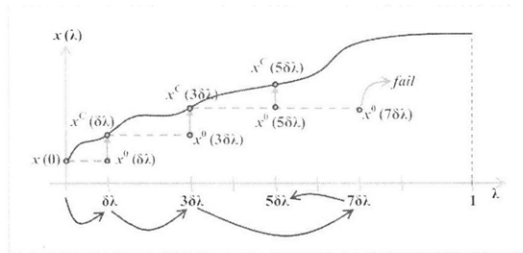
If Newton Converged

$x(\lambda_{prev}) = x(\lambda), \lambda = \lambda + \delta\lambda, \delta\lambda = 2\delta\lambda$

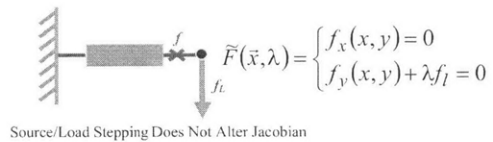
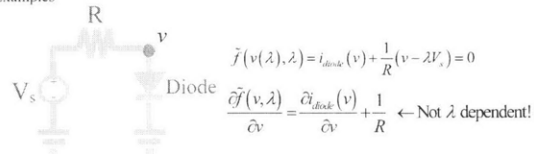
Else

$\delta\lambda = \frac{1}{2}\delta\lambda, \lambda = \lambda_{prev} + \delta\lambda$

}



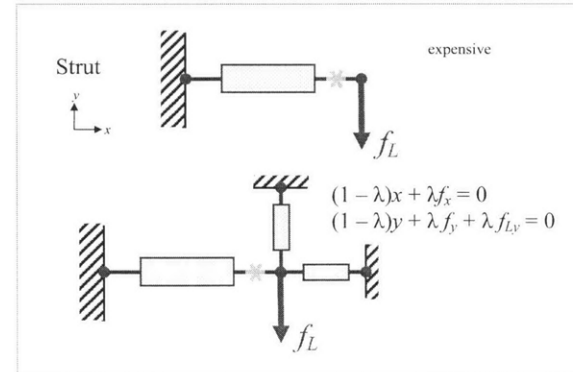
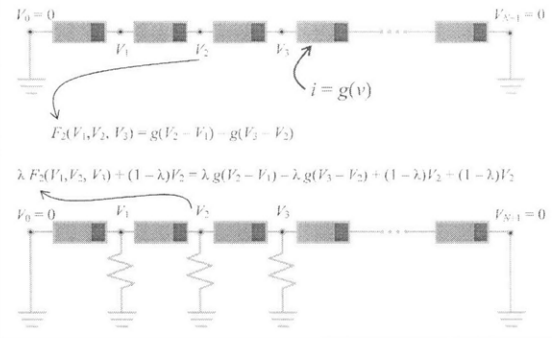
❖ Examples



Source/Load Stepping Does Not Alter Jacobian

Improving Continuation Efficiency

Nonlinear Circuit



❖ Jacobian Altering Scheme – Description

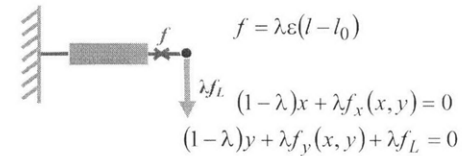
$$\tilde{F}(x(\lambda), \lambda) = \lambda F(x(\lambda)) + (1 - \lambda)x(\lambda)$$

Observations

- $\lambda = 0$ $\tilde{F}(x(0), 0) = x(0) = 0$
 $\frac{\partial \tilde{F}(x(0), 0)}{\partial x} = \mathbf{I}$ Problem is easy to solve and Jacobian is definitely nonsingular.
- $\lambda = 1$ $\tilde{F}(x(1), 1) = F(x(1))$
 $\frac{\partial \tilde{F}(x(1), 1)}{\partial x} = \frac{\partial F(x(1))}{\partial x}$ Back to the original problem and original Jacobian

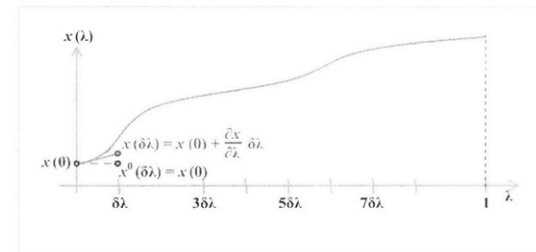
$F(v) = \frac{v - v_s}{R} + i_d(v) = 0$

$\tilde{F}(v(\lambda), \lambda) = \lambda \frac{v - v_s}{R} + \lambda i_d(v) + (1 - \lambda)v = 0$

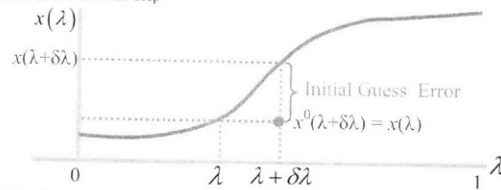
$$\underbrace{\frac{v - v_s}{R_\lambda}}_{\text{Remove Resistors then scale them slowly}} + \underbrace{\lambda i_d(v)}_{\text{Remove non-linear components then bring them back slowly}} + \underbrace{\frac{v}{1 - \lambda}}_{\text{Insert resistors to ground in each mode } R_\lambda = \frac{1}{1 - \lambda} \text{ then remove them slowly}} = 0$$


❖ Basic Algorithm

- Solve $\tilde{F}(x(\lambda), \lambda) = 0$, $x(\lambda_{prev}) = x(0)$
- $\delta\lambda = 0.01$, $\lambda = \delta\lambda$
- While $\lambda < 1$ {
- $x^0(\lambda) = x(\lambda_{prev}) + ?$
- Try to solve $\tilde{F}(x(\lambda), \lambda) = 0$ with Newton
- If Newton Converged
- $x(\lambda_{prev}) = x(\lambda)$, $\lambda = \lambda + \delta\lambda$, $\delta\lambda = 2\delta\lambda$
- Else
- $\delta\lambda = \frac{1}{2}\delta\lambda$, $\lambda = \lambda_{prev} + \delta\lambda$
- }

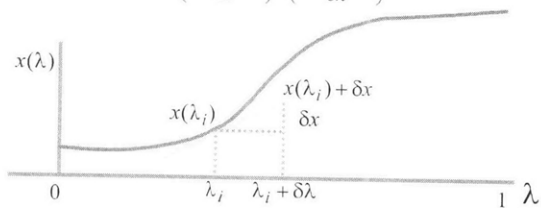


- ❖ Initial Guess for each Step



- ❖ Update Improvement

$$x^0(\lambda + \delta\lambda) = x(\lambda) - \left(\frac{\partial \tilde{F}(x(\lambda), \lambda)}{\partial x} \right)^{-1} \left(\frac{\partial \tilde{F}(x(\lambda), \lambda)}{\partial \lambda} \right) \delta\lambda$$



$$\begin{aligned} \tilde{F}(x(\lambda + \delta\lambda), \lambda + \delta\lambda) &\approx \tilde{F}(x(\lambda), \lambda) + \frac{\partial \tilde{F}(x(\lambda), \lambda)}{\partial x} (x(\lambda + \delta\lambda) - x(\lambda)) + \frac{\partial \tilde{F}(x(\lambda), \lambda)}{\partial \lambda} \delta\lambda \\ \Rightarrow \frac{\partial \tilde{F}(x(\lambda), \lambda)}{\partial x} \left(\underbrace{x^0(\lambda + \delta\lambda) - x(\lambda)}_{\substack{\text{Better guess for} \\ \text{next step's Newton}}} \right) &= - \frac{\partial \tilde{F}(x(\lambda), \lambda)}{\partial \lambda} \delta\lambda \end{aligned}$$

have from last step's Newton

$$\begin{aligned} \tilde{F}(x, \lambda) &= \lambda F(x) + (1 - \lambda)x \\ \text{Set } \delta \tilde{F} &= 0 \end{aligned}$$

$$\delta \tilde{F} = \boxed{\frac{\partial \tilde{F}}{\partial x}} \delta x + \boxed{\frac{\partial \tilde{F}}{\partial \lambda}} \delta \lambda$$

Have From last step's Newton what is it? Have From last step's Newton

$$\begin{aligned} \text{Solve for } \delta x : \boxed{\frac{\partial \tilde{F}}{\partial x}} \delta x &= - [F(x(\lambda_{prev})) - x(\lambda_{prev})] \delta \lambda \\ \Rightarrow \text{Initial guess for next Newton : } x^0(\lambda) &= x(\lambda_{prev}) + \delta x \end{aligned}$$

3.13 NonLinear Systems - A Case Study

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 13.

Newton-Method Case Study – Simulating an Image Smother

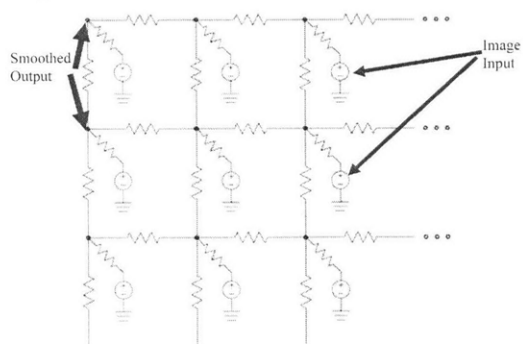
TODAY'S OUTLINE:

- ❖ Image segmentation example description
- ❖ Formulation: node-branch or nodal?
- ❖ What Solver for Newton step: sparse LU or GCR?
 - Convergence of Newton GCR
 - Matrix-free idea
- ❖ Continuation?
 - Jacobian altering scheme (for singular Jacobians)
 - Jacobian altering with update improvement
 - Arc-Length Continuation (for multiple solutions)

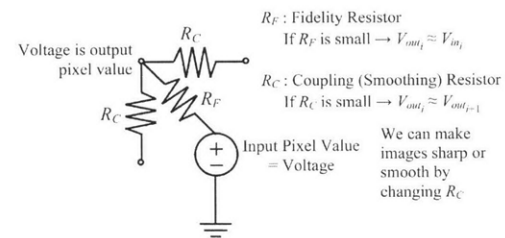
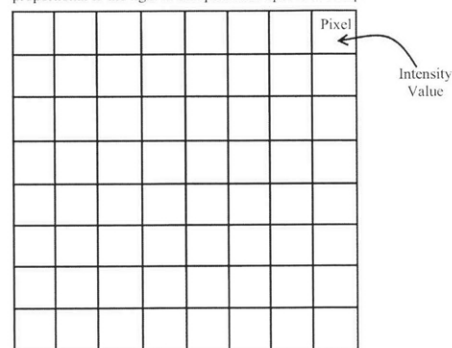
IMAGE EXAMPLE

Simple Smoother

- ❖ Circuit Diagram



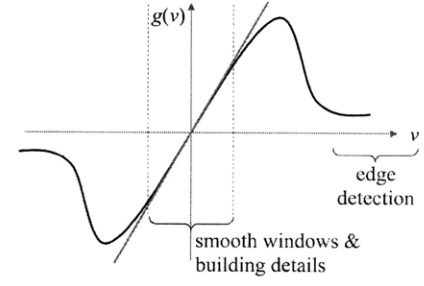
Array of photo diodes or photo transistors that provide some voltage proportional to the light in that particular spot of the chip



❖ Input Image

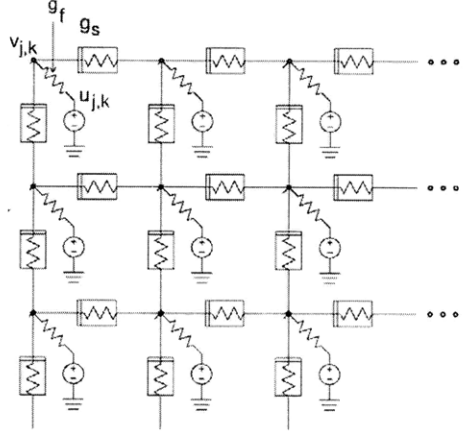


❖ Smoothed Result



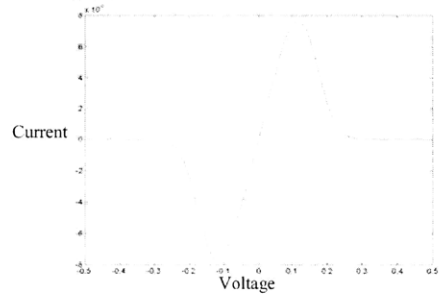
Nonlinear Smoother

❖ Circuit Diagram

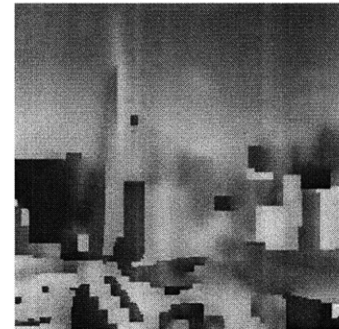
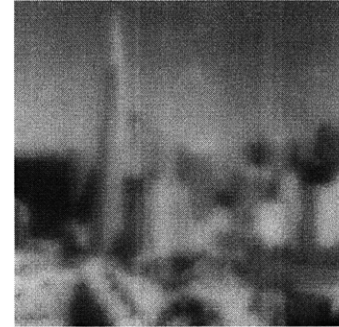
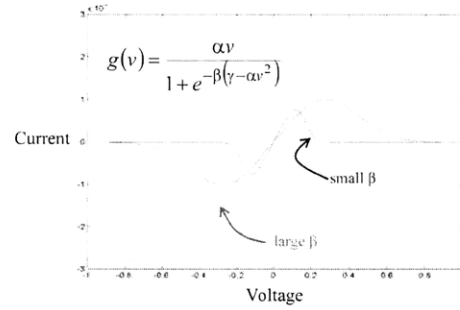


❖ Constitutive Equation

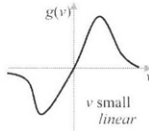
$$i(v) = \frac{\alpha v}{1 + e^{-\beta(\gamma - \alpha v^2)}}$$



○ Varying Beta



$$g(v) = \frac{\alpha v}{1 + e^{-\beta(\gamma - \alpha v^2)}}$$

$$\frac{\partial g}{\partial v} = \frac{\alpha}{1 + e^{-\beta(\gamma - \alpha v^2)}} - \frac{2\alpha^2 \beta v^2 e^{-\beta(\gamma - \alpha v^2)}}{[1 + e^{-\beta(\gamma - \alpha v^2)}]^2}$$


\mathbf{J}_F

$\begin{bmatrix} v g_1 \\ \vdots \\ v g_{m2} \end{bmatrix}$

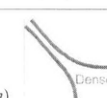
Linear Case

- sparse
- strictly diagonally dominant

Nonlinear Case

- sparsity same as linear
- strict diag. dom only if Δv 's are small

NEWTON SOLVER – SPARSE LU OR GCR?

	LU	GCR
Computation Time	$O(n^{1.5})$	$O(nq)$ $q = \#$ GCR steps
Memory	$O(n)$ 	$O(n)$ Just store vectors Mp
How accurately should we solve the system?	16 digits – no choice	Can stop after 1 digit!
Jacobian needed?	YES	NO

Basic Algorithm

❖ Nested Iteration

$x^0 =$ Initial Guess, $k = 0$
 Repeat {
 Compute $F(x^k), J_F(x^k)$
 Solve (Using GCR)
 $J_F(x^k) \Delta x^{k+1} = -F(x^k)$ for Δx^{k+1}
 $x^{k+1} = x^k + \Delta x^{k+1}$
 $k = k + 1$
 } Until $\|\Delta x^{k+1}\|, \|F(x^{k+1})\|$ small enough

– How Accurately Should We Solve with GCR?

– And do we REALLY need to assemble the Jacobian?

Matrix-free Idea

Consider Applying GCR to the Newton Iterate Equation

$$J_F(x^k) \Delta x^{k+1} = -F(x^k)$$

At each iteration GCR forms a matrix-vector product

$$J_F(x^k) v^j \approx \frac{1}{\epsilon} [F(x^k + \epsilon v^j) - F(x^k)]$$

It is possible to use Newton-GCR without Jacobians!!

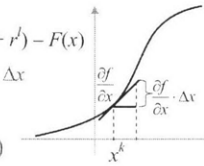
Need to select a good ϵ

l^{th} step of GCR

$$r^0 \rightarrow [J_F] v^l \neq F(x + v^l) - F(x)$$

even if $v^l \leftrightarrow \Delta x$

$$\frac{\partial f}{\partial x} \cdot \Delta x \approx f(x + \Delta x) - f(x)$$



How to choose ϵ ?

$$J_F(x^k) \cdot \Delta x \approx F(x + \Delta x) - F(x)$$

Δx small

choose $\Delta x = \epsilon v^l$ ϵ small

$$J_F(x^k) \epsilon v^l \approx F(x^k + \epsilon v^l) - F(x^k)$$

If too large, not a good approximation
 If too small,
 $F(x^k + \epsilon v^l) \approx F(x^k)$
 numerical problems/
 machine precision

❖ Basic Algorithm – Nested Iteration

$x^0 =$ Initial Guess, $k = 0$

Repeat {

 Compute $F(x^k), J_F(x^k)$

$$r^0 = -F(x^k)$$

 Repeat {

$$\text{Calculate } J_F(x^k) v^j \approx \frac{1}{\epsilon} [F(x^k + \epsilon v^j) - F(x^k)]$$

 M - Orthonormalize that vector

 Update solution and residual r^l

 } Until $\|r^l\|$ small

$$x^{k+1} = x^k + \Delta x^{k+1}$$

$k = k + 1$

} Until $\|\Delta x^{k+1}\|, \|F(x^{k+1})\|$ small enough

Convergence of Newton-GCR

❖ Basic Algorithm – Solve Accuracy Required
After l steps of GCR

$$J_F(x^k) \underbrace{\Delta x^{k+1,l}}_{\substack{\text{Newton} \\ \text{delta from} \\ l \text{ GCR steps}}} = -F(x^k) + \underbrace{r^{k,l}}_{\substack{\text{GCR} \\ \text{Residual}}}$$

- If
- a.) $|J_F^{-1}(x^k)| \leq \beta$ (Inverse is bounded)
 - b.) $|J_F(x) - J_F(y)| \leq L|x - y|$ (Derivative is Lipschitz Cont.)
 - c.) $|r^{k,l}| \leq C|F(x^k)|^2$ (More accurate near convergence)

Then
The Newton-Iterative Method Converges Quadratically

Count the number of iterations
 $\Rightarrow F(x^k) \quad F(x^k + \epsilon, p^l)$

Number of digits – doubles with every Newton step (quadratic convergence)

1 → 2 → 4 → 8 → 16

Number of digits – linear convergence

1 → 2 → 3 → 4 → 5

3 steps GCR – 1 digit $\begin{bmatrix} p^3 \\ p^0 \end{bmatrix} = 0.1$

Accuracy: 8 digits			
Quadratic	4	\times	$(3 \times 8) \times 2 = 192$
	Newton	GCR	eval
Linear	8	\times	$3 \times 2 = 48$
	Newton	GCR	eval

Just make sure $\|r^{k,l}\| \leq C \|F(x^k)\|^2$

1. Need more accuracy near solution
2. But can use as initial guess of GCR the previous Newton iterate which is very close to the solution!

How accurately should we solve GCR?
First idea: so accurate that it does not interfere with the quadratic convergence of the Newton.
For sure we don't want to solve it more accurately than $\|r^{k,l}\|$ because this is an approximation, but maybe we can do even less...

THEOREM	PRACTICALLY
Use as many GCR as needed for quadratic convergence i.e. $\ r^i\ < C \ F(x^k)\ ^2$ so that only need a few Newton steps	Stop GCR when residual is down by one order of magnitude (might need more Newton steps) What we really want to count is not the number of Newton steps but rather the total number of system evaluations (2 per each GCR)
Newton Quadratic 1 → 2 → 4 → 8 → 16 → 32 digits GCR 3 steps per digit 2 evaluations per GCR step e.g. want 8 digits: $4_{\text{Newton steps}} \times (3 \times 8)_{\text{GCR steps}} \times 2_{\text{eval}} = 192$ evaluations e.g. want 16 digits: $5_{\text{Newton steps}} \times (3 \times 16)_{\text{GCR steps}} \times 2_{\text{eval}} = 480$ evaluations	Newton Linear 1 → 2 → 3 → 4 → 5 → 6 digits GCR 3 steps per digit 2 evaluations per GCR step e.g. want 8 digits: $8_{\text{Newton steps}} \times 3_{\text{GCR steps}} \times 2_{\text{eval}} = 48$ evaluations e.g. want 16 digits: $16_{\text{Newton steps}} \times 3_{\text{GCR steps}} \times 2_{\text{eval}} = 96$ evaluations

Proof:

By definition of the Newton-Iterative Method

$$x^{k+1} = x^k - J_F(x^k)^{-1} (F(x^k) + r^{k,l})$$

Approximate Newton Direction

Multidimensional Mean Value Lemma

$$\|F(x) - F(y) - J_F(y)(x - y)\| \leq \frac{\ell}{2} \|x - y\|^2$$

Combining

$$\|F(x^{k+1}) - F(x^k) - J_F(x^k)[F(x^k) + r^{k,l}]\| \leq \frac{\ell}{2} \|J_F(x^k)^{-1} (F(x^k) + r^{k,l})\|^2$$

$$\|F(x^{k+1}) - F(x^k) - F(x^k) + r^{k,l}\| \leq \frac{\ell}{2} \|J_F(x^k)^{-1} (F(x^k) + r^{k,l})\|^2$$

Using the triangle inequality

$$\|F(x^{k+1})\| \leq \frac{\ell}{2} \|J_F(x^k)^{-1} (F(x^k) + r^{k,l})\|^2 + \|r^{k,l}\|$$

Using the Jacobian Bound

$$\|F(x^{k+1})\| \leq \frac{\beta^2 \ell}{2} \|F(x^k)\|^2 + \left(1 + \frac{\beta^2 \ell \|r^{k,l}\|}{2}\right) \|r^{k,l}\|$$

Using the bound on the iterative solver error

$$\|F(x^{k+1})\| \leq \frac{\beta^2 \ell}{2} \|F(x^k)\|^2 + \left(1 + \frac{\beta^2 \ell \|F(x^k)\|^2}{2}\right) C \|F(x^k)\|^2$$

Combining Terms

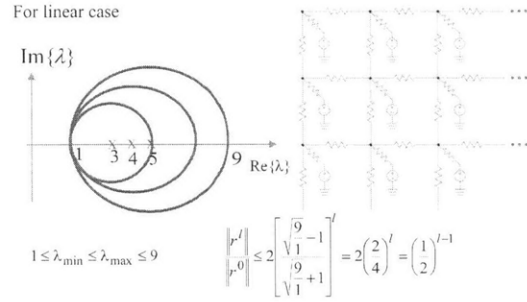
$$\|F(x^{k+1})\| \leq \left[\frac{\beta^2 \ell}{2} + \left(1 + \frac{\beta^2 \ell \|F(x^k)\|^2}{2}\right) C \right] \|F(x^k)\|^2$$

Easily Bounded

$$\Gamma^{k,l} = \underbrace{-F(x^k)}_b - \underbrace{J_F(x^k)}_M \frac{\Delta x^{k,l}}{x^k}$$

❖ Gershgorin & Chebyshev



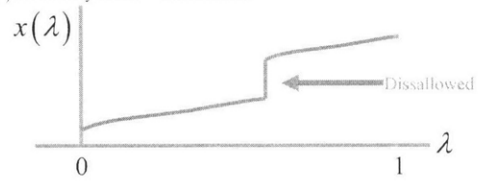


For nonlinear case:
 Could be bad.
 Could be singular.
 Need continuation scheme that fixes it.

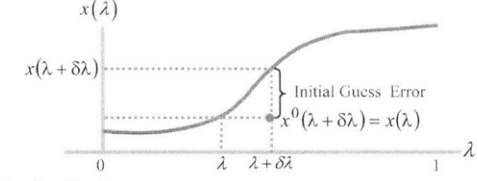
about one digit every 3 iterations:
 about 0.1% error after 3 iterations

DAMPING?
CONTINUATION?

- Basic Concepts**
 Solve $\bar{F}(x(\lambda), \lambda) = 0$ where:
- a.) $\bar{F}(x(0), 0) = 0$ is easy to solve Starts the continuation
 - b.) $\bar{F}(x(1), 1) = F(x)$ Ends the continuation
 - c.) $x(\lambda)$ is sufficiently smooth Hard to ensure!



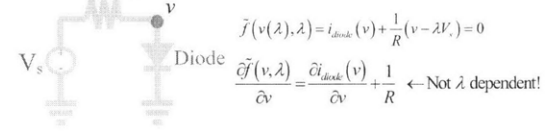
- ❖ Initial Guess for each Step
 Convergence of Newton-GCR
 Basic Algorithm – Nested Iteration

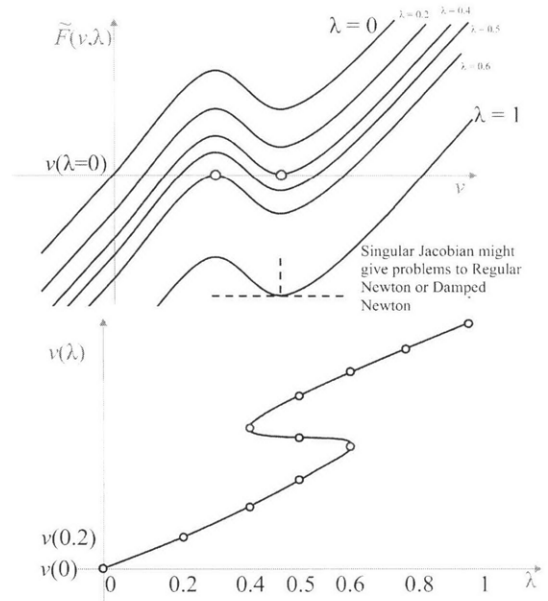
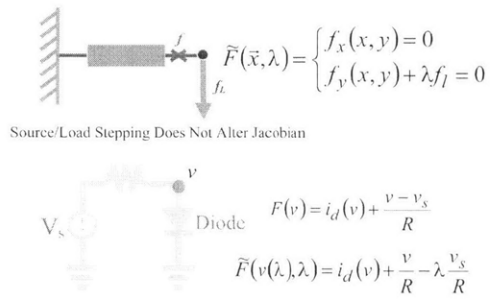


- ❖ Template Algorithm
 Solve $\bar{F}(x(0), 0)$, $x(i_{prev}) = x(0)$
 $\delta\lambda = 0.01$, $\lambda = \delta\lambda$
 While $\lambda < 1$ {
 $x^0(\lambda) = x(\lambda_{prev})$
 Try to Solve $\bar{F}(x(\lambda), \lambda) = 0$ with Newton
 If Newton Converged
 $x(\lambda_{prev}) = x(\lambda)$, $\lambda = \lambda + \delta\lambda$, $\delta\lambda = 2\delta\lambda$
 Else
 $\delta\lambda = \frac{1}{2}\delta\lambda$, $\lambda = \lambda_{prev} + \delta\lambda$
 }

- ❖ Practical Answer
 Look for another continuation scheme
 $\bar{F}(x(\lambda), \lambda) = \lambda F(x(\lambda)) + (1 - \lambda)x(\lambda)$
Observations:
 $\lambda = 0$ $\bar{F}(x(0), 0) = x(0) = 0$ Problem is easy to solve and
 $\frac{\partial \bar{F}(x(0), 0)}{\partial x} = \mathbf{1}$ Jacobian definitely nonsingular
 $\lambda = 1$ $\bar{F}(x(1), 1) = F(x(1))$ Back to the original problem
 $\frac{\partial \bar{F}(x(1), 1)}{\partial x} = \frac{\partial F(x(1))}{\partial x}$ and original Jacobian

- ❖ Source/Load Stepping Examples

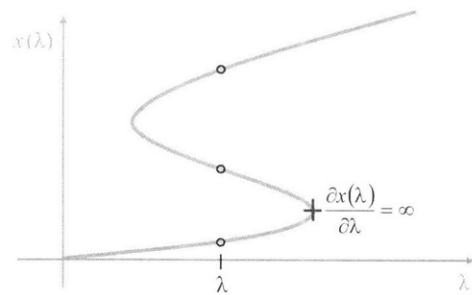
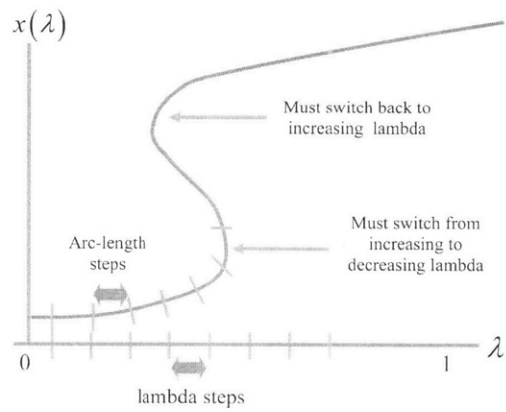




325

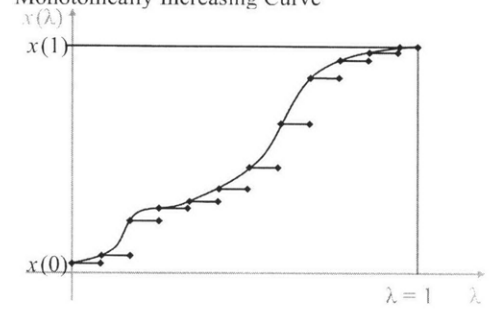
Arc-Length Continuation (for multiple solutions)

❖ Still can have problems

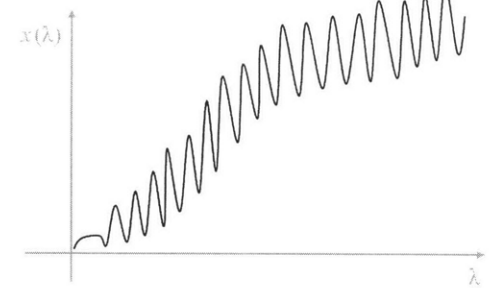


❖ Arc-length steps?

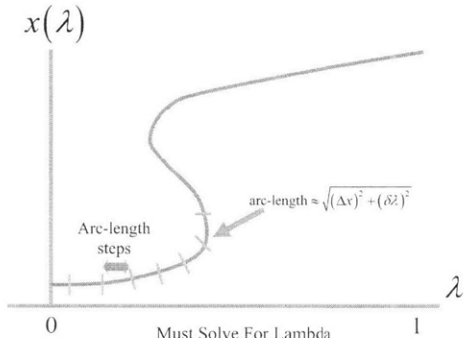
Monotonically Increasing Curve



Avoid snaking curves

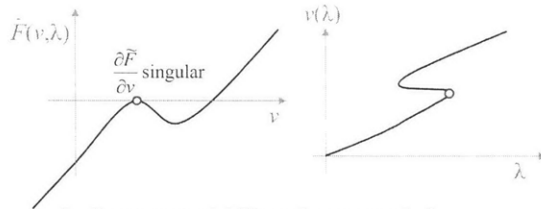


o Pictorially



Must Solve For Lambda
 $\tilde{F}(x, \lambda) = 0$

$$(\lambda - \lambda_{prev})^2 + |x - x(\lambda_{prev})|_2^2 - arc^2 = 0$$



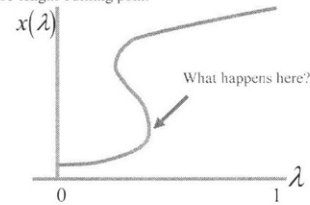
See literature on probability one homotopy methods

- (Detect if close to singular)
- in this case change arc length and hope to jump across the singularity
- (or start decreasing λ but keep moving Δx in the same direction as previous step) etc...
- Better idea: Look for another continuation method.

o By Newton

$$\begin{bmatrix} \frac{\partial \tilde{F}(x^k, \lambda^k)}{\partial x} & \frac{\partial \tilde{F}(x^k, \lambda^k)}{\partial \lambda} \\ 2(x^k - x(\lambda_{prev}))^T & 2(\lambda^k - \lambda_{prev}) \end{bmatrix} \begin{bmatrix} x^{k+1} - x^k \\ \lambda^{k+1} - \lambda^k \end{bmatrix} = - \begin{bmatrix} \tilde{F}(x^k, \lambda^k) \\ (\lambda^k - \lambda_{prev})^2 + |x^k - x(\lambda_{prev})|_2^2 - arc^2 \end{bmatrix}$$

❖ Arc-length Turning point



Upper left-hand Block is singular

$$\begin{bmatrix} \frac{\partial \tilde{F}(x^k, \lambda^k)}{\partial x} & \frac{\partial \tilde{F}(x^k, \lambda^k)}{\partial \lambda} \\ 2(x^k - x(\lambda_{prev}))^T & 2(\lambda^k - \lambda_{prev}) \end{bmatrix}$$

Jacobian Altering Scheme (for singular Jacobians)

❖ Practical Answer

Look for another continuation scheme

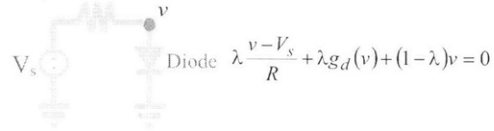
$$\tilde{F}(x(\lambda), \lambda) = \lambda F(x(\lambda)) + (1 - \lambda)x(\lambda)$$

Observations:

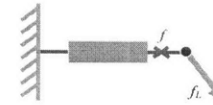
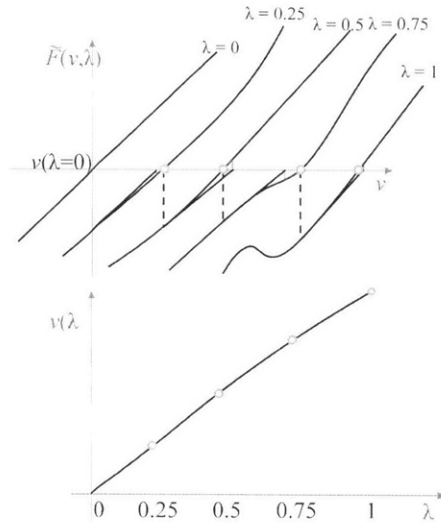
$\lambda = 0$ $\tilde{F}(x(0), 0) = x(0) = 0$
 $\frac{\partial \tilde{F}(x(0), 0)}{\partial x} = \mathbf{I}$ Problem is easy to solve and Jacobian definitely nonsingular

$\lambda = 1$ $\tilde{F}(x(1), 1) = F(x(1))$
 $\frac{\partial \tilde{F}(x(1), 1)}{\partial x} = \frac{\partial F(x(1))}{\partial x}$ Back to the original problem and original Jacobian

❖ Examples

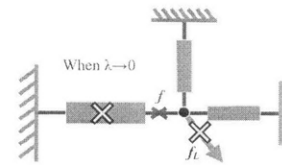


Try $\tilde{F}(v, \lambda) = \lambda F(v, \lambda) + (1 - \lambda)v = \lambda i(v) + \lambda \frac{v}{R} - \lambda \frac{V_s}{R} + (1 - \lambda)v$

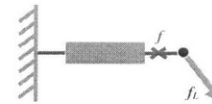


$$\tilde{F}(x(\lambda), \lambda) = \begin{cases} \lambda f_x^*(x, y) + \lambda f_{x,L} + (1 - \lambda)x = 0 \\ \lambda f_y^*(x, y) + \lambda f_{y,L} + (1 - \lambda)y = 0 \end{cases}$$

$\lambda F(x)$ $(1 - \lambda)x(\lambda)$



$$\tilde{F}(x(0), 0) = \begin{cases} x = 0 \\ y = 0 \end{cases}$$



$$\tilde{F}(x(1), 1) = \begin{cases} f_x^*(x, y) + f_{x,L} = 0 \\ f_y^*(x, y) + f_{y,L} = 0 \end{cases}$$

Jacobian Altering with Update Improvement

❖ Basic Algorithm

Solve $\tilde{F}(x(0), 0)$, $x(\lambda_{prev}) = x(0)$

$\delta\lambda = 0.01$, $\lambda = \delta\lambda$

While $\lambda < 1$ {

$x^0(\lambda) = x(\lambda_{prev}) + ?$

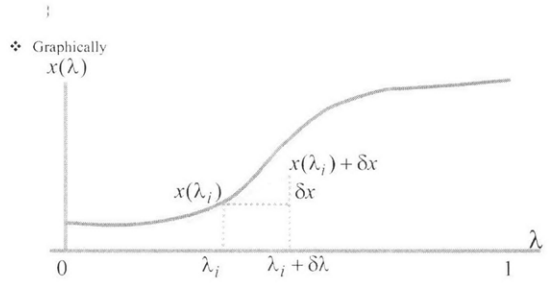
Try to Solve $\tilde{F}(x(\lambda), \lambda) = 0$ with Newton

If Newton Converged

$x(\lambda_{prev}) = x(\lambda)$, $\lambda = \lambda + \delta\lambda$, $\delta\lambda = 2\delta\lambda$

Else

$\delta\lambda = \frac{1}{2}\delta\lambda$, $\lambda = \lambda_{prev} + \delta\lambda$



❖ Update Improvement

$$\tilde{F}(x, \lambda) = \lambda F(x) + (1 - \lambda)x$$

$$\delta \tilde{F} = \left[\frac{\delta \tilde{F}}{\delta x} \right] \delta x + \left[\frac{\delta \tilde{F}}{\delta \lambda} \right] \delta \lambda$$

Have From last step's Newton what is it? Have From last step's Newton

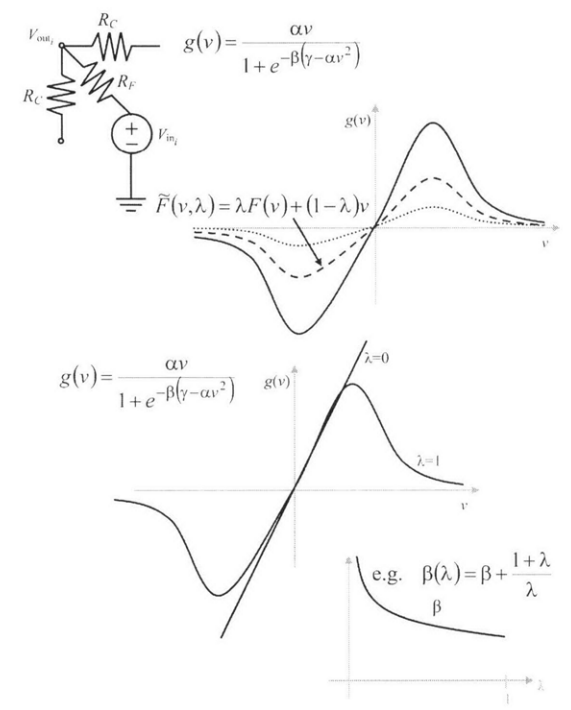
The Jacobian I used to solve the last Newton problem.

Set $\delta \tilde{F} = 0$

Solve for δx : $\left[\frac{\delta \tilde{F}}{\delta x} \right] \delta x = -[F(x(\lambda_{prev})) - x(\lambda_{prev})] \delta \lambda$

⇒ Initial guess for next Newton:

$$x^0(\lambda) = x(\lambda_{prev}) + \delta x$$



INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 14.

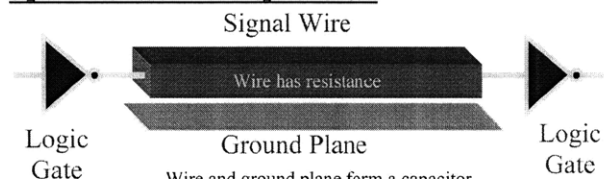
Methods for Ordinary Differential Equations

TODAY'S OUTLINE:

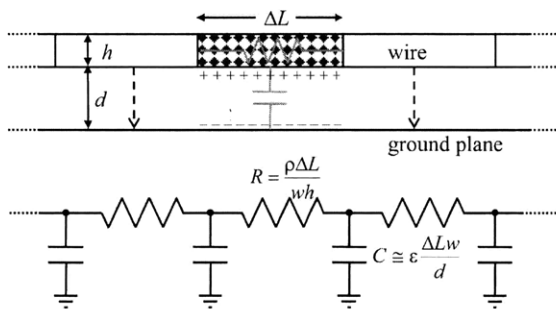
- ❖ Initial Value Problem Examples
 - Signal Propagation (circuits with capacitors)
 - Space Frame Dynamics (struts and masses)
 - Chemical Reaction Dynamics
- ❖ Eigenvalue Analysis for Dynamical Systems
- ❖ Investigate the Simple Finite-Difference Methods
 - Forward-Euler, Backward-Euler, Trapezoidal Rule
 - Examine Properties Experimentally

APPLICATION PROBLEMS

Signal Transmission in an Integrated Circuit



- Wire and ground plane form a capacitor
- Metal Wires carry signals from gate to gate.
 - How long is the signal delayed?

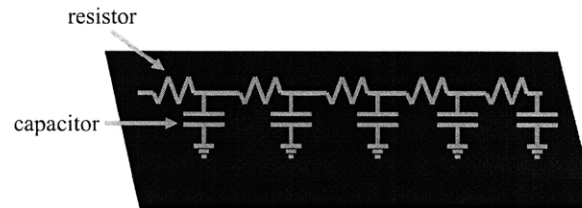


329

Capacitor
 $q \propto v$
 Charge is proportional to voltage

$$q = Cv \xrightarrow{\text{differentiate}} i_c = \frac{dq}{dt} = C \frac{dv_c}{dt}$$

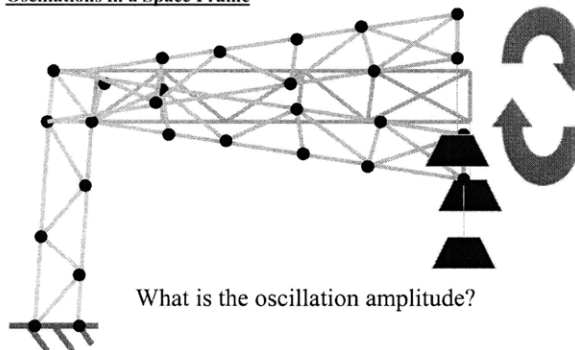
❖ Circuit Model

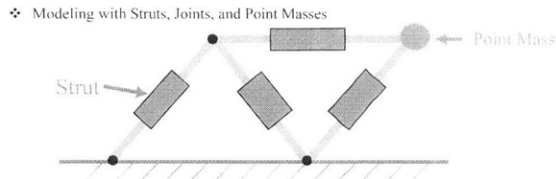
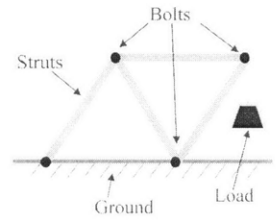
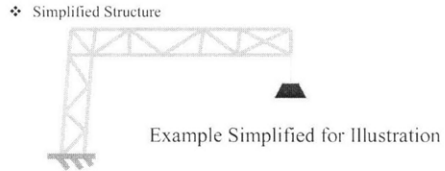


Constructing the Model

- Cut the wire into sections.
- Model wire resistance with resistors.
- Model wire-plane capacitance with capacitors.

Oscillations in a Space Frame

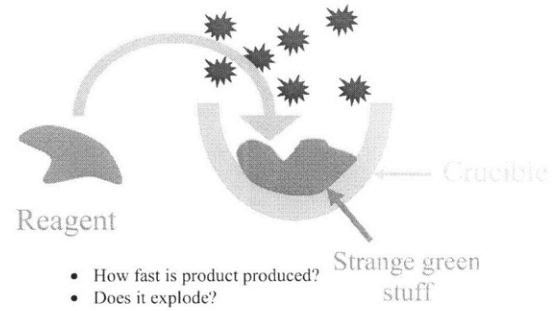




- Constructing the Model
- Replace Metal Beams with Struts.
 - Replace cargo with point mass.

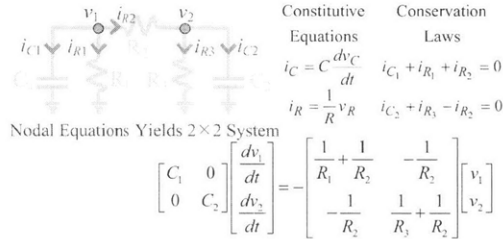
Constitutive Equations	
<p>Capacitor</p> $i_c = C \frac{dv_c}{dt}$ <p>capacitance</p>	<p>Point Mass</p> $f = M \frac{d^2x}{dt^2}$ <p>mass</p> <p>Involves t</p>
<p>Resistor</p> $i_R = \frac{1}{R} v_R$ <p>Independent of t</p>	<p>Strut</p> $f = \epsilon(l_0 - l)$

Chemical Reaction Dynamics

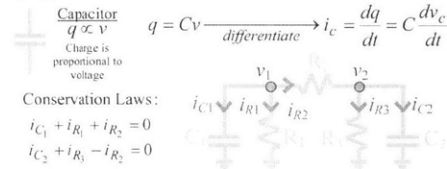


Signal Transmission in an Integrated Circuit

❖ A 2x2 Example



Stamping Procedure for Nodal Formulation in Circuits

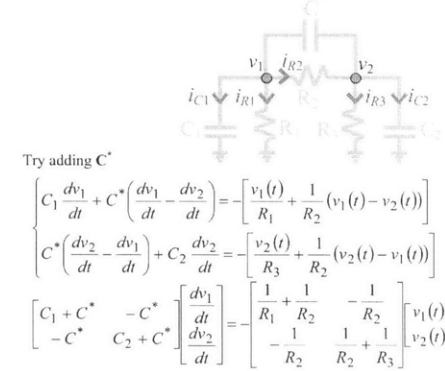


Use Nodal Analysis (substitute branch equations)

$$\begin{cases} C_1 \frac{dv_1}{dt} + \frac{v_1}{R_1} + \frac{1}{R_2} (v_1 - v_2) = 0 \\ C_2 \frac{dv_2}{dt} + \frac{v_2}{R_3} + \frac{1}{R_2} (v_2 - v_1) = 0 \end{cases}$$

$$\underbrace{\begin{bmatrix} C_1 & 0 \\ 0 & C_2 \end{bmatrix}}_{\text{Capacitance Matrix}} \begin{bmatrix} \frac{dv_1}{dt} \\ \frac{dv_2}{dt} \end{bmatrix} = - \underbrace{\begin{bmatrix} \frac{1}{R_1} + \frac{1}{R_2} & -\frac{1}{R_2} \\ -\frac{1}{R_2} & \frac{1}{R_2} + \frac{1}{R_3} \end{bmatrix}}_{\text{Conductance Matrix}} \begin{bmatrix} v_1(t) \\ v_2(t) \end{bmatrix}$$

Stamping for Capacitors is like Stamping for Resistors

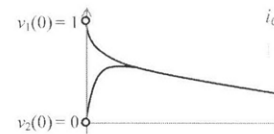


EIGENVALUE ANALYSIS FOR DYNAMICAL SYSTEMS

Signal Transmission in an Integrated Circuit

❖ A 2x2 Example

Let $C_1 = C_2 = 1, R_1 = R_3 = 10, R_2 = 1$



$$\begin{bmatrix} C_1 & 0 \\ 0 & C_2 \end{bmatrix} \begin{bmatrix} \frac{dv_1}{dt} \\ \frac{dv_2}{dt} \end{bmatrix} = - \begin{bmatrix} \frac{1}{R_1} + \frac{1}{R_2} & -\frac{1}{R_2} \\ -\frac{1}{R_2} & \frac{1}{R_2} + \frac{1}{R_3} \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \Rightarrow \frac{dx}{dt} = \underbrace{\begin{bmatrix} -1.1 & 1.0 \\ 1.0 & -1.1 \end{bmatrix}}_A x$$

Eigenvalues and Eigenvectors

$$A = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} -0.1 & 0 \\ 0 & -2.1 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}^{-1}$$

eigenvectors eigenvalues

Aside: Eigenanalysis

Change of Variables $\Leftarrow \mathbf{E}\bar{y} = \bar{x}$

$$y_1 \begin{bmatrix} \mathbf{E}_1 \\ \mathbf{E}_2 \end{bmatrix} + y_2 \begin{bmatrix} \mathbf{E}_2 \\ \mathbf{E}_2 \end{bmatrix} + \dots + y_N \begin{bmatrix} \mathbf{E}_N \\ \mathbf{E}_N \end{bmatrix} = \begin{bmatrix} \bar{x} \\ \bar{x} \end{bmatrix}$$

Columns of \mathbf{E} (eigenvectors) are new basis

e.g.

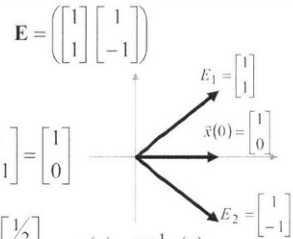
$$\bar{x}(0) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\bar{y}(0) = ?$$

$$y_1(0) \begin{bmatrix} 1 \\ 1 \end{bmatrix} + y_2(0) \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\bar{y}(0) = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix}$$

$$\bar{y}(0) = \mathbf{E}^{-1} \bar{x}(0)$$



Consider an ODE: $\frac{dx(t)}{dt} = \mathbf{A}x(t), \quad x(0) = x_0$

332

$$\text{Eigendecomposition: } \mathbf{A} = \underbrace{\begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ \mathbf{E}_1 & \mathbf{E}_2 & \dots & \mathbf{E}_n \\ \downarrow & \downarrow & \dots & \downarrow \end{bmatrix}}_{\mathbf{E}} \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \dots \\ \lambda_n \end{bmatrix} \underbrace{\begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ \mathbf{E}_1 & \mathbf{E}_2 & \dots & \mathbf{E}_n \\ \downarrow & \downarrow & \dots & \downarrow \end{bmatrix}}_{\mathbf{E}^{-1}}$$

Change of variables: $\mathbf{E}y(t) = x(t) \Leftrightarrow y(t) = \mathbf{E}^{-1}x(t)$

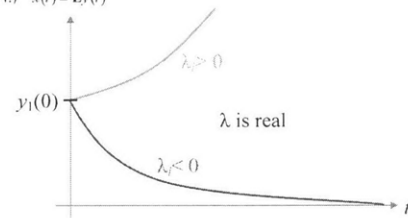
Substituting: $\frac{d\mathbf{E}y(t)}{dt} = \mathbf{A}\mathbf{E}y(t), \quad \mathbf{E}y(0) = x_0$

Multiply by \mathbf{E}^{-1} : $\frac{dy(t)}{dt} = \mathbf{E}^{-1}\mathbf{A}\mathbf{E}y(t) = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} y(t)$ Decoupled Equations!

Decoupling: $\frac{dy_i(t)}{dt} = \lambda_i y_i(t) \Rightarrow y_i(t) = e^{\lambda_i t} y_i(0)$

Steps for Solving $\frac{dx(t)}{dt} = \mathbf{A}x(t), \quad x(0) = x_0$

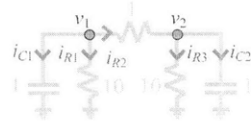
- 1.) Determine \mathbf{E}, λ
- 2.) Compute $y(0) = \mathbf{E}^{-1}x_0$
- 3.) Compute $y(t) = \begin{bmatrix} e^{\lambda_1 t} & & \\ & \ddots & \\ & & e^{\lambda_n t} \end{bmatrix} y(0)$
- 4.) $x(t) = \mathbf{E}y(t)$



$$x(t) = \underset{\text{convergence}}{\mathbf{E}} \begin{bmatrix} e^{\lambda_1 t} & & \\ & e^{\lambda_2 t} & \\ & & \ddots \\ & & & e^{\lambda_n t} \end{bmatrix} \underset{\text{convergence}}{\mathbf{E}^{-1}} x(0)$$

Signal Transmission in an Integrated Circuit

❖ A 2x2 Example



$$A = \begin{bmatrix} -1.1 & 1 \\ 1 & -1.1 \end{bmatrix} \quad E = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \quad \Lambda = \begin{bmatrix} -0.1 & 0 \\ 0 & -2.1 \end{bmatrix}$$

$$\bar{x}(0) = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \bar{y}(0) = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix}$$

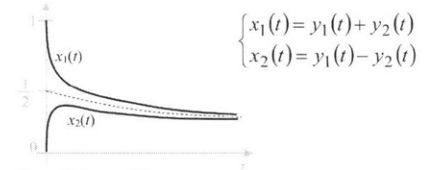
$$\frac{dy_1(t)}{dt} = -0.1y_1(t) \quad \tau = 10 \text{ms} \quad \text{slow}$$

$$y_1(t) = \frac{1}{2} e^{-0.1t} = y_1(0) e^{\lambda_1 t}$$

$$\frac{dy_2(t)}{dt} = -2.1y_2(t) \quad \tau = 0.476 \text{ms} \quad \text{fast}$$

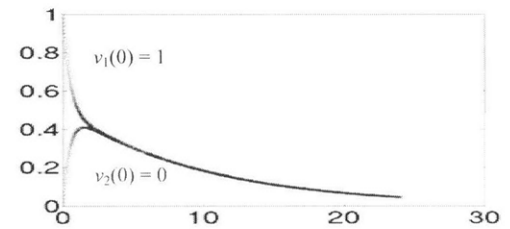
$$y_2(t) = \frac{1}{2} e^{-2.1t} = y_2(0) e^{\lambda_2 t}$$

$$\bar{x}(t) = \mathbf{E}\bar{y}(t) = y_1(t) \begin{bmatrix} 1 \\ 1 \end{bmatrix} + y_2(t) \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} e^{-0.1t} y_1(0) + \begin{bmatrix} 1 \\ -1 \end{bmatrix} e^{-2.1t} y_2(0)$$



for small times $v_2(t)$ is very active (fast mode) and decreases rapidly (the difference gets smaller)

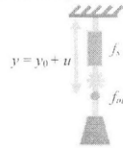
for large times the difference ≈ 0 (fast mode disappears) they both follow the slow mode $x_1(t) \approx y_1(t) \approx x_2(t)$



- Notice two time scale behavior
- v_1 and v_2 come together quickly (fast eigenmode).
 - v_1 and v_2 decay to zero slowly (slow eigenmode).

Struts, Joints and Point Mass Example

❖ A 2x2 Example



Constitutive Equations	Conservation Laws
$f_s = EA_c \frac{y - y_0}{y_0} = \frac{EA_c}{y_0} u$	$f_s + f_m = 0$
$f_m = M \frac{d^2 u}{dt^2}$	

Define v as velocity (du/dt) to yield a 2x2 System

$$\begin{cases} M \frac{d^2 u}{dt^2} + \frac{EA_c}{y_0} u(t) = 0 \\ v(t) = \frac{du}{dt} \end{cases} \rightarrow \begin{cases} M \frac{dv}{dt} = -\frac{EA_c}{y_0} u(t) \\ \frac{du}{dt} = v(t) \end{cases}$$

$$\begin{bmatrix} M & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{dv}{dt} \\ \frac{du}{dt} \end{bmatrix} = \begin{bmatrix} 0 & -\frac{EA_c}{y_0} \\ 1 & 0 \end{bmatrix} \begin{bmatrix} v \\ u \end{bmatrix}$$

Let $M = 1, \frac{EA_c}{y_0} = 1$

$$\begin{bmatrix} M & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{dv}{dt} \\ \frac{du}{dt} \end{bmatrix} = \begin{bmatrix} 0 & -\frac{EA_c}{y_0} \\ 1 & 0 \end{bmatrix} \begin{bmatrix} v \\ u \end{bmatrix} \Rightarrow \frac{dx}{dt} = \begin{bmatrix} 0 & -1.0 \\ 1.0 & 0 \end{bmatrix} x$$

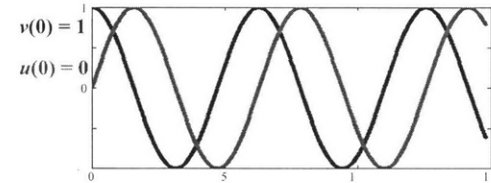
Eigenvalues and Eigenvectors

$$A = \begin{bmatrix} -1 & -1 \\ i & -i \end{bmatrix} \begin{bmatrix} i & 0 \\ 0 & -i \end{bmatrix} \begin{bmatrix} -1 & -1 \\ -1 & -1 \end{bmatrix}^{-1} \quad \begin{aligned} y_1(t) &= e^i y_1(0) = y_1(0)(\cos t + i \sin t) \\ y_2(t) &= e^{-i} y_2(0) = y_2(0)(\cos t - i \sin t) \end{aligned}$$

eigenvectors eigenvalues

$$A \vec{E} = \lambda \vec{E}$$

real complex complex



Note the system has imaginary eigenvalues

- Persistent Oscillation
- Velocity, v, peaks when displacement, u, is zero.

Chemical Reaction Example

❖ A 2x2 Example

Amount of reactant = R, the temperature = T

$$\frac{dT}{dt} = -T + R$$

More reactant causes the temperature to rise, higher temperatures increase heat dissipation causing temperature to fall.

$$\frac{dR}{dt} = -R + 4T$$

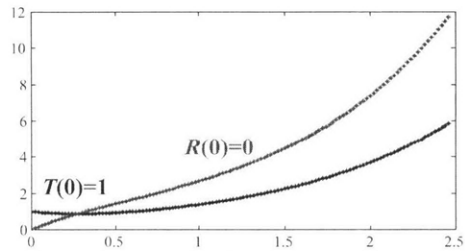
Higher temperature raises reaction rates, increased reactant interferes with reaction and slows rate.

$$\begin{bmatrix} \frac{dT}{dt} \\ \frac{dR}{dt} \end{bmatrix} = \begin{bmatrix} -1 & 1 \\ 4 & -1 \end{bmatrix} \begin{bmatrix} T \\ R \end{bmatrix} \Rightarrow \frac{dx}{dt} = \begin{bmatrix} -1 & 1 \\ 4 & -1 \end{bmatrix} x$$

Eigenvalues and Eigenvectors

$$A = \begin{bmatrix} 1 & -1 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & -3 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ 2 & 2 \end{bmatrix}^{-1} \quad \begin{aligned} y_1(t) &= e^t y_1(0) \\ y_2(t) &= e^{-3t} y_2(0) \end{aligned}$$

eigenvectors eigenvalues

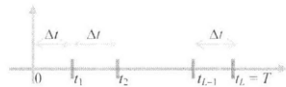


Note the system has a positive eigenvalue
Solutions grow exponentially with time.

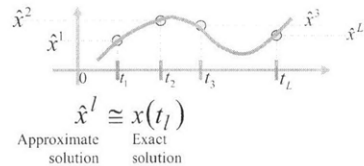
FINITE DIFFERENCE METHODS

Basic Concepts

FIRST: discretize time



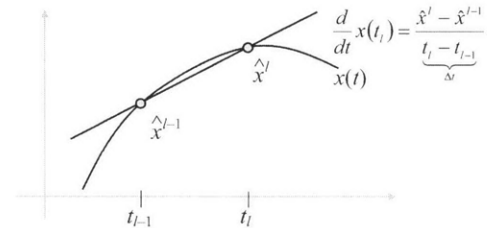
SECOND: Represent $x(t)$ using values at t_i



THIRD: Approximate $\frac{dx}{dt}$ using the discrete \hat{x}^l

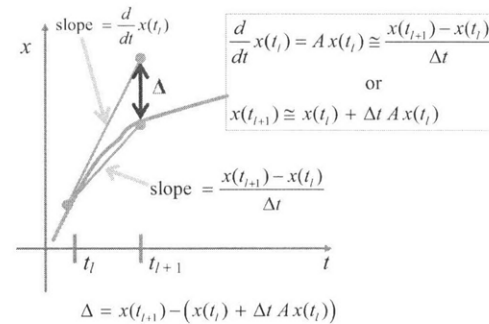
Example: $\frac{d}{dt}x(t) \approx \frac{\hat{x}^{l+1} - \hat{x}^l}{\Delta t_{l+1}} \approx \frac{d}{dt}x(t_{l+1})$

335



$$\frac{x(t_{l+1}) - x(t_l)}{\Delta t} \approx \begin{cases} \frac{dx}{dt}(t_l) & \text{F.E. } \frac{x(t_{l+1}) - x(t_l)}{\Delta t} = Ax(t_l) \\ \frac{dx}{dt}(t_{l+1}) & \text{B.E. } \frac{x(t_{l+1}) - x(t_l)}{\Delta t} = Ax(t_{l+1}) \\ \frac{1}{2} \left[\frac{dx}{dt}(t_l) + \frac{dx}{dt}(t_{l+1}) \right] & \text{TRAP } \frac{x(t_{l+1}) - x(t_l)}{\Delta t} = \dots \\ & = \frac{1}{2} [Ax(t_l) + Ax(t_{l+1})] \end{cases}$$

Forward Euler
o Approximation



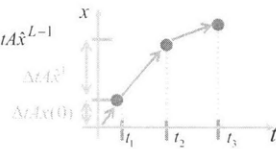
Algorithm

$$x(t_1) \cong \hat{x}^1 = x(0) + \Delta t Ax(0)$$

$$x(t_2) \cong \hat{x}^2 = \hat{x}^1 + \Delta t A \hat{x}^1$$

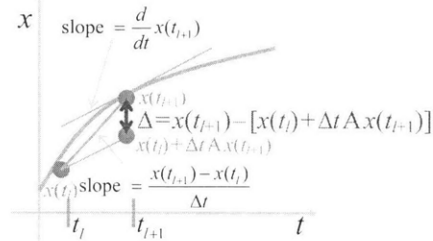
⋮

$$x(t_L) \cong \hat{x}^L = \hat{x}^{L-1} + \Delta t A \hat{x}^{L-1}$$



Backward Euler

Approximation



$$\frac{d}{dt} x(t_{j+1}) = Ax(t_{j+1}) \cong \frac{x(t_{j+1}) - x(t_j)}{\Delta t}$$

or

$$x(t_{j+1}) \cong x(t_j) + \Delta t Ax(t_{j+1})$$

Algorithm

Solve with LU or iterative methods

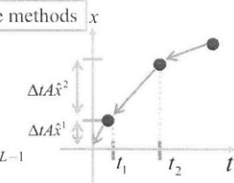
$$x(t_1) \cong \hat{x}^1 = x(0) + \Delta t A \hat{x}^1$$

$$[\mathbf{I} - \Delta t \mathbf{A}] \hat{x}^1 = x(0)$$

$$x(t_2) \cong \hat{x}^2 = [\mathbf{I} + \Delta t \mathbf{A}]^{-1} \hat{x}^1$$

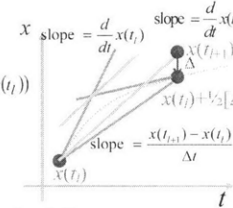
⋮

$$x(t_L) \cong \hat{x}^L = [\mathbf{I} + \Delta t \mathbf{A}]^{-1} \hat{x}^{L-1}$$



Trapezoidal Rule

$$\begin{aligned} \frac{1}{2} \left(\frac{d}{dt} x(t_{j+1}) + \frac{d}{dt} x(t_j) \right) &= \frac{1}{2} (Ax(t_{j+1}) + Ax(t_j)) \\ &\cong \frac{x(t_{j+1}) - x(t_j)}{\Delta t} \end{aligned}$$



$$x(t_{j+1}) \cong x(t_j) + \frac{1}{2} \Delta t Ax(t_{j+1}) + x(t_j)$$

$$\Delta = (x(t_{j+1}) - \frac{1}{2} \Delta t Ax(t_j) - (x(t_j) + \frac{1}{2} \Delta t Ax(t_{j+1})))$$

Algorithm

Solve with LU or iterative methods

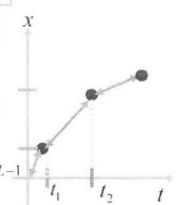
$$x(t_1) \cong \hat{x}^1 = x(0) + \frac{\Delta t}{2} (Ax(0) + A \hat{x}^1)$$

$$\Rightarrow [\mathbf{I} - \frac{\Delta t}{2} \mathbf{A}] \hat{x}^1 = [\mathbf{I} + \frac{\Delta t}{2} \mathbf{A}] x(0)$$

$$x(t_2) \cong \hat{x}^2 = [\mathbf{I} - \frac{\Delta t}{2} \mathbf{A}]^{-1} [\mathbf{I} + \frac{\Delta t}{2} \mathbf{A}] \hat{x}^1$$

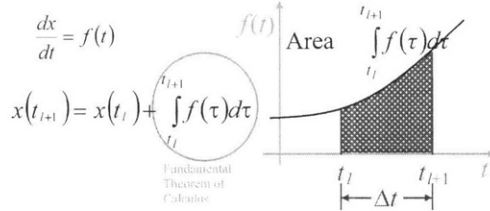
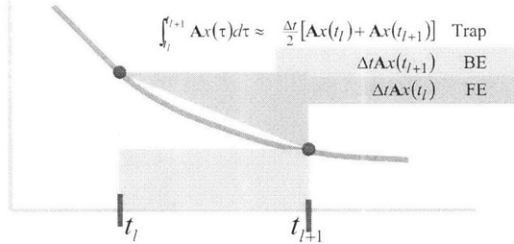
⋮

$$x(t_L) \cong \hat{x}^L = [\mathbf{I} - \frac{\Delta t}{2} \mathbf{A}]^{-1} [\mathbf{I} + \frac{\Delta t}{2} \mathbf{A}] \hat{x}^{L-1}$$

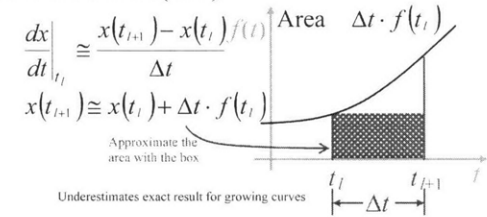


❖ Numerical Integration View

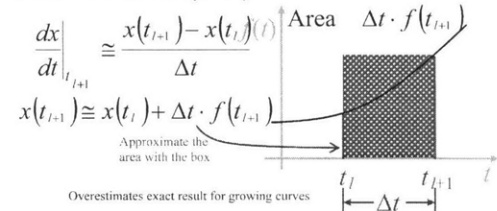
$$\frac{d}{dt}x(t) = Ax(t) \Rightarrow x(t_{i+1}) = x(t_i) + \int_{t_i}^{t_{i+1}} Ax(\tau) d\tau$$



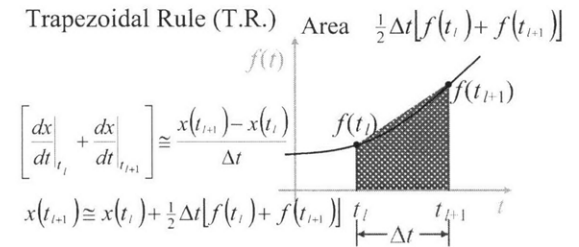
Forward Euler (F.E.)

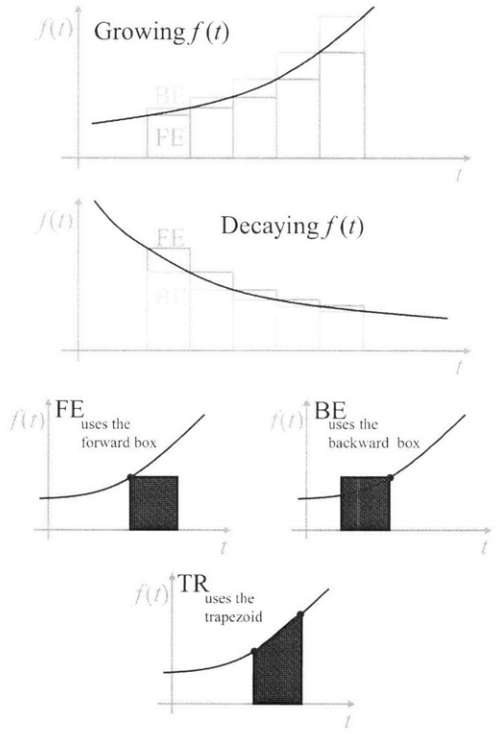


Backward Euler (B.E.)



Trapezoidal Rule (T.R.)

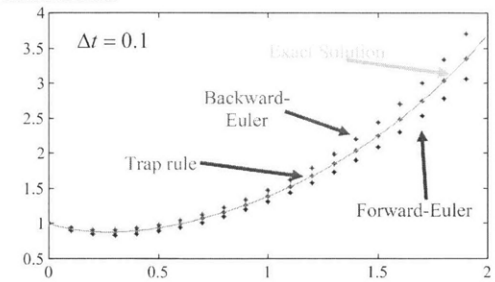




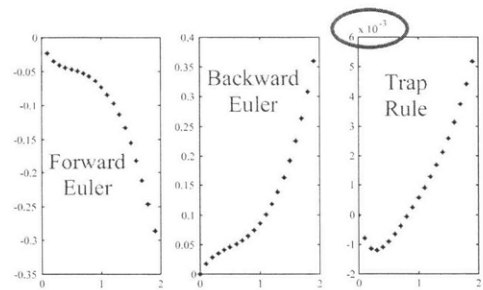
- ❖ Summary
 - Trap Rule, Forward-Euler, Backward-Euler
 - All are one-step methods
 - \hat{x}^l is computed using only \hat{x}^{l-1} , not \hat{x}^{l-2} , \hat{x}^{l-3} , etc.
 - Forward-Euler is simplest
 - No equation solution \Rightarrow explicit method
 - Boxcar approximation to integral
 - Backward-Euler is more expensive
 - Equation solution each step \Rightarrow implicit method
 - Trapezoidal Rule might be more accurate
 - Equation solution each step \Rightarrow implicit method
 - Trapezoidal approximation to integral

Properties – Numerical Experiments

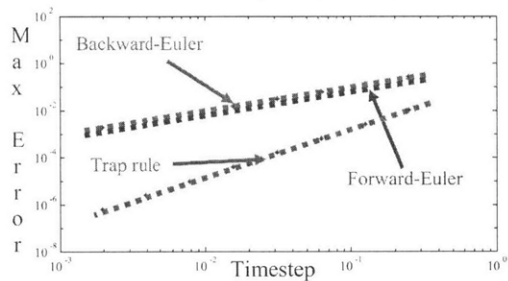
- ❖ Unstable Reaction



FE and BE results have larger errors than Trap Rule, and the errors grow with time.

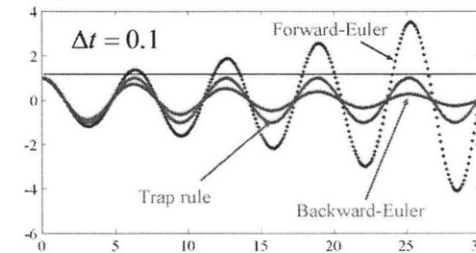


All methods have errors which grow exponentially



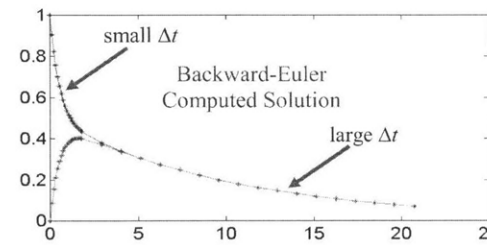
For FE and BE, $error \propto \Delta t$ For Trap, $error \propto (\Delta t)^2$

❖ Oscillating Strut and Mass

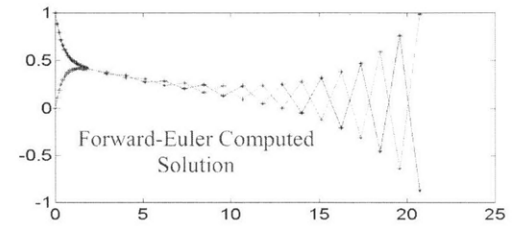


Why does FE result grow, BE result decay and the Trap rule preserves oscillations?

❖ Two timescale RC Circuit



With Backward-Euler it is easy to use small timesteps for the fast dynamics and then switch to large timesteps for the slow decay



The Forward-Euler is accurate for small timesteps, but goes unstable when the timestep is enlarged

	F.E. explicit	B.E. implicit	Trap implicit
COST	cheap	expensive	expensive
ACCURACY	ok	ok	great
CONVERGENCE	linear	linear	quadratic
OSCILLATION	generates	dissipates	conserves
ADJUST Δt	unstable	ok	ok

❖ Summary

- Convergence
 - Did the computed solution approach the exact solution?
 - Why did the trap rule approach faster than BE or FE?
- Energy Preservation
 - Why did BE produce a decaying oscillation?
 - Why did FE produce a growing oscillation?
 - Why did trap rule maintain oscillation amplitude?
- Two timeconstant (stiff) problems
 - Why did FE go unstable when the timestep increased?

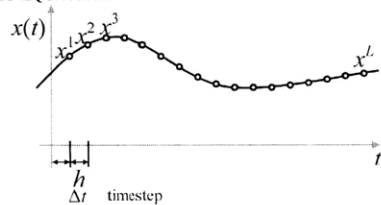
INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 15.
Multistep Methods I

TODAY'S OUTLINE:

- ❖ Multistep Methods
 - e.g. FE, BE, TR
 - Convergence Definitions
 - Local Truncation Error (LTE) and Global Error (GE)
 - Example: LTE and GE for F.E.
 - Minimizing LTE (consistency)
 - Minimizing the accumulation of errors (stability)

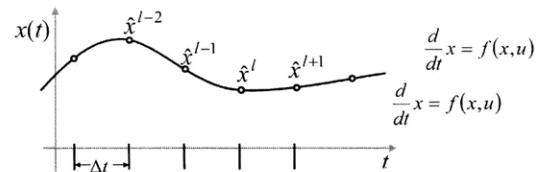
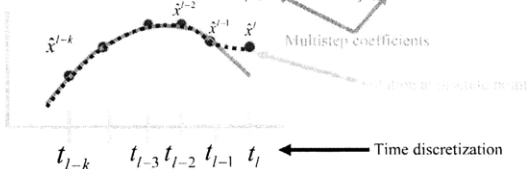
BASIC EQUATIONS



General Notation

Nonlinear Differential Equation: $\frac{d}{dt}x(t) = f(x(t), u(t))$

k-Step Multistep Approach: $\sum_{j=0}^k \alpha_j \hat{x}^{l-j} = \Delta t \sum_{j=0}^k \beta_j f(\hat{x}^{l-j}, u(t_{l-j}))$



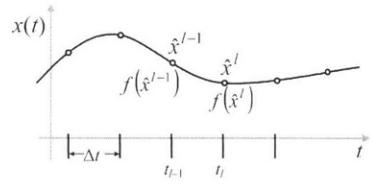
Forward Euler $\hat{x}^{l+1} = \hat{x}^l + \Delta t \cdot \dot{x}(l\Delta t) = \hat{x}^l + \Delta t \cdot f(\hat{x}^l)$

Backward Euler $\hat{x}^{l+1} = \hat{x}^l + \Delta t \cdot \dot{x}(l\Delta t) = \hat{x}^l + \Delta t \cdot f(\hat{x}^{l+1})$

Trapezoidal Rule $\hat{x}^{l+1} = \hat{x}^l + \Delta t \cdot \left[\frac{1}{2} f(\hat{x}^l) + \frac{1}{2} f(\hat{x}^{l+1}) \right]$

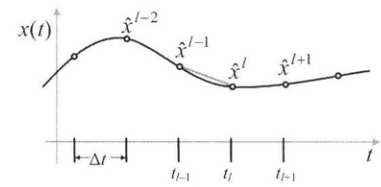
Common Algorithms

Multistep Equation:	$\sum_{j=0}^k \alpha_j \hat{x}^{l-j} = \Delta t \sum_{j=0}^k \beta_j f(\hat{x}^{l-j}, u(t_{l-j}))$
Forward-Euler Approximation:	$x(t_l) \approx x(t_{l-1}) + \Delta t f(x(t_{l-1}), u(t_{l-1}))$
FE Discrete Equation:	$\hat{x}^l - \hat{x}^{l-1} = \Delta t f(x(t_{l-1}), u(t_{l-1}))$
Multistep Coefficients:	$k=1, \alpha_0=1, \alpha_1=-1, \beta_0=0, \beta_1=1$
BE Discrete Equation:	$\hat{x}^l - \hat{x}^{l-1} = \Delta t f(\hat{x}^l, u(t_l))$
Multistep Coefficients:	$k=1, \alpha_0=1, \alpha_1=-1, \beta_0=1, \beta_1=0$
Trap Discrete Equation:	$\hat{x}^l - \hat{x}^{l-1} = \frac{\Delta t}{2} [f(\hat{x}^l, u(t_l)) + f(\hat{x}^{l-1}, u(t_{l-1}))]$
Multistep Coefficients:	$k=1, \alpha_0=1, \alpha_1=-1, \beta_0=\frac{1}{2}, \beta_1=\frac{1}{2}$



$$\alpha_0 \hat{x}^l + \alpha_1 \hat{x}^{l-1} = \Delta t [\beta_0 f(\hat{x}^l) + \beta_1 f(\hat{x}^{l-1})]$$

Forward Euler	1	-1	0	1
Backward Euler	1	-1	1	0
Trapezoidal Rule	1	-1	1/2	1/2



$$\frac{\hat{x}^l - \hat{x}^{l-1}}{\Delta t} = ?$$

Forward Euler

$$\frac{\hat{x}^l - \hat{x}^{l-1}}{\Delta t} = \frac{dx((l-1)\Delta t)}{dt} = f(\hat{x}(t_{l-1})) = f(\hat{x}^{l-1})$$

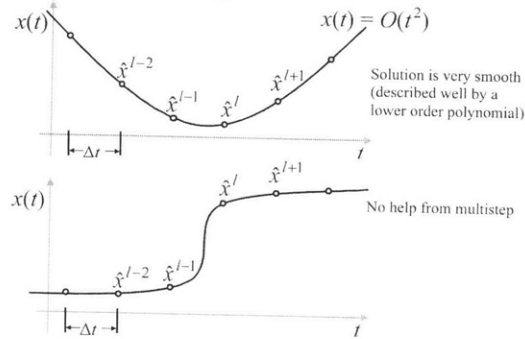
Backward Euler

$$\frac{\hat{x}^l - \hat{x}^{l-1}}{\Delta t} = \frac{dx(l\Delta t)}{dt} = f(\hat{x}^l)$$

Trapezoidal Rule

$$\frac{\hat{x}^l - \hat{x}^{l-1}}{\Delta t} = \frac{1}{2} f(\hat{x}^{l-1}) + \frac{1}{2} f(\hat{x}^l)$$

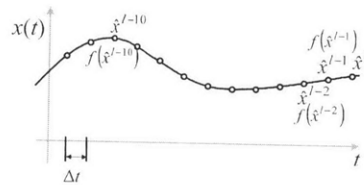
Definitions and Observations



$$\text{Multistep Equation } \sum_{j=0}^k \alpha_j \hat{x}^{l-j} = \Delta t \sum_{j=0}^k \beta_j f(\hat{x}^{l-j}, u(t_{l-j}))$$

1. If $\beta_0 \neq 0$ the multistep method is implicit
2. A k -step multistep method uses k previous x 's and f 's
3. A normalization is needed, $\alpha_0 = 1$ is common
4. A k -step method has $2k+1$ free coefficients

How does one pick good coefficients?



$$\dot{x} = \lambda x \quad x(0) = x_0$$

$k = 1$ Multistep Method

$$\alpha_0 \hat{x}^l + \alpha_1 \hat{x}^{l-1} = \Delta t (\beta_0 \lambda \hat{x}^l + \beta_1 \lambda \hat{x}^{l-1})$$

Forward Euler $\hat{x}^l = \hat{x}^{l-1} + \Delta t \cdot \lambda \hat{x}^{l-1}$ $\alpha_0 = 1 \quad \beta_0 = 0$
 $\alpha_1 = -1 \quad \beta_1 = 1$

When does $x(t)$ grow?

$$\lambda > 0 \quad x(t) = x_0 e^{\lambda t}$$

What about F.E.?

$$\hat{x}^l = (1 + \Delta t \cdot \lambda) \hat{x}^{l-1}$$

\hat{x}^l grows if $(1 + \Delta t \cdot \lambda) > 1$

Simplified Problems for Analysis

Scalar ODE: $\frac{d}{dt} v(t) = \lambda v(t), v(0) = v_0 \quad \lambda \in \mathbb{C}$

Why such a simple Test Problem?

- Nonlinear Analysis has many *unrevealing* subtleties
- Scalar is equivalent to vector for multistep methods

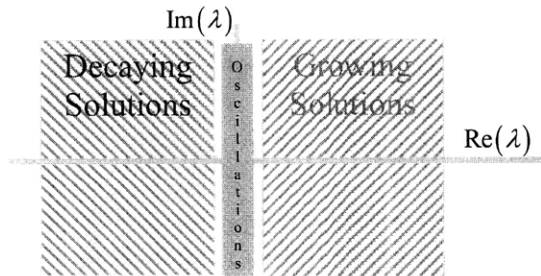
$$\frac{d}{dt} x(t) = Ax(t) \xrightarrow{\text{multistep discretization}} \sum_{j=0}^k \alpha_j \hat{x}^{l-j} = \Delta t \sum_{j=0}^k \beta_j A \hat{x}^{l-j}$$

$$\text{Let } Ey(t) = x(t) \longrightarrow \sum_{j=0}^k \alpha_j \hat{y}^{l-j} = \Delta t \sum_{j=0}^k \beta_j E^{-1} A E \hat{y}^{l-j}$$

$$\text{Decoupled Equations} \longrightarrow \sum_{j=0}^k \alpha_j \hat{y}^{l-j} = \Delta t \sum_{j=0}^k \beta_j \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \hat{y}^{l-j}$$

$$\text{Scalar Multistep formula: } \sum_{j=0}^k \alpha_j \hat{v}^{l-j} = \Delta t \sum_{j=0}^k \beta_j \lambda \hat{v}^{l-j}$$

Must consider ALL $\lambda \in \mathbb{C}$

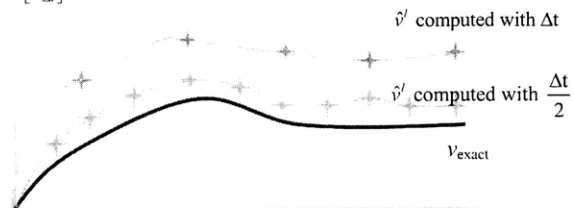


CONVERGENCE ANALYSIS

Convergence Definition

Definition: A multistep method for solving initial value problems on $[0, T]$ is said to be convergent if given any initial condition

$$\max_{t \in [0, T]} \|\hat{v}^l - v(t\Delta t)\| \rightarrow 0 \text{ as } \Delta t \rightarrow 0$$



Order-p Convergence

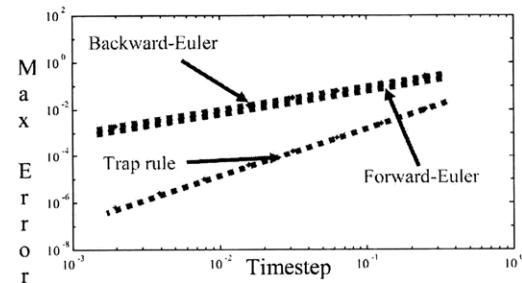
Definition: A multistep method for solving initial value problems on $[0, T]$ is said to be convergent if given any initial condition

$$\max_{t \in [0, T]} \|\hat{v}^l - v(t\Delta t)\| \leq C(\Delta t)^p$$

for all Δt less than a given Δt_0

Forward- and Backward-Euler are order 1 convergent;
Trapezoidal Rule is order 2 convergent

Reaction Equation Example



For FE and BE, $error \propto \Delta t$ For Trap, $error \propto (\Delta t)^2$

LOCAL TRUNCATION ERROR (LTE) AND GLOBAL ERROR (GE)

Convergence Analysis

- ❖ Two Conditions for Convergence
 1. Local Condition: "One step" errors are small (consistency)
Typically verified using Taylor Series
 2. Global Condition: The single step errors do not grow too quickly (stability)
All one-step ($k = 1$) methods are stable in this sense.
Multi-step ($k > 1$) methods require careful analysis

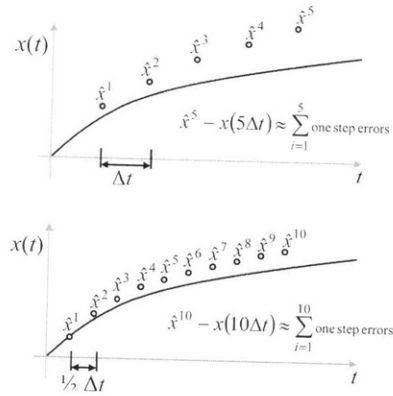
$$\hat{x}^l = \hat{x}^{l-1} + \Delta t \cdot \lambda \hat{x}^{l-1} \Rightarrow \hat{x}^l = (1 + \Delta t \cdot \lambda) \hat{x}^{l-1}$$

$$x(t\Delta t) = x((l-1)\Delta t) + \Delta t \cdot \lambda x((l-1)\Delta t) + O(\Delta t^2)$$

Forward Euler Global Error

$$E^l = E^{l-1} + \Delta t \cdot \lambda E^{l-1} + e^l$$

$$E^l - (1 + \Delta t \cdot \lambda) E^{l-1} = e^l$$



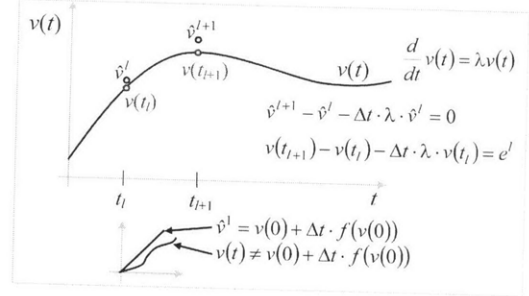
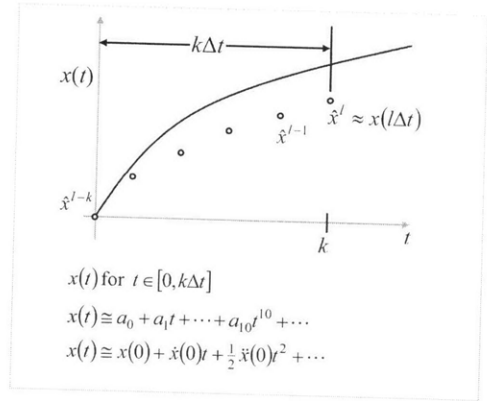
❖ Global Error Equation

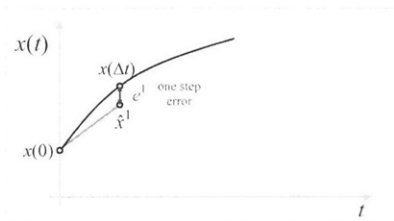
Multistep formula : $\sum_{j=0}^k \alpha_j \hat{v}^{l-j} - \Delta t \sum_{j=0}^k \beta_j \lambda \hat{v}^{l-j} = 0$

Exact solution Almost satisfies Multistep Formula : $\sum_{j=0}^k \alpha_j v(t_{l-j}) - \Delta t \sum_{j=0}^k \beta_j \frac{d}{dt} v(t_{l-j}) = e^l$

Global Error : $E^l \equiv v(t_l) - \hat{v}^l$ Local Truncation Error (LTE)

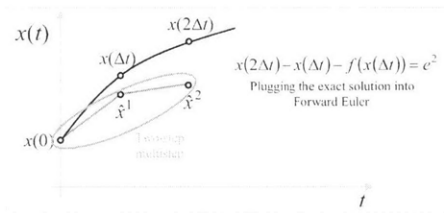
Difference equation relates LTE to Global error $(\alpha_0 - \lambda \Delta t \beta_0) E^l + (\alpha_1 - \lambda \Delta t \beta_1) E^{l-1} + \dots + (\alpha_k - \lambda \Delta t \beta_k) E^{l-k} = e^l$





$$\left. \begin{aligned} & \alpha_0 \left[\underbrace{v(t_i) - \hat{v}^i}_{E^i} \right] - \Delta t \beta_0 \lambda \left(\underbrace{v(t_i) - \hat{v}^i}_{E^i} \right) + \dots \\ & \dots + \alpha_1 \left[\underbrace{v(t_{i-1}) - \hat{v}^{i-1}}_{E^{i-1}} \right] - \Delta t \beta_1 \lambda \left(\underbrace{v(t_{i-1}) - \hat{v}^{i-1}}_{E^{i-1}} \right) + \dots \end{aligned} \right\} k = e^l$$

$$(\alpha_0 - \Delta t \beta_0 \lambda) E^i + (\alpha_1 - \Delta t \beta_1 \lambda) E^{i-1} + \dots = e^l$$



EXAMPLE: LTE AND GE FOR FORWARD-EULER
Convergence Analysis

❖ One Step Methods
 Definition: A one-step method for solving initial value problems on an interval $[0, T]$ is said to be consistent if for any A and any initial condition

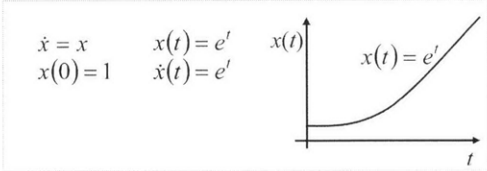
$$\frac{|\hat{x}^1 - x(\Delta t)|}{\Delta t} \rightarrow 0 \text{ as } \Delta t \rightarrow 0$$

$$\hat{x}^l = \hat{x}^{l-1} + \Delta t [\beta_0 f(\hat{x}^l) + \beta_1 f(\hat{x}^{l-1})]$$

Taylor Series Expansion
 $v((l+1)\Delta t) = v(l\Delta t) + \Delta t \frac{dv}{dt}(l\Delta t) + \frac{\Delta t^2}{2} \frac{d^2v}{dt^2}(l\Delta t) + \dots$
 $v((l+1)\Delta t) = v(l\Delta t) + \Delta t \frac{dv}{dt}(l\Delta t) + \frac{\Delta t^2}{2} \frac{d^2v}{dt^2}(\tau)$
 $\tau \in [l\Delta t, (l+1)\Delta t]$

❖ Consistency for Forward Euler
 Forward-Euler definition
 $\hat{v}^{l+1} - \hat{v}^l - \Delta t \lambda \hat{v}^l = 0$
 Substituting the exact $v(t)$ and expanding
 $v((l+1)\Delta t) - v(l\Delta t) - \Delta t \frac{dv}{dt}(l\Delta t) = \frac{(\Delta t)^2}{2} \frac{d^2v}{dt^2}(\tau)$
 $\frac{d}{dt} v = \lambda v$
 $\tau \in [l\Delta t, (l+1)\Delta t]$

where e^l is the LTE and is bounded by
 $|e^l| \leq C(\Delta t)^2$, where $C = 0.5 \max_{\tau \in [0, T]} \left| \frac{d^2v(\tau)}{dt^2} \right|$

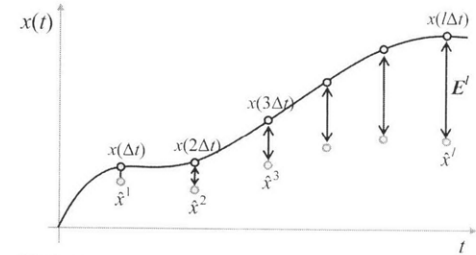


$$\begin{aligned}
 E^0 &= 0 \\
 E^1 &= e^1 \\
 E^2 &= (1 + \Delta t \lambda) E^1 + e^2 = (1 + \Delta t \lambda) e^1 + e^2 \\
 E^3 &= (1 + \Delta t \lambda)^2 e^1 + (1 + \Delta t \lambda) e^2 + e^3 \\
 E^l &\approx \sum_{i=1}^l e^i
 \end{aligned}$$

If $|E^l|$ roughly $\sum e^i \leq \sum C(\Delta t)^2$

$$\left| E^l \right|_{\text{roughly}} \approx \frac{T}{\Delta t} \cdot \underbrace{C(\Delta t)^2}_{\substack{\text{size of errors} \\ \text{shrinking factor}}} \cdot \frac{1}{2} \max_{t \in [0, T]} \frac{d^2 v(t)}{dt^2}$$

- ❖ Global Error Equation
 - Forward-Euler definition
 - $v^{l+1} = v^l + \Delta t \lambda v^l$
 - Using the LTE definition
 - $v((l+1)\Delta t) = v(l\Delta t) + \Delta t \lambda v(l\Delta t) + e^l$
 - Subtracting yields global error equation
 - $E^{l+1} = (1 + \Delta t \lambda) E^l + e^l$
 - Using magnitudes and the bound on e^l
 - $|E^{l+1}| = |1 + \Delta t \lambda| |E^l| + |e^l| \leq (1 + \Delta t |\lambda|) |E^l| + C(\Delta t)^2$
- ❖ A helpful bound on difference equations
 - A lemma bounding difference equation solutions
 - If $|u^{l+1}| \leq (1 + \epsilon) |u^l| + b, u^0 = 0, \epsilon > 0$
 - Then $|u^l| \leq \frac{e^{\epsilon l}}{\epsilon} b$
 - $u^1 \leq (1 + \epsilon) u^0 + b = b$
 - $u^2 \leq (1 + \epsilon) b + b$
 - $u^3 \leq (1 + \epsilon)^2 b + (1 + \epsilon) b + b$



One-Step Methods

- ❖ Hints for the Proof
 - To prove, first write u^l as a power series and sum
 - $|u^l| \leq \sum_{j=0}^{l-1} (1 + \epsilon)^j |b| = \frac{1 - (1 + \epsilon)^l}{1 - (1 + \epsilon)} |b|$
 - To finish, note $(1 + \epsilon) \leq e^\epsilon \Rightarrow (1 + \epsilon)^l \leq e^{\epsilon l}$
 - $|u^l| \leq \frac{1 - (1 + \epsilon)^l}{1 - (1 + \epsilon)} |b| = \frac{(1 + \epsilon)^l - 1}{\epsilon} |b| \leq \frac{e^{\epsilon l}}{\epsilon} |b|$
 - $u^0 = 0 \quad b = 1 \quad \epsilon = 0.1$
 - $|u^1| \leq b = 1$
 - $|u^2| \leq 1.1 \cdot 1 + 1$
 - $|u^3| \leq 1.1 \cdot |u^2| + 1 = 1.1 \cdot (1.1 \cdot 1 + 1) + 1 = 1.1 \cdot 1.1 + 1.1 + 1 = 1.1 \cdot 1.1 \cdot 1 + 1.1 \cdot 1 + 1$
 - $|u^l| \leq \underbrace{(1 + \epsilon)^l}_{\text{contribution from first step}} b + (1 + \epsilon)^{l-1} b + \dots + b$
- ❖ Back to Forward Euler Convergence Analysis
 - Applying the lemma and canceling terms
 - $|E^{l+1}| \leq \left(1 + \frac{\Delta t |\lambda|}{\epsilon} \right) |E^l| + \frac{C(\Delta t)^2}{b} \leq \frac{e^{|\Delta t \lambda|}}{\Delta t |\lambda|} C(\Delta t)^2$
 - Finally noting that $\Delta t \leq T$,
 - $\max_{t \in [0, T]} |E^l| \leq e^{\lambda T} \frac{C}{\lambda} \Delta t$

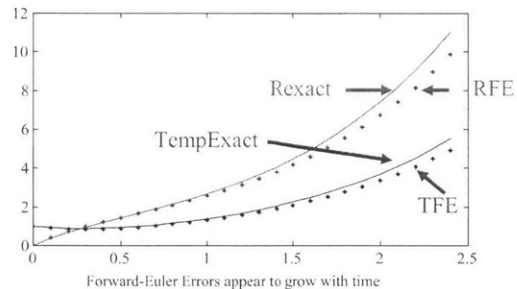
Convergence Analysis

- ❖ Observations about the forward-Euler analysis

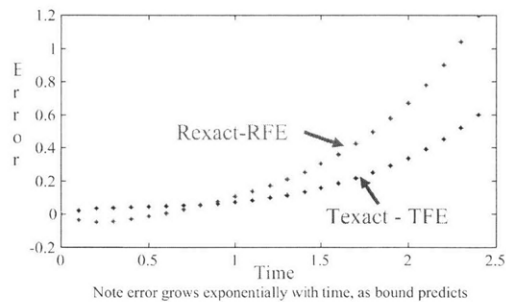
$$\max_{t \in [0, L]} |E^t| \leq e^{\lambda T} \frac{C}{|\lambda|} \Delta t$$

- Forward-Euler is order 1 convergent
- Bound grows exponentially with time interval
- C related to exact solution's second derivative

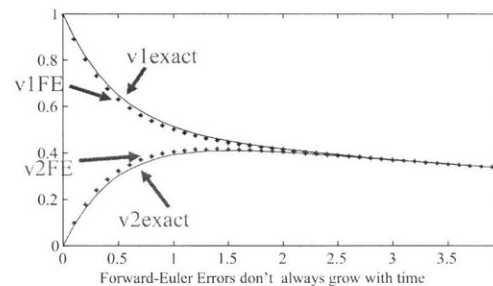
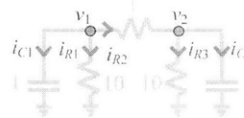
- ❖ Exact and Forward-Euler (FE) Plots for Unstable Reactions



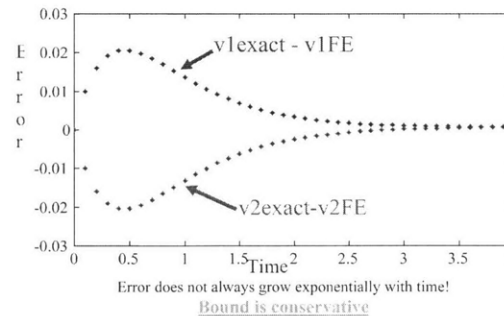
- ❖ Forward-Euler Errors for Solving Reaction Equation



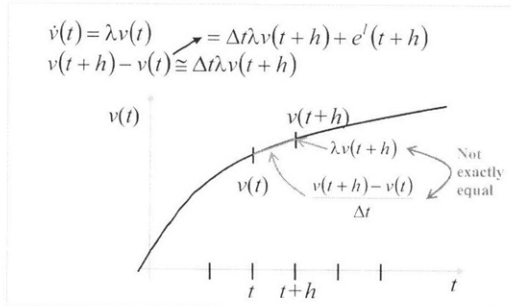
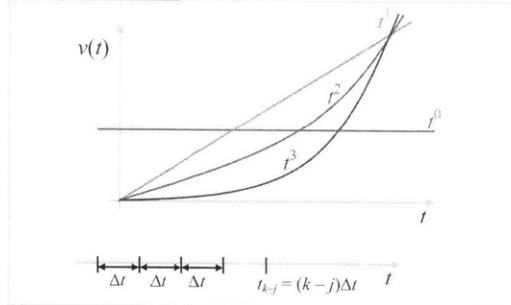
- ❖ Exact and forward-Euler (FE) plots for circuit.



- ❖ Forward-Euler (FE) errors for solving circuit equation.



MINIMIZING LTE (CONSISTENCY)
Exactness Constraints



$v(t_l) - v(t_{l-1}) - \Delta t \lambda v(t_l) = e^l$
 $\hat{v}^l - \hat{v}^{l-1} - \Delta t \lambda \hat{v}^l = 0$

Local Truncation Error: $\sum_{j=0}^k \alpha_j v(t_{l-j}) - \Delta t \sum_{j=0}^k \beta_j \frac{d}{dt} v(t_{l-j}) = e^l$

If $v(t) = t^p \Rightarrow \frac{d}{dt} v(t) = p t^{p-1}$ Local Truncation Error (LTE)

$$\sum_{j=0}^k \frac{\alpha_j ((k-j)\Delta t)^p}{v(t_{k-j})} - \Delta t \sum_{j=0}^k \frac{\beta_j p ((k-j)\Delta t)^{p-1}}{\frac{d}{dt} v(t_{k-j})} = e^k$$

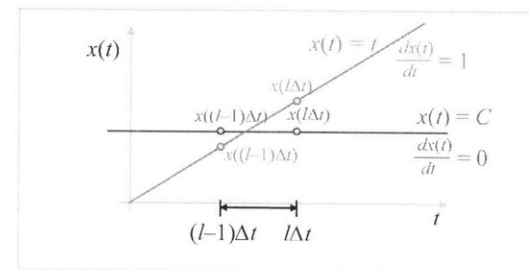
$$\Rightarrow (\Delta t)^p \left[\sum_{j=0}^k \alpha_j (l-j)^p - \sum_{j=0}^k \beta_j p (l-j)^{p-1} \right] = e^k$$

If $\sum_{j=0}^k \alpha_j (k-j)^p - \sum_{j=0}^k \beta_j p (k-j)^{p-1} = 0$ then $e^k = 0$ for $v(t) = t^p$

As any smooth $v(t)$ has a locally accurate Taylor series in t :

If $\sum_{j=0}^k \alpha_j (k-j)^p - \sum_{j=0}^k \beta_j p (k-j)^{p-1} = 0$ for all $p \leq p_0$

Then $\sum_{j=0}^k \alpha_j v(t_{l-j}) - \sum_{j=0}^k \beta_j \frac{d}{dt} v(t_{l-j}) = e^l = \mathcal{O}(\Delta t)^{p_0+1}$



$$\hat{x}^l - \hat{x}^{l-1} + \Delta t f(\hat{x}^l) = 0$$

$$\begin{array}{l} x(t) = 1 \\ 1 - 1 + 0 = 0 \end{array} \qquad \begin{array}{l} x(t) = t \\ l\Delta t - (l-1)\Delta t + \Delta t \cdot 1 = 0 \end{array}$$

Error is zero for both constant and linear functions

❖ Example – Exactness Constraint $k = 2$

Exactness Constraints: $\sum_{j=0}^k \alpha_j (k-j)^p - \sum_{j=0}^k \beta_j p(k-j)^{p-1} = 0$

For $k = 2$, yields a 5×6 system of equations for Coefficients

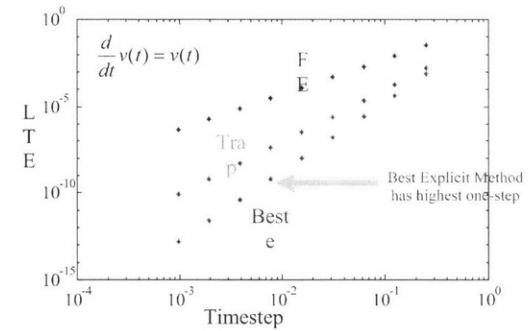
$$\begin{array}{l} p=0 \\ p=1 \\ p=2 \\ p=3 \\ p=4 \end{array} \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & -1 & -1 & -1 \\ 4 & 1 & 0 & -4 & -2 & 0 \\ 8 & 1 & 0 & -12 & -3 & 0 \\ 16 & 1 & 0 & -32 & -4 & 0 \end{bmatrix} \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \alpha_2 \\ \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Note $\sum \alpha_i = 0$ Always

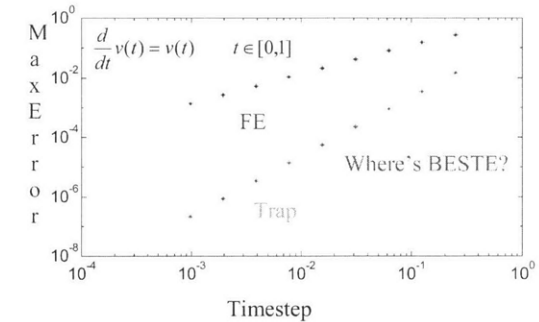
- Forward-Euler $a_0 = 1, a_1 = -1, a_2 = 0, \beta_0 = 0, \beta_1 = 1, \beta_2 = 0$
FE satisfies $p = 0$ and $p = 1$ but not $p = 2 \Rightarrow \text{LTE} = C(\Delta t)^2$
 - Backward-Euler $a_0 = 1, a_1 = -1, a_2 = 0, \beta_0 = 1, \beta_1 = 0, \beta_2 = 0$
BE satisfies $p = 0$ and $p = 1$ but not $p = 2 \Rightarrow \text{LTE} = C(\Delta t)^2$
 - Trapezoidal Rule $a_0 = 1, a_1 = -1, a_2 = 0, \beta_0 = 0.5, \beta_1 = 0.5, \beta_2 = 0$
Trapezoidal satisfies $p = 0, 1, \text{ or } 2$ but not $p = 3 \Rightarrow \text{LTE} = C(\Delta t)^3$
- First introduce a normalization, for example $a_0 = 1$
- Solve for the 2-step method with lowest LTE
 $a_0 = 1, a_1 = 0, a_2 = -1, \beta_0 = 1/3, \beta_1 = 4/3, \beta_2 = 1/3$
Satisfies all five exactness constraints $\Rightarrow \text{LTE} = C(\Delta t)^5$
 - Solve for the 2-step explicit method with lowest LTE
 $a_0 = 1, a_1 = 4, a_2 = -5, \beta_0 = 0, \beta_1 = 4, \beta_2 = 2$
Satisfies all five exactness constraints $\Rightarrow \text{LTE} = C(\Delta t)^4$

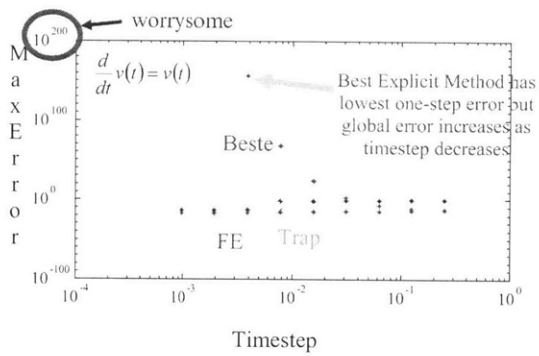
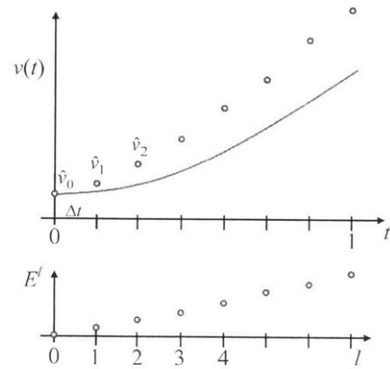
Plots for the FE, Trap, and “Best” Explicit (BESTE)

❖ Local Truncation Error



❖ Global Error





For a "good" method

$$E^l \cong k \sum_{i=1}^l e^i \quad l \propto \frac{T}{\Delta t}$$

FE LTE Δt^2

$$|E^l| \leq kl \max_i |e^i|$$

MINIMIZING THE ACCUMULATION OF ERRORS (STABILITY)

Difference Equation

Why did the "best" 2-step explicit method fail to converge?

Multistep Method Difference Equation

$$(\alpha_0 - \lambda \Delta t \beta_0) E^l + (\alpha_1 - \lambda \Delta t \beta_1) E^{l-1} + \dots + (\alpha_k - \lambda \Delta t \beta_k) E^{l-k} = e^l$$

$v(l\Delta t) - \hat{v}^l$ Global error LTE

We made the LTE so small, how come the Global error is so large?

$\alpha_0 = 1 \quad \beta_0 = 0$ (Explicit)

$$E^l = -(\alpha_1 - \lambda \Delta t \beta_1) E^{l-1} - (\alpha_2 - \lambda \Delta t \beta_2) E^{l-2} - \dots - (\alpha_k - \lambda \Delta t \beta_k) E^{l-k} + e^l$$

$$\begin{bmatrix} E^l \\ E^{l-1} \\ E^{l-2} \\ \vdots \\ E^{l-k+1} \end{bmatrix} = \begin{bmatrix} -(\alpha_1 + \lambda \beta_1) & -(\alpha_2 + \lambda \beta_2) & -(\alpha_3 + \lambda \beta_3) & \dots & -(\alpha_k + \lambda \beta_k) \\ 1 & & & & \\ & 1 & & & \\ & & \ddots & & \\ & & & 1 & 0 \end{bmatrix} \begin{bmatrix} E^{l-1} \\ E^{l-2} \\ \vdots \\ E^{l-k} \end{bmatrix} + \begin{bmatrix} e^l \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Stability Definition

Multistep Method Difference Equation

$$(\alpha_0 - \lambda \Delta t \beta_0) E^l + (\alpha_1 - \lambda \Delta t \beta_1) E^{l-1} + \dots + (\alpha_k - \lambda \Delta t \beta_k) E^{l-k} = e^l$$

Definition: A multistep method is stable if as $\Delta t \rightarrow 0$

$$\max_{l \in [0, T/\Delta t]} |E^l| \leq \frac{C(T)}{\text{interval dependent}} \frac{T}{\Delta t} \max_{l \in [0, T/\Delta t]} |e^l|$$

Stability means: Global error is bounded by a constant times the sum of the L.T.E.'s

Solving Difference Equations

❖ Convolution

Consider a general k^{th} order difference equation

$$a_0 x^l + a_1 x^{l-1} + \dots + a_k x^{l-k} = u^l$$

Which must have k initial conditions

$$x^0 = x_0, x^{-1} = x_1, \dots, x^{-k} = x_k$$

As is clear when the equation is in update form

$$x^l = -\frac{1}{a_0} (a_1 x^{l-1} + \dots + a_k x^{l-k} - u^l)$$

It can be shown that the solution of a difference equation is simply a convolution sum:

$$x \text{ can be related to } u \text{ by } x^l = \sum_{j=0}^l h^{l-j} u^j$$

❖ Calculating h

If $a_0 z^k + a_1 z^{k-1} + \dots + a_k = 0$ has distinct roots $\zeta_1, \zeta_2, \dots, \zeta_k$

$$\text{Then } x^l = \sum_{j=0}^l h^{l-j} u^j \text{ where } h^l = \sum_{j=1}^k \gamma_j (\zeta_j)^l$$

To understand how h is derived, first a simple case

Suppose $x^l = \zeta x^{l-1} + u^l$ and $x^{-1} = 0$

$$x^l = \zeta x^{l-1} + u^l, \quad x^2 = \zeta x^1 + u^2 = \zeta u^1 + u^2$$

$$x^l = \sum_{j=0}^l \zeta^{l-j} u^j$$

time index
exponent

❖ Three Important Observations

If $|\zeta_i| < 1$ for all i , then $|x^l| \leq C \max_j |u^j|$ where C does not depend on l

If $|\zeta_i| > 1$ for any i , then there exists a bounded u^j such that $|x^l| \rightarrow \infty$

If $|\zeta_i| \leq 1$ for all i , and if $|\zeta_i| = 1, \zeta_i$ is distinct then $|x^l| \leq Cl \max_j |u^j|$

$$x^l = \zeta x^{l-1} + u^l \quad x^0 = 0$$

$$x^1 = \zeta x^0 + u^1$$

$$x^2 = \zeta x^1 + u^2 = \zeta u^1 + u^2$$

$$x^l = \zeta^{l-1} u^1 + \zeta^{l-2} u^2 + \zeta^{l-3} u^3 + \dots + u^l$$

❖ Convolution Sum

o When roots are not distinct

$$h^l = \sum_{q=1}^Q \sum_{m=0}^{M_q-1} \gamma_{q,m} (l)^m (\zeta_q)^l \quad x^l = \sum_{j=0}^l h^{l-j} u^j$$

Root multiplicity
Roots of $a_0 z^k + a_1 z^{k-1} + \dots + a_k = 0$

o Bounding Terms

$$x^l = \sum_{q=1}^Q \sum_{m=0}^{M_q-1} \left(\sum_{j=0}^l \gamma_{q,m} (l-j)^m (\zeta_q)^{l-j} u^j \right)$$

$R_{q,m}$

If $|\zeta_i| < 1$, then $|R_{q,m}| \leq C \max_j |u^j|$ Independent of l

If $|\zeta_i| < (1 + \varepsilon)$, then $|R_{q,0}| \leq C \frac{e^{\varepsilon l}}{\varepsilon} \max_j |u^j|$
 Lemma: bounds **distinct** roots

Stability Theorem

Theorem: A multistep method is stable if and only if

Roots of $a_0 z^k + a_1 z^{k-1} + \dots + a_k = 0$ either:

1. Have magnitude less than one
2. Have magnitude equal to one and are distinct

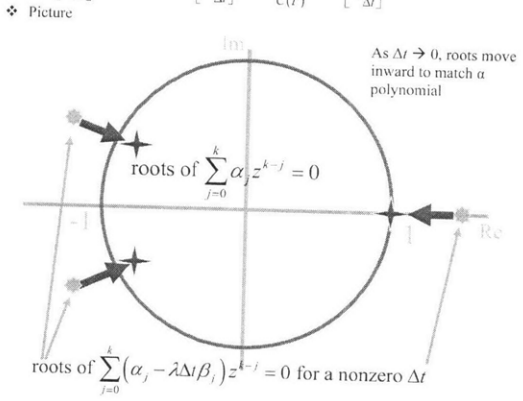
Note: from the previous slides it is easy to see that a multistep method is stable if the same conditions hold for:

$$(a_0 - \lambda \Delta \beta_0) z^k + (a_1 - \lambda \Delta \beta_1) z^{k-1} + \dots + (a_k - \lambda \Delta \beta_k) = 0$$

But the theorem says: forget about the β 's...

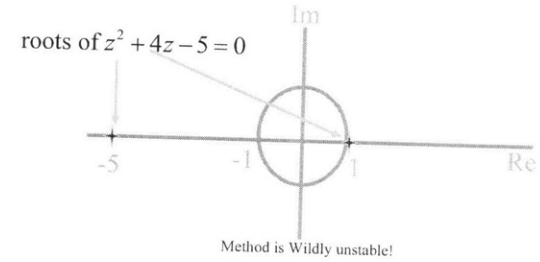
❖ "Proof"
 Given the Multistep Method Difference Equation
 $(\alpha_0 - \lambda \Delta t \beta_0)E^l + (\alpha_1 - \lambda \Delta t \beta_1)E^{l-1} + \dots + (\alpha_k - \lambda \Delta t \beta_k)E^{l-k} = e^l$
 If, as $\Delta t \rightarrow 0$, roots of $(\alpha_0 - \lambda \Delta t \beta_0)z^l + \dots + (\alpha_k - \lambda \Delta t \beta_k)z^{l-k} = 0$
 o less than one in magnitude or
 o are distinct and bounded by $1 + \kappa \Delta t$, $\kappa > 0$
 Then from the aside on difference equations

$$\max_{l \in [0, \frac{T}{\Delta t}]} |E^l| \leq C \frac{e^{\kappa T}}{\Delta t} \max_{l \in [0, \frac{T}{\Delta t}]} |e^l| \leq \frac{C e^{\kappa T} T}{C(T) \Delta t} \max_{l \in [0, \frac{T}{\Delta t}]} |e^l|$$



BESTE Method

Best explicit 2-step method
 $a_0 = 1, a_1 = 4, a_2 = -5, \beta_0 = 0, \beta_1 = 4, \beta_2 = 2$



Dahlquist's First Stability Barrier

For a stable, explicit k -steps method, the maximum degree of a polynomial t^n that can be integrated exactly is less than or equal to k : $p_0 \leq k$ (note there are $2k$ coefficients).

For a stable, implicit k -steps method, (with $2k + 1$ coeff)
 $p_0 \leq k + 2$ if k is even
 $p_0 \leq k + 1$ if k is odd

STABLE:	EXPLICIT	IMPLICIT
$k = 1$ step	$p_0 \leq 1$ (FE)	$p_0 \leq 2$ (TR)
$k = 2$ steps	$p_0 \leq 2$	$p_0 \leq 4$
$k = 3$ steps	$p_0 \leq 3$	$p_0 \leq 4$

Conditions for Convergence, Stability, and Consistency

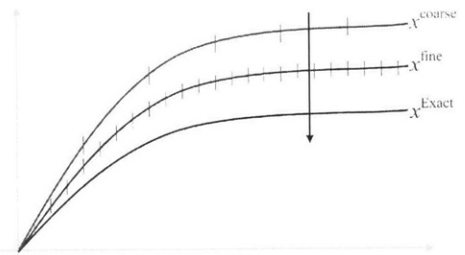
- Local Condition: One step errors are small (consistency)
 Exactness Constraints up to p_0 (p_0 must be > 0)

$$\Rightarrow \max_{l \in [0, \frac{T}{\Delta t}]} |e^l| \leq C_1 (\Delta t)^{p_0+1}$$
 for $\Delta t < \Delta t_0$
- Global Condition: One step errors grow slowly (stability)

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 16.
Multistep Methods II: Small and Large Time-Steps Stability

- TODAY'S OUTLINE:**
- ❖ SMALL Timestep Issues for Multistep Methods
 - Stability
 - Difference Equations
 - Stability & Consistency implies Convergence
 - ❖ LARGE Timestep Issues
 - Absolute Stability for two time-scale examples.
 - Oscillators



STABILITY FOR SMALL TIMESTEPS

Difference Equation

Why did the "best" two-step explicit method fail to converge?
Multistep Method Difference Equation:

$$(\alpha_0 - \lambda \Delta t \beta_0)E^l + (\alpha_1 - \lambda \Delta t \beta_1)E^{l-1} + \dots + (\alpha_k - \lambda \Delta t \beta_k)E^{l-k} = e^l$$

$$v(l\Delta t) - \hat{v}^l \quad \text{Global Error} \quad \text{LTE}$$

We made the LTE so small, how come the Global error is so large?

Stability Definition

Multistep Method Difference Equation
 $(\alpha_0 - \lambda \Delta t \beta_0)E^l + (\alpha_1 - \lambda \Delta t \beta_1)E^{l-1} + \dots + (\alpha_k - \lambda \Delta t \beta_k)E^{l-k} = e^l$

Definition: A multistep method is stable if as $\Delta t \rightarrow 0$

$$\max_{l \in [0, \frac{T}{\Delta t}]} |E^l| \leq \frac{C(T)}{\Delta t} \max_{l \in [0, \frac{T}{\Delta t}]} |e^l|$$

interval dependent

354

Stability means: Global Error is bounded by a constant times the sum of the LTE's

Solving Difference Equations

❖ Convolution
Consider a general k^{th} order difference equation

$$a_0 x^l + a_1 x^{l-1} + \dots + a_k x^{l-k} = u^l$$

Which must have k initial conditions

$$x^0 = x_0, \quad x^{-1} = x_1, \quad \dots \quad x^{-k} = x_k$$

as is clear when the equation is in update form

$$x^l = -\frac{1}{a_0} (a_1 x^{l-1} + \dots + a_k x^{l-k} - u^l)$$

it can be shown that the solution of a difference equation is simply a convolution sum:

$$x \text{ can be related to } u \text{ by } x^l = \sum_{j=0}^l h^{l-j} u^j$$

BE $\dot{v}(t) = \lambda v$

$\hat{v}^l - \hat{v}^{l-1} - \Delta t \lambda \hat{v}^l = 0$

$v(l\Delta t) - v((l-1)\Delta t) - \Delta t \lambda v(l\Delta t) = e^l$

$$(1 - \Delta t \lambda) E^l - E^{l-1} = e^l$$

$$E^l = \frac{1}{1 - \Delta t \lambda} E^0 + \frac{1}{1 - \Delta t \lambda} e^1$$

$$E^2 = \frac{1}{1 - \Delta t \lambda} E^1 + \frac{1}{1 - \Delta t \lambda} e^2 = \left(\frac{1}{1 - \Delta t \lambda}\right)^2 e^1 + \frac{1}{1 - \Delta t \lambda} e^2$$

$$E^l = \left(\frac{1}{1 - \Delta t \lambda}\right)^l e^1 + \left(\frac{1}{1 - \Delta t \lambda}\right)^{l-1} e^2 + \dots + \frac{1}{1 - \Delta t \lambda} e^l$$

$$E^l = \sum_{j=1}^l \left(\frac{1}{1 - \Delta t \lambda}\right)^{l-j} \underbrace{\left(\frac{1}{1 - \Delta t \lambda}\right)}_u e^j$$

❖ Calculating h
 If $a_0 z^k + a_1 z^{k-1} + \dots + a_k = 0$ has distinct roots $\zeta_1, \zeta_2, \dots, \zeta_k$

Then $x^l = \sum_{j=0}^l h^{lj} u^j$ where $h^l = \sum_{j=1}^k \gamma_j (\zeta_j)^l$

To understand how h is derived, first a simple case

Suppose $x^l = \zeta x^{l-1} + u^l$ and $x^{-1} = 0$

$x^1 = \zeta x^0 + u^1 = u^1, \quad x^2 = \zeta x^1 + u^2 = \zeta u^1 + u^2$

$x^l = \sum_{j=0}^l \zeta^{l-j} u^j$
time index ↙ ↘
 ↗ ↘
 time exponent index

❖ Three Important Observations

If $|\zeta_i| < 1$ for all i , then $|x^l| \leq C \max_j |u^j|$ where C does not depend on l

If $|\zeta_i| > 1$ for all i , then there exists a bounded u^l such that $|x^l| \rightarrow \infty$

If $|\zeta_i| \leq 1$ for all i , and i $|\zeta_i| = 1, \zeta_i$ is distinct then $|x^l| \leq Cl \max_j |u^j|$

$$\begin{aligned} x^l &= \lambda x^{l-1} + u^l \\ x^{-1} &= 0 \\ x^0 &= u^0 \\ x^1 &= \lambda x^0 + u^1 = \lambda u^0 + u^1 \\ x^2 &= \lambda x^1 + u^2 = \lambda^2 u^0 + \lambda u^1 + u^2 \\ x^l &= \sum_{i=0}^l \lambda^{l-i} u^i \end{aligned}$$

❖ Convolution Sum

o When roots are not distinct

$h^l = \sum_{q=1}^Q \sum_{m=0}^{M_q-1} \gamma_{q,m} (l)^m (\zeta_q)^l$ $x^l = \sum_{j=0}^l h^{l-j} u^j$

Root Multiplicity ↙ ↘
 ↗ ↘
 Roots of $a_0 z^k + a_1 z^{k-1} + \dots + a_k = 0$

o Bounding Terms

$$x^l = \sum_{q=1}^Q \sum_{m=0}^{M_q-1} \left[\sum_{j=0}^l \gamma_{q,m} (l-j)^m (\zeta_q)^{l-j} u^j \right]$$

$R_{q,m}$

If $|\zeta_i| < 1$, then $|R_{q,m}| \leq C \max_j |u^j|$
Independent of l

If $|\zeta_i| < (1 + \varepsilon)$, then $|R_{q,0}| \leq C \frac{e^{\varepsilon l}}{\varepsilon} \max_j |u^j|$
 Lemma: bounds **distinct** roots

Stability Theorem

Theorem: A multistep method is stable if and only if

Roots of $a_0 z^k + a_1 z^{k-1} + \dots + a_k = 0$ either:

- Have magnitude less than one
- Have magnitude equal to one and are distinct

Note: from the previous slides it is easy to see that a multistep method is stable if the same conditions hold for:

$(\alpha_0 - \lambda \Delta \beta_0) z^k + (\alpha_1 - \lambda \Delta \beta_1) z^{k-1} + \dots + (\alpha_k - \lambda \Delta \beta_k) = 0$

But the theorem says: forget about the betas...

❖ "Proof"

Given the Multistep Method Difference Equation

$(\alpha_0 - \lambda \Delta \beta_0) E^l + (\alpha_1 - \lambda \Delta \beta_1) E^{l-1} + \dots + (\alpha_k - \lambda \Delta \beta_k) E^{l-k} = e^l$

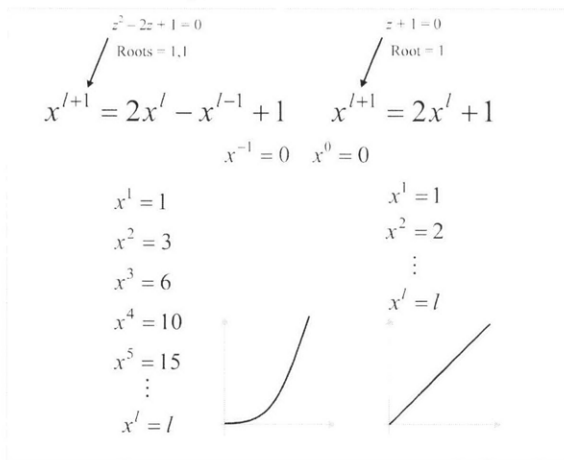
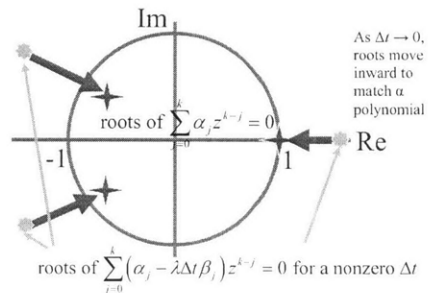
If, as $\Delta t \rightarrow 0$, roots of $(\alpha_0 - \lambda \Delta \beta_0) z^k + \dots + (\alpha_k - \lambda \Delta \beta_k) = 0$

- o Less than one in magnitude or
- o Are distinct and bounded by $1 + \kappa \Delta t, \kappa > 0$

Then from the aside on difference equations

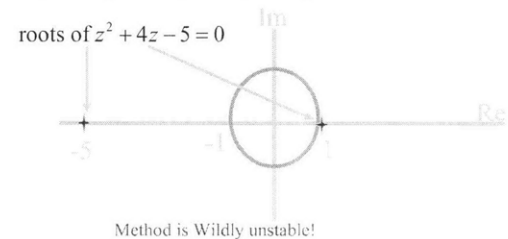
$$\max_{l \in \left[0, \frac{T}{\Delta t}\right]} |E^l| \leq C \frac{e^{\kappa T}}{\Delta t} \max_{l \in \left[0, \frac{T}{\Delta t}\right]} |e^l| \leq \frac{C e^{\kappa T}}{C(T)} \frac{T}{\Delta t} \max_{l \in \left[0, \frac{T}{\Delta t}\right]} |e^l|$$

❖ Picture



BESTE Method

Best explicit 2-step method
 $\alpha_0 = 1, \alpha_1 = 4, \alpha_2 = -5, \beta_0 = 0, \beta_1 = 4, \beta_2 = 2$



Dahlquist's First Stability Barrier

For a stable, explicit k -steps method, the maximum degree of a polynomial f' that can be integrated exactly is less than or equal to k : $p_0 \leq k$ (note there are $2k$ coefficients).

For a stable, implicit k -steps method, (with $2k+1$ coefficients)
 $p_0 \leq k+2$ if k is even
 $p_0 \leq k+1$ if k is odd

STABLE:	EXPLICIT:	IMPLICIT:
$k=1$ step	$p_0 \leq 1$ (FE)	$p_0 \leq 2$ (TR)
$k=2$ steps	$p_0 \leq 2$	$p_0 \leq 4$
$k=3$ steps	$p_0 \leq 3$	$p_0 \leq 4$

Conditions for Convergence, Stability, and Consistency

❖ Local Condition: One step errors are small (consistency)

Exactness Constrains up to p_0 (p_0 must be > 0)

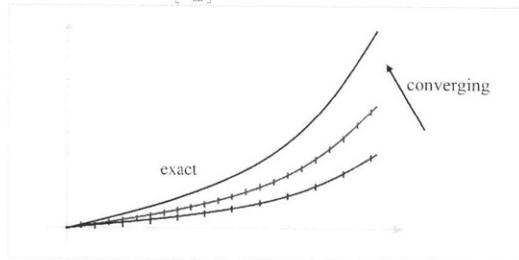
$$\Rightarrow \max_{l \in [0, \frac{T}{\Delta t}]} \|E^l\| \leq C_1 (\Delta t)^{p_0+1} \text{ for } \Delta t < \Delta t_0$$

❖ Global Condition: One step errors grow slowly (stability)

roots of $\sum_{j=0}^k \alpha_j z^{k-j} = 0$ Inside the unit circle or on the unit circle and distinct

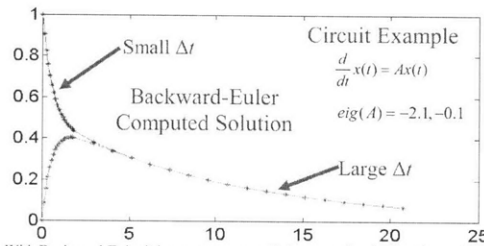
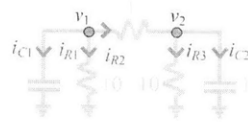
$$\Rightarrow \max_{l \in [0, \frac{T}{\Delta t}]} \|E^l\| \leq C_2 \frac{T}{\Delta t} \max_{l \in [0, \frac{T}{\Delta t}]} \|e^l\|$$

Convergence Result : $\max_{t \in [0, \frac{T}{\Delta t}]} \|E^t\| \leq CT(\Delta t)^{p_0}$



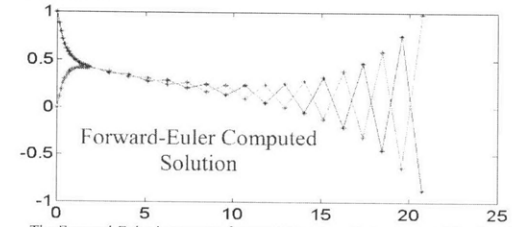
LARGE TIMESTEP ISSUES FOR MULTISTEP METHODS
Absolute Stability for two time-scale examples

❖ Backward-Euler



With Backward-Euler it is easy to use small timesteps for the fast dynamics and then switch to large timesteps for the slow decay.

❖ Forward-Euler



The Forward-Euler is accurate for small timesteps, but goes unstable when the timestep is enlarged.

❖ FE, BE, and Trap on the scalar ODE problem

Scalar ODE : $\frac{d}{dt} v(t) = \lambda v(t), v(0) = v_0, \lambda \in \mathbb{C}$

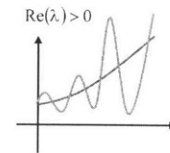
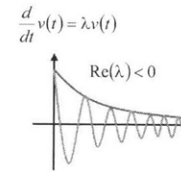
Forward - Euler : $v^{l+1} = v^l + \Delta t \lambda v^l = (1 + \Delta t \lambda) v^l$

If $|1 + \Delta t \lambda| > 1$ the solution grows even if $\lambda < 0$

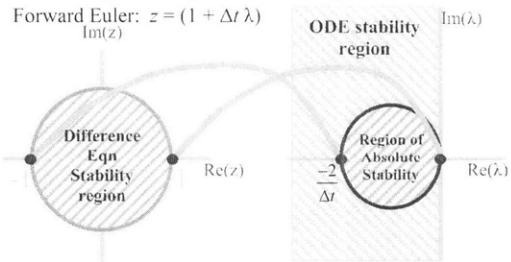
Backward - Euler : $v^{l+1} = v^l + \Delta t \lambda v^{l+1} = \frac{1}{1 - \Delta t \lambda} v^l$

If $|\frac{1}{1 - \Delta t \lambda}| < 1$ the solution decays even if $\lambda < 0$

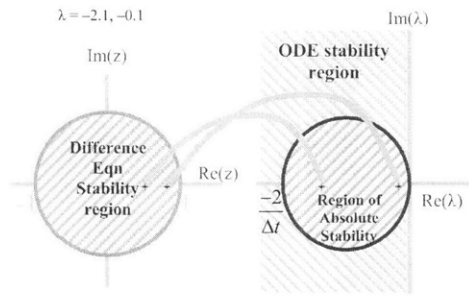
Trapezoidal Rule : $v^{l+1} = v^l + \frac{1}{2} \Delta t \lambda (v^{l+1} + v^l) \Rightarrow v^{l+1} = \frac{(1 + \frac{1}{2} \Delta t \lambda)}{(1 - \frac{1}{2} \Delta t \lambda)} v^l$



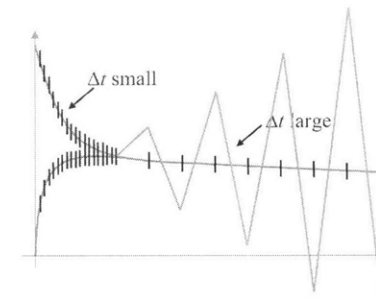
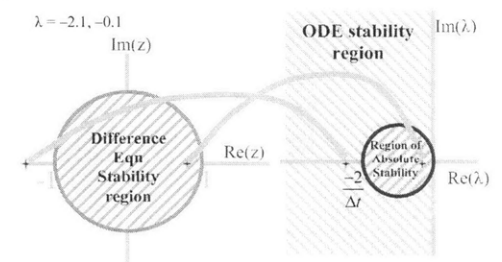
❖ Forward-Euler large timestep region of absolute stability



○ Circuit Example
▪ $\Delta t = 0.1$

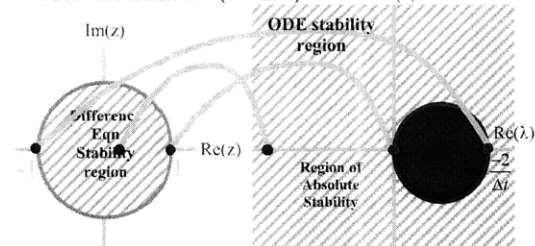


▪ $\Delta t = 1.0$

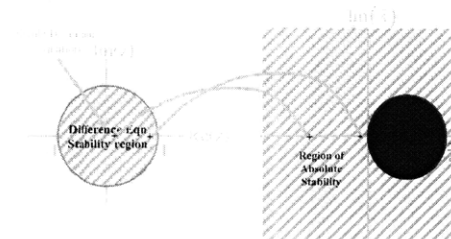


❖ Backward-Euler large timestep region of absolute stability

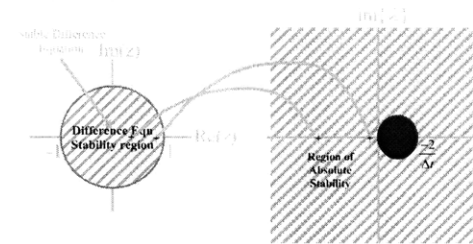
Backward Euler: $z = (1 + \Delta t \lambda)^{-1}$



- Circuit Example
 - $\Delta t = 0.1, \lambda = -2.1, -0.1$



- $\Delta t = 1.0, \lambda = -2.1, -0.1$



❖ Stability Definitions

Region of Absolute Stability for a Multistep method:

Values of $\lambda \Delta t$ where roots of $\sum_{j=0}^k (\alpha_j - \lambda \Delta t \beta_j) z^{k-j} = 0$ are inside the unit circle.

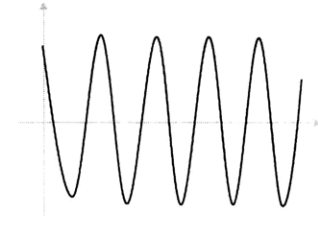
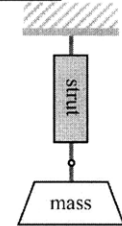
A-Stable:

A method is A-stable if its region of absolute stability includes the entire left-half of the complex plane

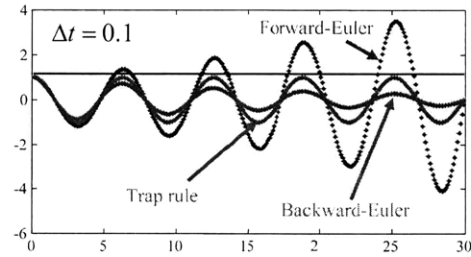
Dahlquist's Second Stability Barrier:

There are no A-stable multistep methods of convergence order greater than 2, and the trap rule is the most accurate.

Oscillators



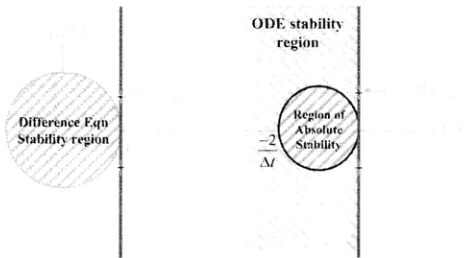
❖ Numerical Experiments – oscillating strut and mass



Why does FE result grow, BE result decay and the Trap rule preserve oscillations

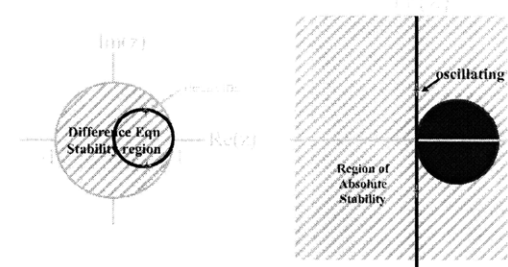
❖ FE Large Timestep Oscillator Example

$$z = 1 + \Delta t \lambda$$



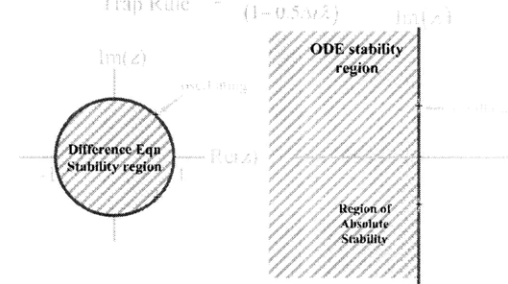
❖ BE Large Timestep Oscillator Example

$$\text{Backward Euler } z = \frac{1}{1 - \Delta t \lambda}$$



❖ Trap Large Timestep Oscillator Example

$$\text{Trap Rule } z = \frac{1 + 0.5\Delta t \lambda}{1 - 0.5\Delta t \lambda}$$



❖ Large Timestep Issues

- Two Time-Constant Stable Problem (Circuit)
 - FE: stability, not accuracy, limited timestep size
 - BE: was A-stable, any timestep could be used.
 - Trap Rule most accurate A-stable m-step method
- Oscillator Problem
 - Forward-Euler generated an unstable difference equation regardless of timestep size.

- Backward-Euler generated a stable (decaying) difference equation regardless of timestep size.
- Trapezoidal rule mapped the imaginary axis to the unit circle, regardless of timestep size:
 - Decaying ODE are mapped to stable difference equations
 - Unstable ODE are mapped to unstable difference equations
 - Oscillating ODE are mapped to oscillating difference equations

SUMMARY

- ❖ Small Timestep Issues for Multistep Methods
 - Local truncation error and exactness
 - Difference equation stability
 - Stability + consistency implies convergence
- ❖ Large Timestep Issues
 - Absolute stability for two time-scale examples
 - Oscillators

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 17.

Methods for Computing Periodic Steady-State

TODAY'S OUTLINE:

- ❖ Periodic Steady-State Problems
 - Application examples and simple cases
- ❖ Review Full Time Integration Methods
- ❖ Finite-Difference Methods for Periodic Steady-State
- ❖ Shooting Methods
 - State Transition Function
 - Sensitivity Matrix
 - Matrix Free Approach
- ❖ Spectral Methods (Harmonic Balance)
 - Galerkin Methods
 - Collocation Methods

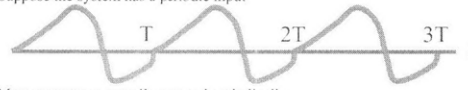
PERIODIC STEADY-STATE PROBLEMS

Basics

❖ Definition

$$\frac{dx(t)}{dt} = F\left(\underbrace{x(t)}_{\text{state}}\right) + \underbrace{u(t)}_{\text{input}}$$

o Suppose the system has a periodic input

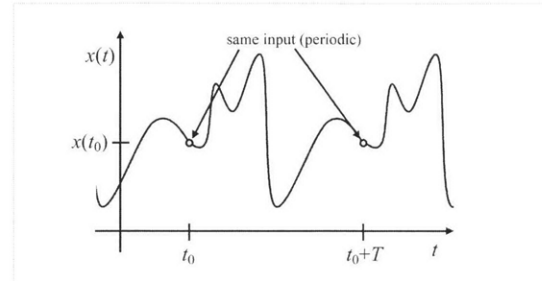


o Many systems eventually respond periodically
 $x(t+T) = x(t)$ for $t \gg 0$

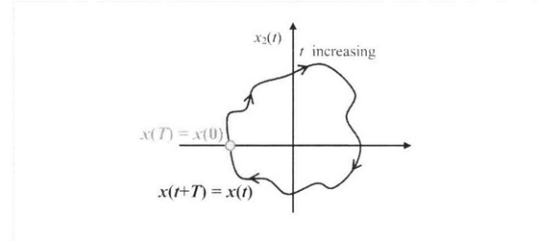
❖ Interesting Property

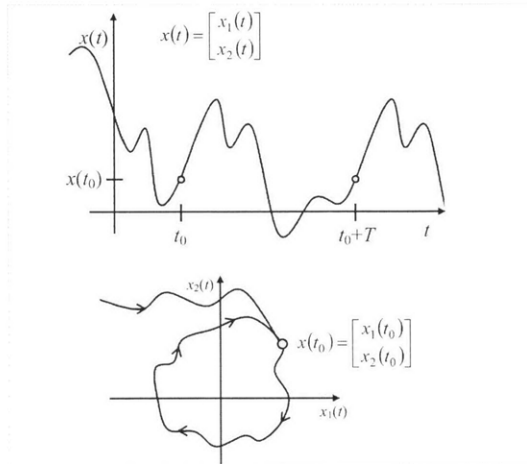
$$\dot{x} = F(x(t), u(t)) \quad x(0) = x_0 \Rightarrow x(t) \text{ is unique}$$

Assuming F is "nice"
smooth enough



- o If x satisfies a differential equation which has a unique solution for any initial condition
 $\frac{dx(t)}{dt} = F(x(t)) + u(t)$
- o Then if u is periodic with period T and $x(t_0+T) = x(t_0)$ for some t_0
 $\Rightarrow x(t+T) = x(t)$ for all $t > t_0$



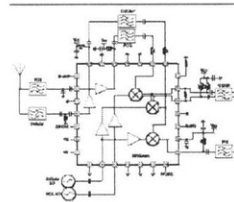


Application Examples

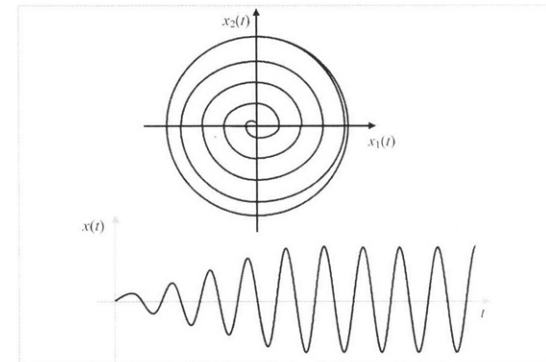
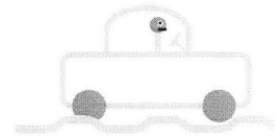
- ❖ Swaying Bridge
 - Periodic Input
 - Wind
 - Response
 - Oscillating Platform
 - Desired Info
 - Oscillation Amplitude

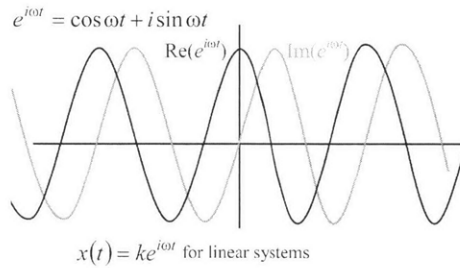


- ❖ Communication Integrated Circuit
 - Periodic Input
 - Received signal at 900MHz
 - Response
 - Filtered Demodulated Signal
 - Desired Info
 - Distortion

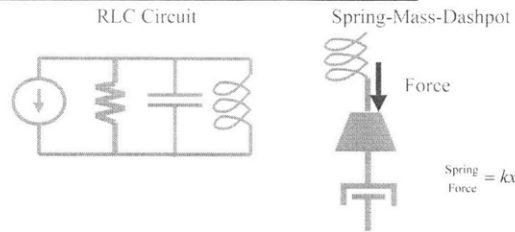


- ❖ Automobile Vibration
 - Periodic Input
 - Regularly Spaced Road Bumps
 - Response
 - Car Shakes
 - Desired Info
 - Shake amplitude





Simple Example [RLC Filter, Spring+Mass+Dashpot]



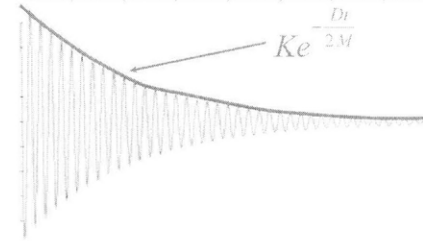
- Both described by second-order ODE

$$\mathbf{M} \frac{d^2 x}{dt^2} + \mathbf{D} \frac{dx}{dt} + x = \underbrace{u(t)}_{\text{input}}$$
- $u(t) = 0$ lightly damped ($\mathbf{D} \ll \mathbf{M}$) response

$$x(t) \approx \mathbf{K} e^{-\frac{\mathbf{D}}{2\mathbf{M}}t} \cos\left(\frac{t}{\sqrt{\mathbf{M}}} + \phi\right)$$

$u(t) = e^{j\omega t} \xrightarrow{\text{LTI}} x(t) = \mathbf{X}_0 e^{j\omega t} \quad \mathbf{X}_0 = ?$
 $j\omega \mathbf{X}_0 e^{j\omega t} = \mathbf{A} \mathbf{X}_0 e^{j\omega t} + b e^{j\omega t} \quad \forall t$
 $(j\omega \mathbf{I} - \mathbf{A}) \mathbf{X}_0 = b$
 $\mathbf{X}_0 = (j\omega \mathbf{I} - \mathbf{A})^{-1} b$
 $x(t) = (j\omega \mathbf{I} - \mathbf{A})^{-1} b e^{j\omega t} = \mathbf{X}_0 e^{j\omega t}$
 $\frac{dx}{dt} = j\omega \mathbf{X}_0 e^{j\omega t}$

- A lightly damped system oscillates many times before settling to a steady-state



Computing Steady-State

$\frac{dx}{dt} \Big|_{t=l\Delta t} \cong \frac{x(l\Delta t) - x((l-1)\Delta t)}{\Delta t} \quad \text{B.E. formula}$
 $\frac{dx}{dt} \Big|_{t=l\Delta t} \cong F \left(\frac{x(l\Delta t)}{f(x(l\Delta t), u(l\Delta t))} + u(l\Delta t) \right) \quad \text{O.D.E.}$

- ❖ Frequency Domain Approach
 - Sinusoidally excited linear time-invariant system

$$\frac{dx}{dt} = Ax(t) + bu(t) \quad \text{input } u(t) = e^{j\omega t}$$
 - Steady-state solution simple to determine

$$x(t) = (j\omega I - A)^{-1} b e^{j\omega t}$$
 Not useful for nonlinear or time-varying systems

Phaser Analysis

input = $e^{j\omega t}$ complex constant

$\dot{x} = Ax + bu(t)$

Assume $x(t) = He^{j\omega t}$ Jumping to periodic steady state

$j\omega He^{j\omega t} = AHe^{j\omega t} + be^{j\omega t}$

Works well for linear problems

Brut Force

$x(0) = 0 \quad u(t) = \sin t$

Forward Euler $\hat{x}^1 = \frac{x(0) + \Delta t(Ax(0) + u(0))}{0}$

$\hat{x}^2 = \hat{x}^1 + \Delta t(A\hat{x}^1 + u(\Delta t))$

$x(t)$ periodic steady state

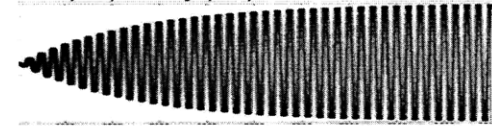
REVIEW FULL TIME INTEGRATION METHODS
Computing Steady State

$$\frac{\hat{x}^l - \hat{x}^{l-1}}{\Delta t} = F(\hat{x}^l) + u(l\Delta t) \quad \text{Solve for } \hat{x}^l \text{ using Newton}$$

$$\alpha_0 = 1 \quad \alpha_1 = -1 \quad \beta_0 = 1$$

- ❖ Time-Integrate Until Steady-State Achieved

$$\frac{dx(t)}{dt} = F(x(t)) + u(t) \Rightarrow \hat{x}^l = \hat{x}^{l-1} + \Delta t [F(\hat{x}^l) + u(l\Delta t)]$$
- ❖ Need many timepoints for lightly damped case!



Solve with Backward-Euler

- ❖ Nonlinear System

$$\frac{dx(t)}{dt} = F \left(\underbrace{x(t)}_{\text{state}} \right) + \underbrace{u(t)}_{\text{input}} \quad \underbrace{x(0) = x_0}_{\text{Initial Condition}}$$

- ❖ Backward Euler Equation for timestep l

$$\hat{x}^l - \hat{x}^{l-1} = \Delta t [F(\hat{x}^l) + u(l\Delta t)]$$
 How do we solve the backward-Euler Equation?

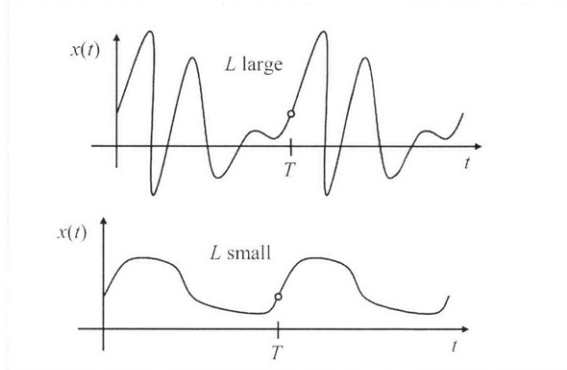
$$\hat{x}^l - \Delta t \left[F(\hat{x}^l) + u(l\Delta t) \right] - \hat{x}^{l-1} = 0$$

unknown
known
known

$$H(\hat{x}^l) = \hat{x}^l - \Delta t [F(\hat{x}^l) + u(l\Delta t)] - \hat{x}^{l-1}$$

Newton Iteration

$$\left[I - \Delta t \frac{\partial F(\hat{x}^{l,k})}{\partial x} \right]_{j(\hat{x}^{l,k})} [\hat{x}^{l,k} - \hat{x}^{l,k-1}] = -H(\hat{x}^{l,k})$$



$$\left. \begin{aligned} \hat{x}^1 &= x(0) + \Delta t [F(x(0)) + u(0)] \\ \hat{x}^2 &= \hat{x}^1 + \Delta t [F(\hat{x}^1) + u(\Delta t)] \\ &\vdots \end{aligned} \right\} F.E.$$

$$\begin{aligned} \hat{x}^l - \hat{x}^{l-1} - \Delta t [F(\hat{x}^l) + u(\Delta t)] &= 0 \\ \Rightarrow \hat{x}^l - \Delta t F(\hat{x}^l) - \underbrace{\hat{x}^{l-1} - \Delta t u(\Delta t)}_{\text{known}} &= 0 \\ \left(\mathbf{I} - \Delta t \frac{\partial F}{\partial x} \right) \hat{x}^{l,j} + 1 - \hat{x}^{l,j} &= - \left(\hat{x}^{l,j} - \Delta t F(\hat{x}^{l,j}) - \hat{x}^{l,j-1} - \Delta t u(\Delta t) \right) \\ \mathbf{J}_H \Delta x &= -\mathbf{H} \end{aligned}$$

Implicit Methods

❖ Backward-Euler Example

Forward - Euler	Backward - Euler
$x(t_1) \approx \hat{x}^1 = x(0) + \Delta t \cdot f(x(0), u(0))$	$x(t_1) \approx \hat{x}^1 = x(0) + \Delta t \cdot f(\hat{x}^1, u(t_1))$
$x(t_2) \approx \hat{x}^2 = \hat{x}^1 + \Delta t \cdot f(\hat{x}^1, u(t_1))$	$x(t_2) \approx \hat{x}^2 = \hat{x}^1 + \Delta t \cdot f(\hat{x}^2, u(t_2))$
\vdots	\vdots

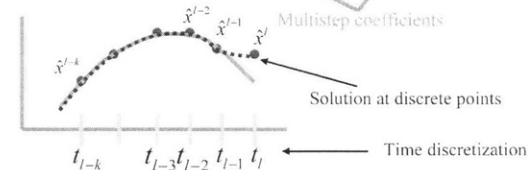
Forward - Euler	Backward - Euler
$x(t_L) \approx \hat{x}^L = \hat{x}^{L-1} + \Delta t \cdot f(\hat{x}^{L-1}, u(t_{L-1}))$	$x(t_L) \approx \hat{x}^L = \hat{x}^{L-1} + \Delta t \cdot f(\hat{x}^L, u(t_L))$
Requires just function evaluations	Nonlinear equation solution at each step

Stepwise Nonlinear equation solution needed whenever $\beta_0 \neq 0$

❖ Multi-Step Methods

Nonlinear Differential Equation: $\frac{d}{dt}x(t) = f(x(t), u(t))$

k - Step Multistep Approach: $\sum_{j=0}^k \alpha_j \hat{x}^{l-j} = \Delta t \sum_{j=0}^k \beta_j f(\hat{x}^{l-j}, u(t_{l-j}))$



❖ Solution with Newton

$$\left[\begin{array}{c} \vdots \\ \text{Jacobian} \end{array} \right] \Delta \hat{x}^{l,j+1} = - \left[\begin{array}{c} \vdots \\ \vdots \end{array} \right]$$

Rewrite the multistep equation

$$\alpha_0 \hat{x}^l - \Delta t \beta_0 f(\hat{x}^l, u(t_l)) + \underbrace{\sum_{i=1}^k \alpha_i \hat{x}^{l-i} - \Delta t \sum_{i=1}^k \beta_i f(\hat{x}^{l-i}, u(t_{l-i}))}_{\text{Independent of } \hat{x}^l} = 0$$

Solve with Newton

$$\left[\alpha_0 I - \Delta t \beta_0 \frac{\partial f(\hat{x}^{l,j}, u(t_l))}{\partial x} \right] (\hat{x}^{l,j+1} - \hat{x}^{l,j}) = - \left[\alpha_0 \hat{x}^{l,j} - \Delta t \beta_0 f(\hat{x}^{l,j}, u(t_l)) + b \right]$$

Here j is the Newton iteration index

Newton Iteration: $\left[\alpha_0 I - \Delta t \beta_0 \frac{\partial f(\hat{x}^{l,j}, u(t_l))}{\partial x} \right] (\hat{x}^{l,j+1} - \hat{x}^{l,j}) = -F(\hat{x}^{l,j})$

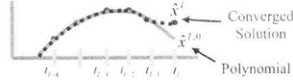
$$\text{Newton } \begin{cases} F(x) = 0 \\ \frac{\partial F}{\partial x}(\tilde{x}^l) \Delta x^{l,j+1} = -F(\tilde{x}^l) \end{cases} \quad \lambda F(x) + (1-\lambda)x(\lambda) = 0$$

$$F(x) = \alpha_0 \tilde{x}^l - \Delta t \beta_0 f(\tilde{x}^l, u(t_l)) + b$$

$$\left[\frac{\partial F}{\partial x} \right] = \alpha_0 \mathbf{I} - \Delta t \beta_0 \left[\frac{\partial f}{\partial x} \right] \quad \left. \vphantom{\left[\frac{\partial F}{\partial x} \right]} \right\} \text{Multistep}$$

$$\frac{dx}{dt} = f(x(t), u(t))$$

Solution with Newton is very efficient



Easy to generate a good initial guess using polynomial fitting
 Jacobian becomes very diagonally dominant for small timesteps

- easy to factor
- or easy to solve iteratively

$$\alpha_0 \mathbf{I} - \Delta t \beta_0 \frac{\partial f(\tilde{x}^l, u(t_l))}{\partial x} \Rightarrow \alpha_0 \mathbf{I} \text{ as } \Delta t \rightarrow 0$$

$$l = 0 \quad x^0 = x(t = 0)$$

$$l \leftarrow l + 1$$

REPEAT (Multistep Integration Loop: time)

$$\tilde{x}^{l,0} \leftarrow \tilde{x}^{l-1} \quad \text{OR F.E.} \quad \tilde{x}^{l,0} \leftarrow \tilde{x}^{l-1} + \Delta t f(\tilde{x}^{l-1}, u(t_l))$$

$$b \leftarrow \alpha_1 \tilde{x}^{l-1} + \dots + \alpha_k \tilde{x}^{l-k} - \Delta t [\beta_1 f(\tilde{x}^{l-1}, u^{l-1}) + \dots + \beta_k f(\tilde{x}^{l-k}, u^{l-k})]$$

REPEAT (Newton Loop)

$$\left[\alpha_0 \mathbf{I} - \Delta t \beta_0 \frac{\partial f}{\partial x}(\tilde{x}^{l,j}, u(t_l)) \right] \Delta \tilde{x}^{l,j+1} = - \left[\alpha_0 \tilde{x}^{l,j} - \Delta t \beta_0 f(\tilde{x}^{l,j}, u(t_l)) + b \right]$$

Solve for $\Delta \tilde{x}^{l,j+1}$ with LU or GCR

$$\tilde{x}^{l,j+1} \leftarrow \tilde{x}^{l,j} + \Delta \tilde{x}^{l,j+1}$$

until $\Delta \tilde{x}^{l,j+1} \parallel \alpha_0 \tilde{x}^{l,j+1} - \Delta t \beta_0 f(\tilde{x}^{l,j+1}) + b \parallel$ small

$$\tilde{x}^l \leftarrow \tilde{x}^{l,j+1}$$

$$l \leftarrow l + 1$$

$$t = t + \Delta t$$

until $t \geq t_{end}$

LU or GCR?

In GCR: at each step need matrix vector product

$$\left[\alpha_0 \mathbf{I} - \Delta t \beta_0 \frac{\partial f}{\partial x}(\tilde{x}^{l,j}, u(t_l)) \right] r^k \cong$$

$$\cong \alpha_0 r^k - \frac{\Delta t \beta_0}{\epsilon} \left[f(\tilde{x}^{l,j} + \epsilon r^k, u(t_l)) - f(\tilde{x}^{l,j}, u(t_l)) \right]$$

GCR

$$M r^k = \alpha_0 r^k - \Delta t \beta_0 \frac{f(x + \epsilon r^k) - f(x)}{\epsilon}$$



Using GCR for

$$\underbrace{\left[\alpha_0 I - \Delta t \beta_0 \frac{\partial f}{\partial x}(\hat{x}^{l,j}, u(t_j)) \right]}_{\text{Jacobian}} \Delta \hat{x}^{l,j+1} = - \left[\dots \right]$$

At each step of GCR need matrix Vector product

$$\underbrace{\left[\alpha_0 I - \Delta t \beta_0 \frac{\partial f}{\partial x}(\hat{x}^{l,j}, u(t_j)) \right]}_{\text{Jacobian}} \vec{r}^k \cong \underbrace{\alpha_0 \vec{r}^k - \Delta t \beta_0}_{\text{M}} \underbrace{f(\hat{x}^{l,j} + \epsilon \vec{r}^k, u(t_j)) - f(\hat{x}^{l,j}, u(t_j))}_{\epsilon}$$

For small ϵ

$$\frac{\partial f}{\partial x}(\hat{x}^{l,j}, u(t_j)) \Delta x \cong f(\hat{x}^{l,j} + \Delta x, u(t_j)) - f(\hat{x}^{l,j}, u(t_j))$$

for small $\Delta x = \epsilon \vec{r}^k$

Total number of time steps $\cong \frac{100 \text{ msec}}{0.1 \text{ nsec}} = 10^9$ steps

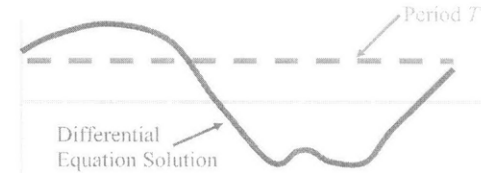
- Each step solve Newton (very fast e.g. 5 iterations)
- Each Newton iteration solve system (e.g. $N = 1000$, sparse $O(N)$ say 10,000 FLOPS)
- Each time step 50,000 FLOPS say 1 GFLOP/sec

$$\frac{50,000 \text{ FLOPS}}{1 \times 10^9 \frac{\text{FLOPS}}{\text{sec}}} = 50 \mu\text{sec per time step}$$

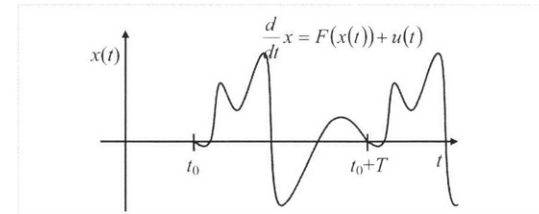
$50 \mu\text{sec} \times 10^9 = 50,000 \text{ sec} \cong 14 \text{ h} \cong \frac{1}{2} \text{ day!!!}$

FINITE-DIFFERENCE METHODS FOR PERIODIC STEADY-STATE Boundary-Value Problem

❖ Basic Formulation



N Differential Equations: $\frac{d}{dt} x_i(t) = F_i(x(t))$
 N Periodicity Constraints: $x_i(T) = x_i(0)$



❖ Finite Difference Methods

○ Linear Example Problem

$$\frac{dx(t)}{dt} = Ax(t) + \underbrace{u(t)}_{\text{input}} \quad t \in [0, T] \quad \underbrace{x(T) = x(0)}_{\text{periodicity constraint}}$$

Discretize e.g. with Backward-Euler

$$\begin{aligned} \hat{x}^1 &= x(0) + \Delta t \cdot (A \hat{x}^1 + u(\Delta t)) \\ \hat{x}^2 &= \hat{x}^1 + \Delta t \cdot (A \hat{x}^2 + u(2\Delta t)) \\ &\vdots \\ \hat{x}^L &= \hat{x}^{L-1} + \Delta t \cdot (A \hat{x}^L + u(L\Delta t)) \end{aligned} \quad \Delta t = \frac{T}{L}$$

Periodicity implies $\hat{x}^L = \hat{x}^0$

$$\begin{matrix} & \xrightarrow{N \times L} & \\ \begin{matrix} \uparrow \\ N \times L \\ \downarrow \end{matrix} & \begin{bmatrix} \mathbf{I} - \Delta t \mathbf{A} & 0 & \cdots & 0 & -\mathbf{I} \\ -\mathbf{I} & \mathbf{I} - \Delta t \mathbf{A} & & & 0 \\ 0 & -\mathbf{I} & \mathbf{I} - \Delta t \mathbf{A} & & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & -\mathbf{I} & \mathbf{I} - \Delta t \mathbf{A} \end{bmatrix} & \begin{bmatrix} \hat{x}^1 \\ \hat{x}^2 \\ \hat{x}^3 \\ \vdots \\ \hat{x}^L \end{bmatrix} = \Delta t \begin{bmatrix} u(\Delta t) \\ u(2\Delta t) \\ u(3\Delta t) \\ \vdots \\ u(L\Delta t) \end{bmatrix} \end{matrix}$$

Matrix is almost lower triangular

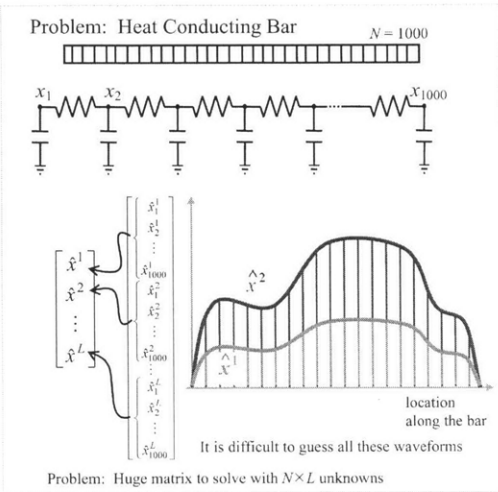
$\hat{x}^0 = \hat{x}^L$
 $\hat{x}^1 = \hat{x}^0 + \Delta t(\mathbf{A}\hat{x}^1 + u(\Delta t))$
 $\Rightarrow \hat{x}^1 = \hat{x}^L + \Delta t(\mathbf{A}\hat{x}^1 + u(\Delta t))$
 $\Rightarrow (\mathbf{I} - \Delta t \mathbf{A})\hat{x}^1 - \hat{x}^L = \Delta t u(\Delta t)$
 $\hat{x}^2 = \hat{x}^1 + \Delta t(\mathbf{A}\hat{x}^2 + u(2\Delta t))$

Backward Euler

$$\begin{bmatrix} \mathbf{I} - \Delta t \mathbf{A} & & & & -\mathbf{I} \\ -\mathbf{I} & \mathbf{I} - \Delta t \mathbf{A} & & & \\ & & \ddots & \ddots & \\ & & & -\mathbf{I} & \mathbf{I} - \Delta t \mathbf{A} \end{bmatrix} \begin{bmatrix} \hat{x}^1 \\ \hat{x}^2 \\ \vdots \\ \hat{x}^L \end{bmatrix} = \Delta t \begin{bmatrix} u(\Delta t) \\ u(2\Delta t) \\ \vdots \\ u(L\Delta t) \end{bmatrix}$$

$\hat{x}^L = \hat{x}^{L-1} + \Delta t(\mathbf{A}\hat{x}^L + u(L\Delta t))$
 $\hat{x}^1 = \hat{x}^L + \Delta t(\mathbf{A}\hat{x}^1 + u(\Delta t))$
 $\hat{x}^2 = \hat{x}^1 + \Delta t(\mathbf{A}\hat{x}^2 + u(2\Delta t))$

Forward Euler

$$\begin{bmatrix} \mathbf{I} & & & & \\ -\mathbf{I} - \Delta t \mathbf{A} & \mathbf{I} & & & \\ & & \ddots & \ddots & \\ & & & -\mathbf{I} - \Delta t \mathbf{A} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \hat{x}^1 \\ \hat{x}^2 \\ \vdots \\ \hat{x}^L \end{bmatrix} = \Delta t \begin{bmatrix} u(\Delta t) \\ u(2\Delta t) \\ \vdots \\ u(L\Delta t) \end{bmatrix}$$


o Nonlinear Problem

$$\frac{dx(t)}{dt} = F(x(t)) + \underbrace{u(t)}_{\text{input}} \quad t \in [0, T] \quad \underbrace{x(T) = x(0)}_{\text{periodicity constraint}}$$

Discretize e.g. with Backward-Euler

$$H_{FD} \begin{bmatrix} \hat{x}^1 \\ \hat{x}^2 \\ \vdots \\ \hat{x}^L \end{bmatrix} = \begin{bmatrix} \hat{x}^1 - \hat{x}^L - \Delta t [F(\hat{x}^1) + u(\Delta t)] \\ \hat{x}^2 - \hat{x}^1 - \Delta t [F(\hat{x}^2) + u(2\Delta t)] \\ \vdots \\ \hat{x}^L - \hat{x}^{L-1} - \Delta t [F(\hat{x}^L) + u(L\Delta t)] \end{bmatrix} = 0$$

Solve the huge system using Newton's Method

$x(t)$
 $x(t_1) = \phi(x_{t_0}, t_0, t_1)$
 x_{t_0}
 t_0
 t_1
 t

Your ODE Integrator
 Your Friend's ODE Integrator
 A Commercial ODE Integrator
 $\phi(\quad)$

x_{t_0}, t_0, t_1 → $x(t_1)$

$\phi: (x_{t_0}, t_0, t_1) \rightarrow x(t_1)$ The state transition function is a good model for

$\phi(x_{t_0}, t_0, t_1) \rightarrow x(t_1) = x_{t_0} e^{\lambda(t_1 - t_0)}$

$\frac{d}{dt} x(t) = \lambda x(t)$
 $x(t) = \mathbf{A} e^{\lambda t}$
 $x_{t_0} = \mathbf{A} e^{\lambda t_0} \rightarrow \mathbf{A} = x_{t_0} e^{-\lambda t_0}$
 $x(t) = x_{t_0} e^{\lambda(t - t_0)}$ $x(t_1) = x_{t_0} e^{\lambda(t_1 - t_0)}$

SHOOTING METHODS
State Transition Function

❖ Basic Definitions

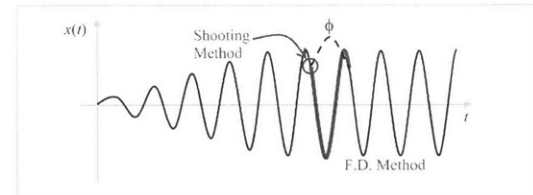
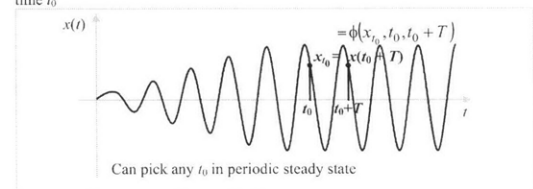
Start with $\frac{dx(t)}{dt} = F(x(t)) + u(t)$

And assume $x(t)$ is unique given $x(0)$.

The O.D.E. defines a State-Transition Function

$\Phi(x_{t_0}, t_0, t_1) = x(t_1)$

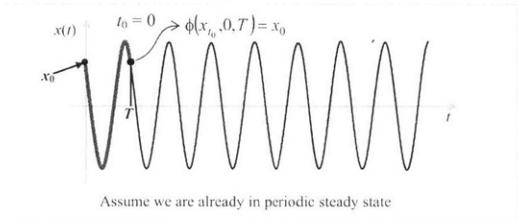
Where $x(t_1)$ is the O.D.E. solution at time t_1 when starting from state $x_{t_0} = x(t_0)$ at time t_0



❖ State Transition Function Example

$\frac{dx(t)}{dt} = \lambda x(t)$

$\Phi(x_{t_0}, t_0, t_1) = e^{\lambda(t_1 - t_0)} x_{t_0}$



❖ Abstract Formulation
 Solve $\overbrace{x(t_0 + T)} = x_{t_0}$ for a t_0 of your choice
 $\Phi(x_0, t_0, t_0 + T) = x_0$
 e.g. pick $t_0 = 0$ $H(x_0) = \Phi(x_0, 0, T) - x_0 = 0$

Evaluate $H(x^k)$
 Integrate $\frac{d}{dt}x(t) = F(x(t)) + u(t)$
 $x(0) = x^k$
 Over the interval $[0, T]$

x_0 unknown
 $H(x_0) = 0$
 $\frac{\Phi(x_0, 0, T) - x_0}{H(x_0)} = 0$

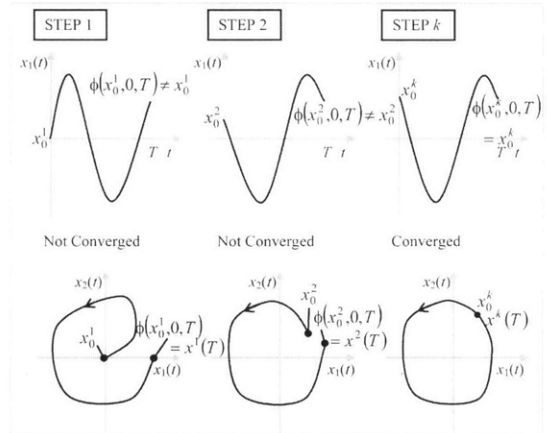
REPEAT
 $\left[\frac{\partial \Phi}{\partial x_0} \Big|_{x_0^k} - \mathbf{I} \right] \Delta x_0^{k+1} = -[\Phi(x_0^k, 0, T) - x_0^k]$
 UNTIL $|\Delta x_0^{k+1}|, |H(x_0^k)|$ small

Solve $H(x_0) = 0$ for x_0 using Newton's method
 $J_H(x_0^k) (x_0^{k+1} - x_0^k) = -H(x_0^k)$

where $J_H(x_0^k) = \frac{\partial \Phi(x_0^k, 0, T)}{\partial x_0} - \mathbf{I}$

$k \leftarrow 0$
 $x^k \leftarrow x_0$ Note this is not so different to guess anymore:
 it is just one distribution of temperatures.

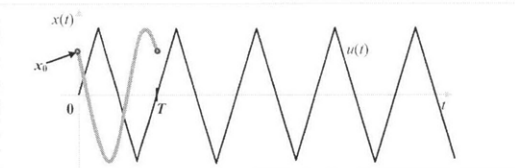
REPEAT
 $\left[\frac{\partial \Phi}{\partial x_0} \Big|_{x_0^k} - \mathbf{I} \right] \Delta x_0^{k+1} = -[\Phi(x_0^k, 0, T) - x_0^k]$
 $x_0^{k+1} \leftarrow x_0^k + \Delta x_0^{k+1}$
 $k \leftarrow k + 1$
 UNTIL $|\Delta x_0^k|, |\Phi(x_0^{k-1}, 0, T) - x_0^{k-1}|$ small



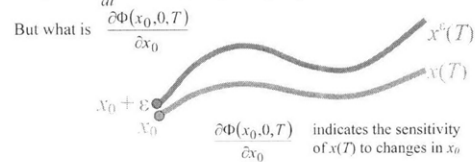
$$\begin{bmatrix} \frac{\partial F}{\partial x} \end{bmatrix}_{x_0} \bar{\varepsilon} \cong F(x_0 + \bar{\varepsilon}) - F(x_0) \quad \Delta x \text{ small}$$

$$\begin{bmatrix} \frac{\partial \Phi}{\partial x_0} \end{bmatrix}_{(x_0, 0, T)} \bar{\varepsilon} \cong \underbrace{\Phi(x_0 + \bar{\varepsilon}, 0, T)}_{x^e(T)} - \underbrace{\Phi(x_0, 0, T)}_{x(T)}$$

Matrix



❖ Computing Newton
 In order to calculate $H(x_0) = \Phi(x_0, 0, T) - x_0$
 Integrate $\frac{dx(t)}{dt} = F(x(t)) + u(t)$ on $[0, T]$



$$\begin{bmatrix} \frac{\partial \Phi}{\partial x} \end{bmatrix} \begin{bmatrix} x_1^e(0) - x_1(0) \\ \vdots \\ x_N^e(0) - x_N(0) \end{bmatrix} \cong \begin{bmatrix} x_1^e(T) - x_1(T) \\ \vdots \\ x_N^e(T) - x_N(T) \end{bmatrix}$$

If we choose $\bar{\varepsilon} = \varepsilon_1 \hat{e}_1 = \begin{bmatrix} \varepsilon_1 \\ 0 \\ 0 \\ \vdots \end{bmatrix}$ we get the first column of $\frac{\partial \Phi}{\partial x}(x_0, 0, T)$

First column of $\frac{\partial \Phi}{\partial x} \cong \begin{bmatrix} x_1^{\varepsilon_1}(T) - x_1(T) \\ \varepsilon_1 \\ x_2^{\varepsilon_1}(T) - x_2(T) \\ \vdots \\ x_N^{\varepsilon_1}(T) - x_N(T) \end{bmatrix}$

Difference is 2^{nd} column of $\frac{\partial \Phi}{\partial x}$
 Difference is 1^{st} column of $\frac{\partial \Phi}{\partial x}$

This is not efficient we need to do a lot of integrations to get the entire matrix (e.g. $M = 1000$ columns)

$$\left[\frac{\partial \phi}{\partial x_0}(x_0, 0, T) \right] \bar{\varepsilon} \cong \phi(x_0 + \bar{\varepsilon}, 0, T) - \phi(x_0, 0, T) = x^{\bar{\varepsilon}}(T) - x(T)$$

First Column use $\bar{\varepsilon} = \varepsilon_1 \hat{e}_1$

$$\left[\frac{\partial \phi}{\partial x_0}(x_0, 0, T) \right]_1 \varepsilon_1 \hat{e}_1 \cong \phi(x_0 + \varepsilon_1 \hat{e}_1, 0, T) - \phi(x_0, 0, T) = x^{\varepsilon_1}(T) - x(T)$$

Second Column use $\bar{\varepsilon} = \varepsilon_2 \hat{e}_2$

$$\left[\frac{\partial \phi}{\partial x_0}(x_0, 0, T) \right]_2 \varepsilon_2 \hat{e}_2 \cong \phi(x_0 + \varepsilon_2 \hat{e}_2, 0, T) - \phi(x_0, 0, T) = x^{\varepsilon_2}(T) - x(T)$$

Not Efficient
Need N simulations
of one period (one
for each column)

Sensitivity Matrix

❖ Perturbation

$$\frac{\partial \Phi(x, 0, T)}{\partial x} \approx \begin{bmatrix} x_1^{\varepsilon_1}(T) - x_1(T) & \dots & x_1^{\varepsilon_N}(T) - x_1(T) \\ \varepsilon_1 & \dots & \varepsilon_N \\ \vdots & \ddots & \vdots \\ x_N^{\varepsilon_1}(T) - x_N(T) & \dots & x_N^{\varepsilon_N}(T) - x_N(T) \\ \varepsilon_1 & \dots & \varepsilon_N \end{bmatrix}$$

Is $\frac{\partial \phi}{\partial x}(x_0, 0, T)$ DENSE or SPARSE?

Try $x_0 = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$ Integrate over one period get $\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$
 $\phi(x_0, 0, T) = x(T)$

Try perturbing initial condition $x_0 + \varepsilon_1 \hat{e}_1 = \begin{bmatrix} \varepsilon_1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$

Integrate over one period get $\begin{bmatrix} * \\ * \\ \vdots \\ * \end{bmatrix}$ all nonzero
 $\phi(x_0 + \varepsilon_1 \hat{e}_1, 0, T) = x^{\varepsilon_1}(T)$

So the matrix $\frac{\partial \phi}{\partial x}$ is DENSE.

For $i = 1$ to N

$x_0^{\varepsilon_i} = x_0 + \varepsilon_i \hat{e}_i$ Integrate $x = F(x)$
 Compute $x(T) = \Phi(x_0^{\varepsilon_i}, 0, T)$ $x(0) = x_0^{\varepsilon_i}$ on $[0, T]$
 Column of $\frac{\partial \Phi}{\partial x}$
 N Columns $\rightarrow N$ Integrations

Is the sensitivity matrix dense or sparse?

$$\frac{d}{dt} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_N \end{bmatrix} = \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & & -1 & 2 \\ & & & & -1 & 2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_N \end{bmatrix}$$

In general the sensitivity matrix will be dense

$\frac{\partial \Phi}{\partial x}$ by perturbation
For $i = 1$ to N
Perturb $x(0)$
Compute $x^{(i)}(T)$ \rightarrow $\frac{\partial \Phi}{\partial x}$

❖ Efficient Sensitivity Evaluation

Differentiate the first step of Backward-Euler

$$\frac{\partial}{\partial x(0)} (\hat{x}^1 - x(0) - \Delta t (F(\hat{x}^1) + u(\Delta t))) = 0$$

$$\Rightarrow \frac{\partial \hat{x}^1}{\partial x(0)} - \frac{\partial x(0)}{\partial x(0)} - \Delta t \frac{\partial F(\hat{x}^1)}{\partial x} \frac{\partial \hat{x}^1}{\partial x(0)} = 0$$

$$\Rightarrow \left(I - \Delta t \frac{\partial F(\hat{x}^1)}{\partial x} \right) \frac{\partial \hat{x}^1}{\partial x(0)} = \frac{\partial x(0)}{\partial x(0)}$$

Applying the same trick on the l^{th} step

$$\Rightarrow \left(I - \Delta t \frac{\partial F(\hat{x}^l)}{\partial x} \right) \frac{\partial \hat{x}^l}{\partial x(0)} = \frac{\partial \hat{x}^{l-1}}{\partial x(0)}$$

$$\frac{\partial \Phi(x, 0, T)}{\partial x} \approx \prod_{l=1}^L \left(I - \Delta t \frac{\partial F(\hat{x}^l)}{\partial x} \right)^{-1}$$

B.E. $\hat{x}^1 - x(0) - \Delta t (F(\hat{x}^1) + u(\Delta t)) = 0$

Apply Newton $\left(I - \Delta t \frac{\partial F(\hat{x}^{1,k})}{\partial x} \right) \Delta \hat{x}^{1,k+1} = \dots$

$$\frac{d}{dt} x(t) = Ax(t) \Rightarrow x(t) = e^{At} x(0)$$

$$\frac{d\Phi}{dx(0)} \cong (I - \Delta t A)^{-l}$$

$$\frac{d\Phi}{dx(0)} = e^{A^T T} = e^{A^T (l\Delta t)}$$

$$\begin{aligned} \frac{\partial \Phi}{\partial x} = \frac{\partial \hat{x}^l}{\partial x_0} &= \left[I - \Delta t \frac{\partial F}{\partial x}(\hat{x}^l) \right]^{-1} \frac{\partial \hat{x}^{l-1}}{\partial x(0)} \\ &= \left[I - \Delta t \frac{\partial F}{\partial x}(\hat{x}^l) \right]^{-1} \left[I - \Delta t \frac{\partial F}{\partial x}(\hat{x}^{l-1}) \right]^{-1} \frac{\partial \hat{x}^{l-2}}{\partial x(0)} \\ &\quad \vdots \\ &= \underbrace{\left[I - \Delta t \frac{\partial F}{\partial x}(\hat{x}^l) \right]^{-1} \left[I - \Delta t \frac{\partial F}{\partial x}(\hat{x}^{l-1}) \right]^{-1} \dots \left[I - \Delta t \frac{\partial F}{\partial x}(\hat{x}^1) \right]^{-1}}_I I \end{aligned}$$

In practice we solve l systems (each with size N)

❖ Observations on Sensitivity Matrix

Newton at each timestep uses same matrices

$$\frac{\partial \Phi(x, 0, T)}{\partial x} \approx \prod_{l=1}^L \underbrace{\left(I - \Delta t \frac{\partial F(\hat{x}^l)}{\partial x} \right)^{-1}}_{\text{Timestep Newton Jacobian}}$$

Formula simplifies in the linear case

$$\frac{\partial \Phi(x, 0, T)}{\partial x} \approx (I - \Delta t A)^{-L}$$

$$\frac{dx}{dt} = Ax(t) + u(t) \quad F(x(t)) = Ax(t)$$

$$x(t) = e^{At} x(0) \quad T = \Delta t \cdot l$$

$$\left[\frac{\partial \Phi}{\partial x} \right] \Delta x(0) \cong \Phi(x(0) + \Delta x(0), 0, T) - \Phi(x(0), 0, T)$$

$$\cong e^{A\Delta t} [x(0) + \Delta x] - e^{A\Delta t} x(0)$$

$$= [e^{A\Delta t}] \Delta x(0)$$

$$= [e^{-A\Delta t}]^T \Delta x(0)$$

$$\left[\frac{\partial \Phi}{\partial x} \right] \cong [\mathbf{I} - \mathbf{A}\Delta t]^{-l} \quad \text{First Order Expansion (B.E.)}$$

Matrix-Free Approach

❖ Basic Setup

REPEAT

$$\left[\begin{array}{c} \frac{\partial \Phi}{\partial x_0} \\ \frac{\partial \Phi}{\partial x_0} \end{array} - \mathbf{I} \right] \Delta x_0^{k+1} = - \left[\frac{\Phi(x_0^k, 0, T) - x_0^k}{h} \right]$$

$$x_0^{k+1} \leftarrow x_0^k + \Delta x_0^{k+1}$$

$$k \leftarrow k + 1$$

UNTIL $\|\Delta x_0^k\|, \|\Phi(x_0^{k-1}, 0, T) - x_0^{k-1}\|$ small

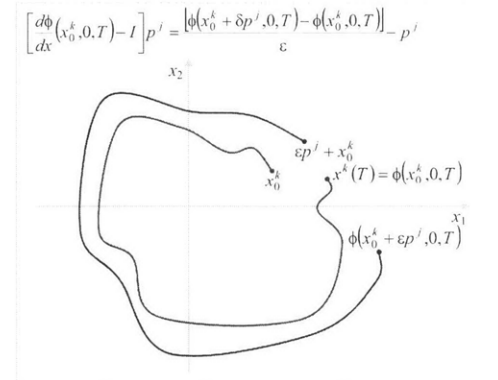
Start with $\frac{dx(t)}{dt} = F(x(t)) + u(t)$

$$H(x_0) = \Phi(x_0, 0, T) - x_0 = 0$$

Solve $H(x_0) = 0$ for x_0 using Newton's method

$$J_H(x_0^k) (x_0^{k+1} - x_0^k) = -H(x_0^k)$$

$$\text{where } J_H(x_0^k) = \frac{\partial \Phi(x_0^k, 0, T)}{\partial x_0} - \mathbf{I}$$



❖ Matrix-Vector Product

Solve Newton equation with Krylov-subspace method

$$\left[\frac{\partial \Phi(x_0^k, 0, T)}{\partial x_0} - \mathbf{I} \right] (x_0^{k+1} - x_0^k) = x_0^k - \Phi(x_0^k, 0, T)$$

Matrix-Vector Product Computation

$$\left[\frac{\partial \Phi(x_0^k, 0, T)}{\partial x_0} - \mathbf{I} \right] p^j \cong \frac{\Phi(x_0^k + \epsilon p^j, 0, T) - \Phi(x_0^k, 0, T)}{\epsilon} - p^j$$

Krylov method search direction

$$Ap = \left[\frac{\partial \Phi}{\partial x} - \mathbf{I} \right] p = \frac{\partial \Phi}{\partial x} p - p$$

Newton Iteration

$$\left[\frac{\partial \phi}{\partial x} - I \right] [\Delta x^{j+1}(0)] = - \underbrace{H(x^j(0))}_{\substack{\text{call} \\ \text{integration}}}$$

Use Krylov

$$\left[\frac{\partial \phi}{\partial x} - I \right] p \cong \frac{\phi(x^j(0) + \varepsilon p) - \phi(x^j(0))}{\varepsilon} - p$$

Matrix-Vector Product

- 1.) Computing P.S.S. $x(T) = x(0)$
- 2.) Use Shooting $\phi(x(0), 0, T) - x(0) = 0$
- 3.) Use Newton $\left(\frac{\partial \phi}{\partial x} - I \right) \Delta x^{k+1}(0) = -(\phi(x^k, 0, T) - x^k)$
- 4.) Use Krylov to Solve Newton $\left(\frac{\partial \phi}{\partial x} - I \right) p$

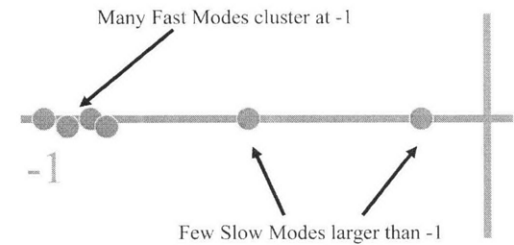
search direction

❖ Convergence for GCR

Example
 $\frac{dx}{dt} - Ax = 0 \quad \text{eig}(A) \text{ real and negative}$
 Shooting-Newton Jacobian
 $\frac{\partial \phi(x, 0, T)}{\partial x} - I = e^{AT} - I$

$\text{eig}(e^{AT} - I)$ versus $\text{eig}(A)$

$e^{\lambda T} - 1$
 \swarrow
 $\text{eig}(A)$

$$e^{AT} - I = S \begin{bmatrix} e^{\lambda_1 T} - 1 & & \\ & \ddots & \\ & & e^{\lambda_n T} - 1 \end{bmatrix} S^{-1}$$


Use Spectral Mapping Theorem

$\text{eig}(e^{AT} - I)$

if λ_i is eigenvalue of $A \rightarrow e^{\lambda_i T} - 1$ is eigenvalue of $e^{AT} - I$

Fast Modes

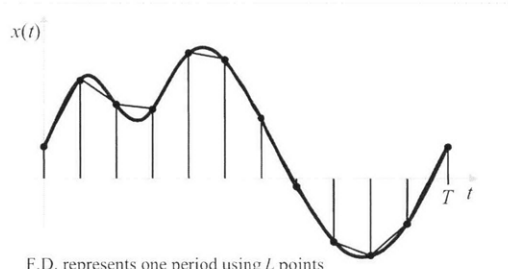
Slow Modes

$e^{\lambda T} - 1$

a few slow λ modes

all fast modes are in one cluster

Typically there are many fast modes (each strut) and few slow modes (the whole frame)

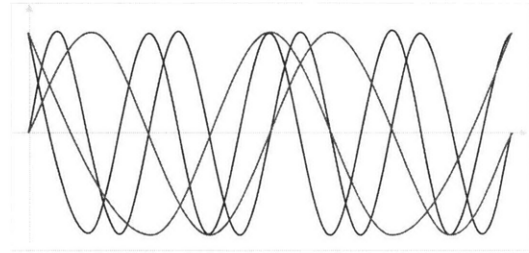
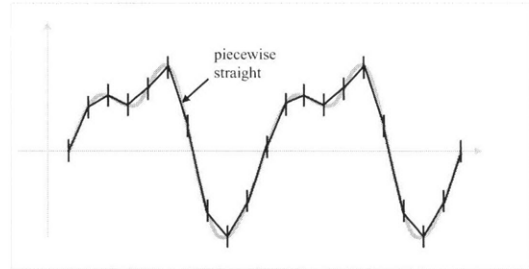


F.D. represents one period using L points
 Maybe we can use other ways to represent the shape
 Periodic \rightarrow use Fourier Coefficients

$$x(t) = \dots + x_{-2}e^{-j2\pi t/T^2} + x_{-1}e^{-j2\pi t/T} + x_0 + x_1e^{j2\pi t/T} + x_2e^{j2\pi t/T^2} + \dots$$

\swarrow frequency $\frac{1}{T}$ fundamental
 \swarrow frequency $\frac{2}{T}$ fundamental \searrow 2nd Harmonic etc...

Works well (need few coefficients) if function smooth harmonics
 Let's solve for x_j



SPECTRAL METHODS (HARMONIC BALANCE)

Fourier Representation

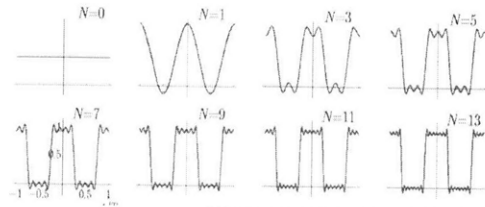
- ❖ Truncation Approximation
 - Periodic Function \rightarrow Fourier Series

$$x(t) = \sum_{l=-\infty}^{\infty} X_l e^{-j2\pi l t/T}$$
 - Approximate a function with truncated series

$$x(t) \approx \sum_{l=-L}^L X_l e^{-j2\pi l t/T}$$

❖ Square Wave Example

$$x_N(t) = \frac{1}{2} + \sum_{n=1}^N \left(\frac{\sin(n\pi/2)}{n\pi/2} \right) \cos(2\pi n t/T)$$



©1997 by Alan V. Oppenheim & Alan S. Willsky

❖ Annoyance for Real Functions

- Real $x \rightarrow$ Fourier Coefficients complex conjugate
- Can rewrite series with fewer unknowns

$$x(t) = \sum_{l=1}^{\infty} \underbrace{X_l e^{-j2\pi l t/T} + X_l^* e^{+j2\pi l t/T}}_{\text{Real}} + \underbrace{X_0}_{l=0}$$

$$\langle f(t), g(t) \rangle_T = \int_0^T g(t) f(t) dt$$

$$l = m \int_0^T e^{j2\pi m t/T} e^{-j2\pi m t/T} dt = \int_0^T dt = T = \left\| e^{j2\pi m t/T} \right\|^2$$

$$\|f(t)\|^2 = \langle f(t), g(t) \rangle_T$$

$$x(t) = \dots + X_{-1} e^{-j2\pi t/T} + X_0 + X_1 e^{j2\pi t/T} + \dots$$

$\{\hat{Q}_1, \hat{Q}_2, \dots, \hat{Q}_n\}$ orthonormal vectors

$$\bar{b} = y_1 \hat{Q}_1 + y_2 \hat{Q}_2 + \dots + y_n \hat{Q}_n \quad y_i = \langle \bar{b}, \hat{Q}_i \rangle$$

$y_1 = \langle \bar{b}, \hat{Q}_1 \rangle$ $y_2 = \langle \bar{b}, \hat{Q}_2 \rangle$ Orthogonal Projection

$$\hat{Q}_1 = \frac{\bar{v}_1}{\|\bar{v}_1\|}$$

$\{\bar{v}_1, \bar{v}_2, \dots, \bar{v}_n\}$ orthogonal vectors $\alpha_i = \frac{\langle \bar{b}, \bar{v}_i \rangle}{\|\bar{v}_i\|^2}$

$$\bar{b} = \alpha_1 \bar{v}_1 + \alpha_2 \bar{v}_2 + \dots + \alpha_n \bar{v}_n$$

$\alpha_1 = \frac{\langle \bar{b}, \bar{v}_1 \rangle}{\|\bar{v}_1\|^2}$ $\alpha_2 = \frac{\langle \bar{b}, \bar{v}_2 \rangle}{\|\bar{v}_2\|^2}$ Orthogonal Projection

Compare

$$x(t) = \dots + x_{-2} e^{-j2\pi t/T} + x_{-1} e^{-j\pi t/T} + x_0 + x_1 e^{j\pi t/T} + x_2 e^{j2\pi t/T} + \dots$$

$$\bar{b} = \dots + \alpha_{-2} \bar{v}_{-2} + \alpha_{-1} \bar{v}_{-1} + \alpha_0 + \alpha_1 \bar{v}_1 + \alpha_2 \bar{v}_2 + \dots$$

$b \leftrightarrow x(t)$

$$\bar{v}_m \leftrightarrow e^{jm\pi t/T}$$

$$\alpha_m \leftrightarrow x_m$$

$$x(t) = X_1 e^{-j2\pi t/T} + X_2 e^{-j2\pi t/T} + \dots$$

$$X_1 = \frac{1}{T} \int_0^T e^{j2\pi t/T} x(t) dt \quad X_2 = \frac{1}{T} \int_0^T e^{j2\pi t/T} x(t) dt$$

$$x_m = \frac{1}{T} \int_0^T e^{-j2\pi m t/T} x(t) dt$$

- ❖ Orthogonality
 - Terms in Fourier Series are orthogonal

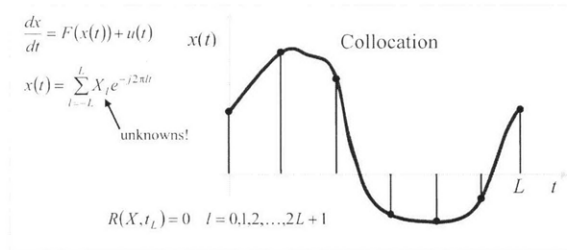
$$\int_0^T e^{i2\pi l t/T} e^{-i2\pi m t/T} dt = 0 \quad l \neq m$$
 - Simple formula for computing coefficients

$$\int_0^T e^{-i2\pi m t/T} x(t) dt = \int_0^T e^{i2\pi m t/T} \sum_{l=-\infty}^{\infty} X_l e^{-i2\pi l t/T} dt = T X_m$$
- ❖ Advantages
 - For smooth functions (infinitely cont. diff)

$$\lim_{m \rightarrow \infty} \frac{1}{T} \int_0^T e^{-i2\pi m t/T} x(t) dt = \lim_{m \rightarrow \infty} X_m = O(\epsilon^m)$$
 - Automatically satisfies periodicity

$$x(t+T) = \sum_{l=-L}^L X_l e^{-i2\pi l (t+T)/T} = \sum_{l=-L}^L X_l e^{-i2\pi l t/T} = x(t)$$

Computing Coefficients



- ❖ Residual
 - Plug representation in differential equation

$$\frac{R(\tilde{x}, t)}{\text{Residual}} = \frac{d}{dt} \left(\sum_{l=-L}^L X_l e^{-i2\pi l t/T} \right) - F \left(\sum_{l=-L}^L X_l e^{-i2\pi l t/T} \right) - u(t)$$

$$\frac{d}{dt} \underbrace{x(t)}_{\sum_{l=-L}^L X_l e^{-i2\pi l t/T}} = f(x(t)) + u(t)$$

- Simplify by differentiating representation

$$\frac{R(\tilde{x}, t)}{\text{Residual}} = \sum_{l=-L}^L -\frac{i2\pi l}{T} X_l e^{-i2\pi l t/T} - F \left(\sum_{l=-L}^L X_l e^{-i2\pi l t/T} \right) - u(t)$$

- ❖ Collocation & Galerkin
 - Collocation: residual = 0 at test points

$$\frac{R(\tilde{x}, t)}{\text{Residual}} = 0 \quad l = \{1, \dots, 2L+1\}$$

$$x(t) = \sum_{l=-L}^L X_l e^{i2\pi l t/T} \quad \text{complex}$$

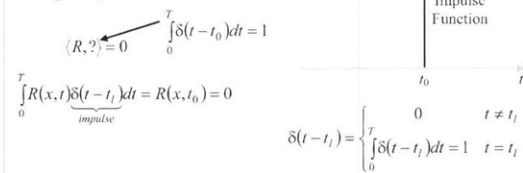
- Galerkin: residual orthogonal to Fourier Terms

$$\int_0^T e^{i2\pi m t/T} \frac{R(\tilde{x}, t)}{\text{Residual}} dt = 0 \quad m = \{-L, \dots, 0, \dots, L\}$$

$R(x, t) = 0$

$\langle R, \psi \rangle = 0 \quad \text{Galerkin } \langle R, v_m \rangle = 0 \quad m = \dots$

Actually also collocation can be written as an orthogonality condition with different vectors



Galerkin Methods

❖ Computing Coefficients

- Linear Galerkin $F(x) = Ax$

$$-\left\{ \int_0^T e^{i2\pi m \frac{t}{T}} \left[\sum_{l=-L}^L -\frac{i2\pi l}{T} X_l e^{-i2\pi l \frac{t}{T}} - F \left(\sum_{l=-L}^L X_l e^{-i2\pi l \frac{t}{T}} \right) - u(t) \right] dt \right\} =$$

$$i2\pi m X_m + \int_0^T e^{i2\pi m \frac{t}{T}} F \left(\sum_{l=-L}^L X_l e^{-i2\pi l \frac{t}{T}} \right) dt + \int_0^T e^{i2\pi m \frac{t}{T}} u(t) dt = 0$$

$$m = \{-L, \dots, 0, \dots, L\}$$

Solve the system using Newton

$$i2\pi m X_m + \int_0^T e^{i2\pi m \frac{t}{T}} A \left(\sum_{l=-L}^L X_l e^{-i2\pi l \frac{t}{T}} \right) dt + \int_0^T e^{i2\pi m \frac{t}{T}} u(t) dt = 0$$

$$\begin{bmatrix} i2\pi L + A & & & \\ & i2\pi(L-1) + A & & \\ & & \ddots & \\ & & & -i2\pi L + A \end{bmatrix} \begin{bmatrix} X_{-L} \\ X_{-(L-1)} \\ \vdots \\ X_L \end{bmatrix} = \begin{bmatrix} U_{-L} \\ U_{-(L-1)} \\ \vdots \\ U_L \end{bmatrix}$$

$$\begin{aligned} j2\pi(-L)x_L + Ax_{-L}T + TU_{-L} &= 0 \\ j2\pi(-(L-1))x_{L-1} + Ax_{-L-1}T + TU_{-(L-1)} &= 0 \\ \frac{j2\pi(-(L-1))}{T}x_{L-1} + Ax_{-L-1} + U_{-(L-1)} &= 0 \end{aligned}$$

Collocation Methods

❖ Computing Coefficients

- Collocation – residual zero at test times

$$\text{Residual } R(\vec{X}, t) = 0 = \sum_{l=-L}^L -\frac{i2\pi l}{T} X_l e^{-i2\pi l \frac{t}{T}} - F \left(\sum_{l=-L}^L X_l e^{-i2\pi l \frac{t}{T}} \right) - u(t_l)$$

$$l = \{1, \dots, 2L+1\}$$

What is this?
DFT Matrix
(Converts to time)

$$\begin{bmatrix} \frac{dx}{dt}(t_1) \\ \frac{dx}{dt}(t_2) \\ \vdots \\ \frac{dx}{dt}(t_{2L+1}) \end{bmatrix}$$

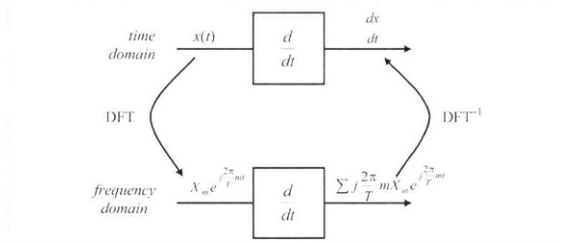
❖ Discrete Fourier Transform

$$\begin{bmatrix} e^{i2\pi \frac{L}{T} t_1} & \dots & e^{-i2\pi \frac{L}{T} t_1} \\ \vdots & \ddots & \vdots \\ e^{i2\pi \frac{L}{T} t_{(2L+1)}} & \dots & e^{-i2\pi \frac{L}{T} t_{(2L+1)}} \end{bmatrix} \begin{bmatrix} X_{-L} \\ X_{-(L-1)} \\ \vdots \\ X_L \end{bmatrix} = \begin{bmatrix} x(t_1) \\ x(t_2) \\ \vdots \\ x(t_{(2L+1)}) \end{bmatrix}$$

Discrete Fourier Transform (DFT)

$$x(t) \approx \sum_{l=-L}^L X_l e^{-i2\pi l \frac{t}{T}}$$

If $t_l = \frac{1}{2L+1}T$ then DFTMatrix has orthogonal columns



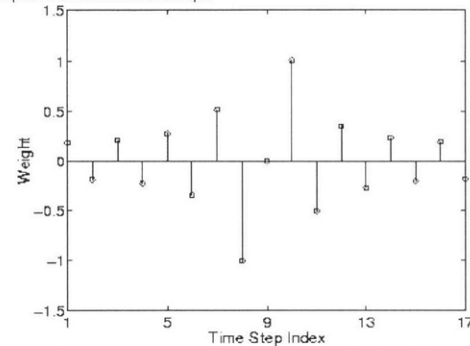
❖ Collocation using Timepoints

$$DFT \begin{bmatrix} i2\pi L \\ T \\ i2\pi(L-1) \\ T \\ \vdots \\ -i2\pi L \\ T \end{bmatrix} (DFT)^{-1} \begin{bmatrix} x(t_1) \\ x(t_2) \\ \vdots \\ x(t_{(2L+1)}) \end{bmatrix} = \begin{bmatrix} F(x(t_1)) \\ F(x(t_2)) \\ \vdots \\ F(x(t_{(2L+1)})) \end{bmatrix} = \begin{bmatrix} u(t_1) \\ u(t_2) \\ \vdots \\ u(t_{(2L+1)}) \end{bmatrix}$$

Spectral Differentiation
 Converting timepoint into Fourier Coeffs,
 Differentiating, and then returning to time

$$\begin{bmatrix} \frac{dx}{dt}(t_1) \\ \frac{dx}{dt}(t_2) \\ \vdots \\ \frac{dx}{dt}(t_{2L+1}) \end{bmatrix} = DFT \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} (DFT)^{-1} \begin{bmatrix} x(t_1) \\ x(t_2) \\ \vdots \\ x(t_{2L+1}) \end{bmatrix}$$

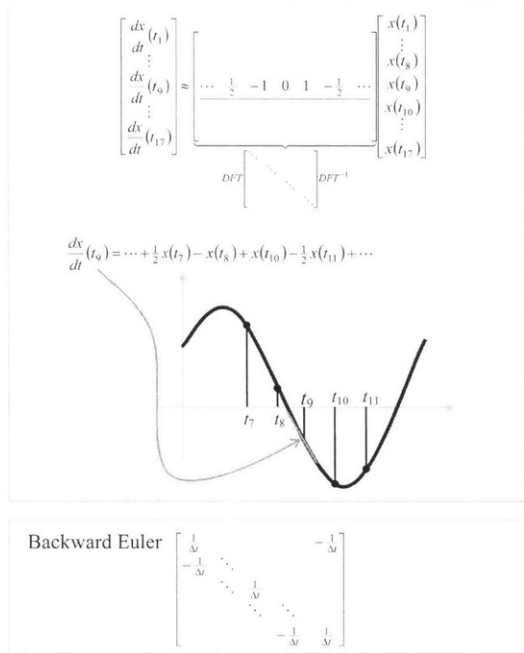
❖ Spectral Differentiation Example



Coeff of row 9 of the spectral differentiation matrix, $T = 17$ and $2L+1=17$

Backward Euler

$$\begin{bmatrix} \frac{dx}{dt}(t_1) \\ \vdots \\ \frac{dx}{dt}(t_9) \\ \vdots \\ \frac{dx}{dt}(t_{17}) \end{bmatrix} \approx \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ \dots & -\frac{1}{\Delta t} & \frac{1}{\Delta t} & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} x(t_1) \\ \vdots \\ x(t_9) \\ \vdots \\ x(t_{17}) \end{bmatrix} \quad \frac{dx}{dt}(t_9) = \frac{-x(t_8) + x(t_9)}{\Delta t}$$



❖ Spectral Collocation versus F-D
Spectral Differentiation:

$$DFT \underbrace{\begin{bmatrix} \frac{i2\pi L}{T} & & & \\ & i2\pi(L-1) & & \\ & & \ddots & \\ & & & -\frac{i2\pi L}{T} \end{bmatrix}}_{\text{Spectral Differentiation}} (DFT)^{-1} \begin{bmatrix} x(t_1) \\ x(t_2) \\ \vdots \\ x(t_{(2L+1)}) \end{bmatrix} = \begin{bmatrix} F(x(t_1)) \\ F(x(t_2)) \\ \vdots \\ F(x(t_{(2L+1)})) \end{bmatrix} = \begin{bmatrix} u(t_1) \\ u(t_2) \\ \vdots \\ u(t_{(2L+1)}) \end{bmatrix}$$

Backward Euler Differentiation:

$$\begin{bmatrix} \frac{1}{\Delta t} & & & -\frac{1}{\Delta t} \\ -\frac{1}{\Delta t} & \frac{1}{\Delta t} & & \\ & -\frac{1}{\Delta t} & \frac{1}{\Delta t} & \\ & & \ddots & \\ & & & -\frac{1}{\Delta t} & \frac{1}{\Delta t} \end{bmatrix} \begin{bmatrix} \hat{x}^1 \\ \hat{x}^2 \\ \vdots \\ \hat{x}^{2L+1} \end{bmatrix} + \begin{bmatrix} F(x(t_1)) \\ F(x(t_2)) \\ \vdots \\ F(x(t_{(2L+1)})) \end{bmatrix} = \begin{bmatrix} u(t_1) \\ u(t_2) \\ \vdots \\ u(t_{(2L+1)}) \end{bmatrix}$$

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 18.

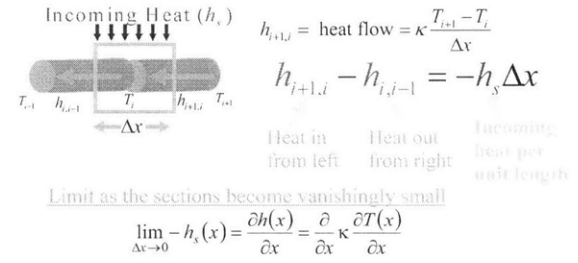
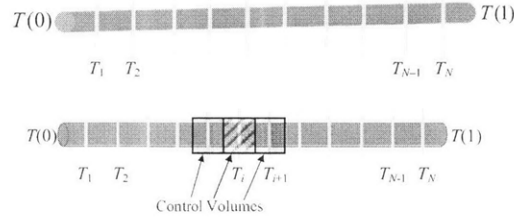
Finite Difference Methods for Boundary Value Problems

TODAY'S OUTLINE:

- ❖ Informal Finite Difference Methods
 - Heat Conducting Bar
- ❖ Comparing FEM and FD in 1-D
- ❖ Convergence Analysis for Heat Equation
 - Local Truncation Error and Consistency
 - Global Error and Stability
- ❖ Formal Generalization to any PDE
 - Consistency + Stability yields Convergence

INFORMAL FINITE DIFFERENCE METHODS

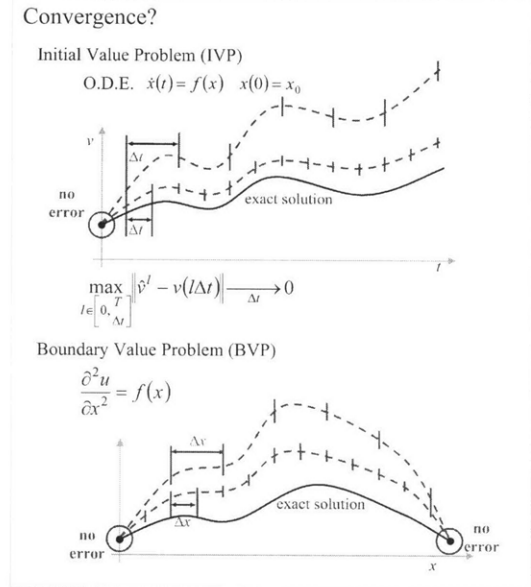
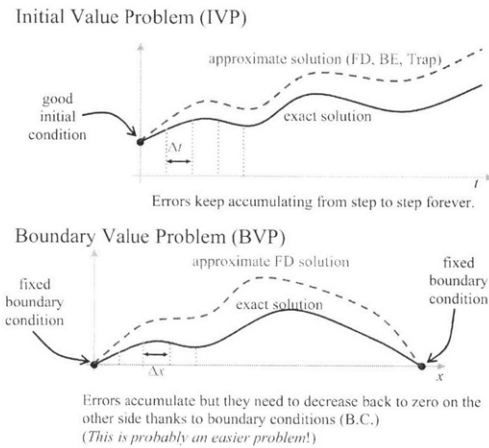
- ❖ Heat Conducting Bar
 1. Cut the bar into short sections
 2. Assign each cut a temperature



Normalized Equation

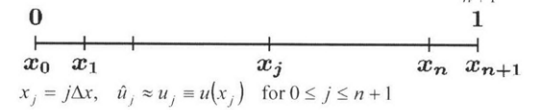
$$\frac{\partial}{\partial x} \kappa \frac{\partial T(x)}{\partial x} = -h_s \Rightarrow -\frac{\partial^2 u(x)}{\partial x^2} = f(x)$$

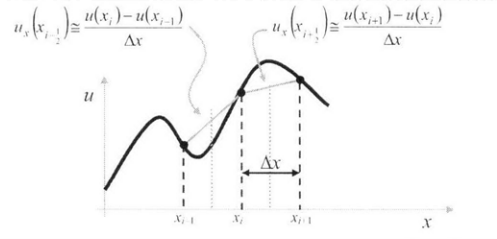
$$-u_{xx}(x) = f(x)$$



❖ Numerical Solution

Subdivide interval (0,1) into $n + 1$ equal subintervals: $\Delta x = \frac{1}{n+1}$

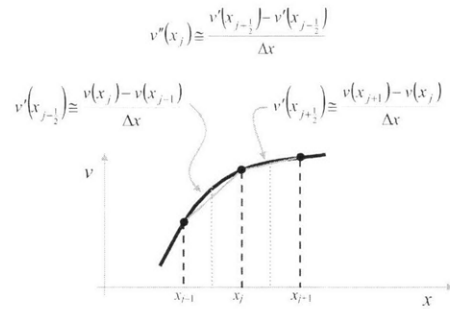




For example...

$$v''(x_j) \approx \frac{1}{\Delta x} \left[\frac{v(x_{j+1}) - v(x_j)}{\Delta x} - \frac{v(x_j) - v(x_{j-1}))}{\Delta x} \right]$$

$$= \frac{v_{j+1} - 2v_j + v_{j-1}}{\Delta x^2} \quad \text{for } \Delta x \text{ small}$$

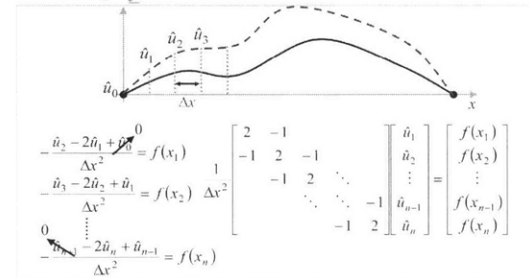


$-u_{xx} = f$ suggests...

$$-\frac{\hat{u}_{j+1} - 2\hat{u}_j + \hat{u}_{j-1}}{\Delta x^2} = f(x_j) \quad 1 \leq j \leq n$$

$$\hat{u}_0 = \hat{u}_{n+1} = 0$$

$$\Rightarrow \mathbf{A}\hat{\mathbf{u}} = \mathbf{f}$$



Initial Value Problem – Backward Euler

$$\hat{x}^j - \hat{x}^{j-1} = \Delta t [\lambda \hat{x}^j + u^j] \quad \leftarrow \text{test problem}$$

solve difference equation

no end point to match

Boundary Value Problem – Finite Difference Method

$$-\hat{u}_{j+1} + 2\hat{u}_j - \hat{u}_{j-1} = \Delta x^2 f(x_j)$$

solve difference equation

an end point to match

Backward Euler

$$\hat{x}^1 - \hat{x}^0 = \Delta t \lambda [\lambda \hat{x}^1 + u(t_1)] \Rightarrow \hat{x}^1 - x(0) = \Delta t \lambda [\lambda \hat{x}^1 + u(t_1)]$$

$$\begin{bmatrix} 1 - \lambda \Delta t & & & & \\ -1 & 1 - \lambda \Delta t & & & \\ & \ddots & \ddots & & \\ & & & -1 & 1 - \lambda \Delta t \end{bmatrix} \begin{bmatrix} \hat{x}^1 \\ \hat{x}^2 \\ \vdots \\ \hat{x}^j \end{bmatrix} = \begin{bmatrix} u(t_1) + x(0) \\ u(t_2) \\ \vdots \\ u(t_j) \end{bmatrix}$$

Lower Triangular input

$$\mathbf{A} = \frac{1}{\Delta x^2} \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix}, \quad \hat{u} = \begin{pmatrix} \hat{u}_1 \\ \hat{u}_2 \\ \vdots \\ \hat{u}_{n-1} \\ \hat{u}_n \end{pmatrix}, \quad f = \begin{pmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_{n-1}) \\ f(x_n) \end{pmatrix}$$

(Symmetric)

$\mathbf{A} \in \mathfrak{R}^{n \times n}$ $\hat{u}, f \in \mathfrak{R}^n$

Is \mathbf{A} nonsingular?

For any $v = [v_1, v_2, \dots, v_n]^T$

$$v^T \mathbf{A} v = \frac{1}{\Delta x^2} \left(v_1^2 + \sum_{i=2}^n (v_i - v_{i-1})^2 + v_n^2 \right)$$

Hence $v^T \mathbf{A} v > 0$, for any $v \neq 0$ (\mathbf{A} is SPD)

$\mathbf{A} \hat{u} = f$: \hat{u} exists and is unique

1.) How to visualize that \mathbf{A} is non-singular.

If I put in some current (=heat) into the network of resistors (=bar) I get a unique set of voltages (=temperatures).

2.) Gershgorin Circle Theorem

We know $\lambda(\mathbf{A}) \geq 0$ not $\lambda(\mathbf{A}) > 0$

If I had some leakage

Note if $\lambda(\mathbf{A}) > 0 \forall \lambda$ positive definite $\Rightarrow \lambda(\mathbf{A}) \neq 0 \forall \lambda$ non-singular

Check $v^T \mathbf{A} v =$ sum of squares always positive unless: $\begin{cases} v_i = 0 \\ v_j = v_{j-1} \forall i \mapsto v = 0 \\ v_n = 0 \end{cases}$

❖ **COMPARING FEM AND FD IN 1-D**

❖ Residual Equation

Partial Differential Equation Form

$$-\frac{\partial^2 u}{\partial x^2} = f \quad u(0) = 0 \quad u(1) = 0$$

Step 1: Choose Basis Functions to represent the solution

$$u(x) \approx u_h(x) = \sum_{j=1}^n \omega_j \underbrace{\varphi_j(x)}_{\text{Basis Functions}}$$

Step 2: Generate equations for the basis functions weights setting residual orthogonal to some test functions

$$\int_0^1 \underbrace{\varphi_l(x)}_{\text{Test Functions}} R(x) dx = 0 \quad R(x) = \sum_{j=1}^n \omega_j \frac{d^2 \varphi_j(x)}{dx^2} + f(x)$$

❖ FEM – Basis Weights

○ Galerkin Scheme

Force the residual to be “orthogonal” to the basis functions

$$\int_0^1 \varphi_l(x) R(x) dx = 0$$

Generates n equations in n unknowns

$$\int_0^1 \varphi_l(x) \left[\sum_{j=1}^n \omega_j \frac{d^2 \varphi_j(x)}{dx^2} + f(x) \right] dx = 0 \quad l \in \{1, \dots, n\}$$

○ Linear System

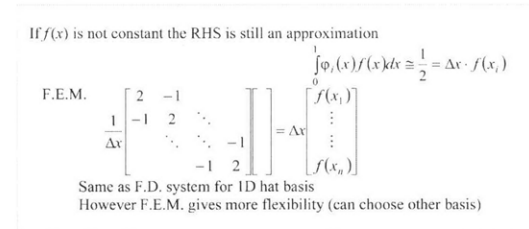
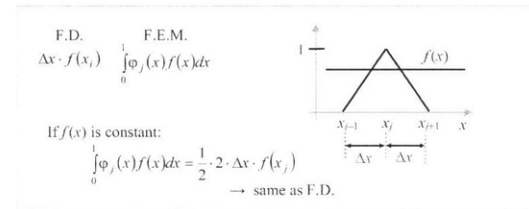
$$\sum_{j=1}^n \omega_j \int_0^1 \underbrace{\frac{d\varphi_j}{dx} \frac{d\varphi_l}{dx}}_{A_{l,j}} dx = \int_0^1 \underbrace{\varphi_l f(x)}_{F_l} dx$$

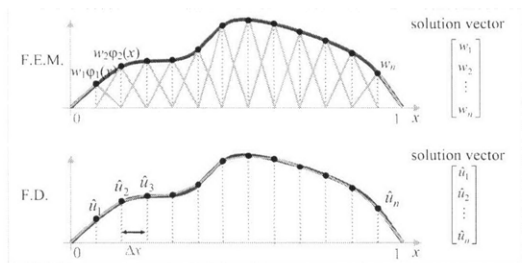
$$\frac{1}{\Delta x} \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & -1 & 2 & \ddots & \\ & & \ddots & \ddots & -1 \\ & & & -1 & 2 \end{bmatrix} \begin{bmatrix} \omega_1 \\ \omega_2 \\ \vdots \\ \omega_n \end{bmatrix} = \begin{bmatrix} F_1 \\ F_2 \\ \vdots \\ F_n \end{bmatrix}$$

❖ Comparing – FD & FEM (hat basis) – 1D problem

$$\text{FEM} \quad \frac{1}{\Delta x} \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & \ddots & & \\ & \ddots & \ddots & -1 & \\ & & -1 & 2 & \\ & & & & \omega_n \end{bmatrix} \begin{bmatrix} \omega_1 \\ \omega_2 \\ \vdots \\ \omega_n \end{bmatrix} = \begin{bmatrix} \int_0^1 \varphi_1(x) f(x) dx \\ \int_0^1 \varphi_2(x) f(x) dx \\ \vdots \\ \int_0^1 \varphi_n(x) f(x) dx \end{bmatrix}$$

$$\text{FD} \quad \frac{1}{\Delta x} \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & \ddots & & \\ & \ddots & \ddots & -1 & \\ & & -1 & 2 & \\ & & & & \hat{u}_n \end{bmatrix} \begin{bmatrix} \hat{u}_1 \\ \hat{u}_2 \\ \vdots \\ \hat{u}_n \end{bmatrix} = \begin{bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_n) \end{bmatrix}$$



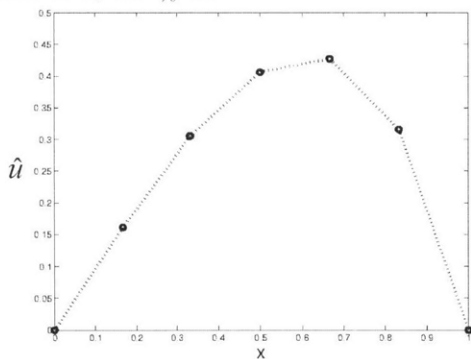


CONVERGENCE ANALYSIS FOR HEAT EQUATION

❖ Local Truncation Error and Consistency

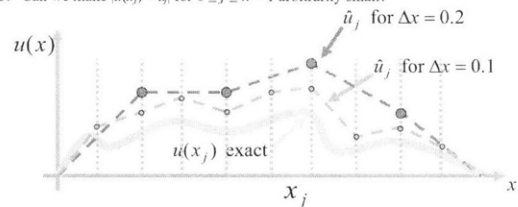
○ Example

$-u_{xx} = (3x + x^2)e^x, \quad x \in (0,1) \quad \text{with} \quad u(0) = u(1) = 0$
 Take $n = 5, \quad \Delta x = 1/6 \dots$



○ Convergence?

1. Does the discrete solution \hat{u} retain the qualitative properties of the continuous solution $u(x)$?
2. Does the solution become more accurate when $\Delta x \rightarrow 0$?
3. Can we make $|u(x_j) - \hat{u}_j|$ for $0 \leq j \leq n+1$ arbitrarily small?



○ Truncation Error

For any $v \in C^4$ we can show that

$$\frac{v(x_{j+1}) - 2v(x_j) + v(x_{j-1}))}{\Delta x^2} = v''(x_j) + \frac{\Delta x^2}{12} v^{(4)}(x_j + \theta \Delta x) \quad -1 \leq \theta \leq 1$$

Take $u \equiv v, \quad (-u'' = f)$

$$\frac{u(x_{j+1}) - 2u(x_j) + u(x_{j-1}))}{\Delta x^2} = f(x_j) + \underbrace{\frac{\Delta x^2}{12} u^{(4)}(x_j + \theta_j \Delta x)}_{\tau_j}$$

$$e_i = e(x_j) = \underbrace{u(x_j)}_{\text{exact solution}} - \underbrace{\hat{u}_j}_{\text{F.D. solution}}$$

$u(x)$ exact solution solves $-u'' = f$ exact equation

\hat{u}_i F.D. solution solves $-\frac{\hat{u}_{i-1} - 2\hat{u}_i + \hat{u}_{i+1}}{\Delta x^2} = f_i$ F.D. equation

Local truncation error (LTE) tells me how well the exact solution solves the F.D. equation.

$$-\left[\frac{u(x_{j-1}) - 2u(x_j) + u(x_{j+1}))}{\Delta x^2} \right] - f_j = \tau_j \quad \text{LTE}$$

hopefully small!!!

How do we calculate the L.T.E. τ_i ?

Use: Taylor Series

$$\frac{u(x_{j+1}) - 2u(x_j) + u(x_{j-1}))}{\Delta x^2} \cong \frac{\partial^2 u}{\partial x^2}$$

$$u(x_{j+1}) = u(x_j) + \Delta x \cdot u'(x_j) + \frac{\Delta x^2}{2} u''(x_j) + \frac{\Delta x^3}{3!} u'''(x_j) + \frac{\Delta x^4}{4!} u^{(4)}(x) \Big|_{x \in [x_{j-1}, x_{j+1}]}$$

$$-2u(x_j) = -2u(x_j)$$

$$u(x_{j-1}) = u(x_j) - \Delta x \cdot u'(x_j) + \frac{\Delta x^2}{2} u''(x_j) - \frac{\Delta x^3}{3!} u'''(x_j) + \frac{\Delta x^4}{4!} u^{(4)}(x) \Big|_{x \in [x_{j-1}, x_{j+1}]}$$

$$u(x_{j+1}) - 2u(x_j) + u(x_{j-1})) = \Delta x^2 u''(x_j) + \frac{\Delta x^4}{12} u^{(4)}(x) \Big|_{x \in [x_{j-1}, x_{j+1}]}$$

Hence $u''(x_j) = \frac{u(x_{j+1}) - 2u(x_j) + u(x_{j-1}))}{\Delta x^2} + \frac{\Delta x^2}{12} u^{(4)}(x) \Big|_{x \in [x_{j-1}, x_{j+1}]}$

Substitution back and get τ_j L.T.E.

$$-\left[u''(x_j) - \frac{\Delta x^2}{12} u^{(4)}(x) \right] - f_j = \tau_j \Rightarrow \tau_j = \frac{\Delta x^2}{12} u^{(4)}(x) \Big|_{x \in [x_{j-1}, x_{j+1}]}$$

Since $u'' = f$

- ❖ Global Error and Stability
 - Discretization Error Analysis – Error Equation
 - Let $e_j = u(x_j) - \hat{u}_j$ be the discretization error.
 - $$\frac{u(x_{j+1}) - 2u(x_j) + u(x_{j-1}))}{\Delta x^2} = f(x_j) + \tau_j$$
 - $$\frac{\hat{u}_{j+1} - 2\hat{u}_j + \hat{u}_{j-1}}{\Delta x^2} = f(x_j)$$
- Subtracting
- $$\frac{-e_{j+1} - 2e_j + e_{j-1}}{\Delta x^2} = \tau_j \quad 1 \leq j \leq n$$
- and $e_0 = e_{n+1} = 0$

What we really want is the global error

$$e_j = u(x_j) - \hat{u}_j$$

$$\frac{u(x_{j+1}) - 2u(x_j) + u(x_{j+1}))}{\Delta x^2} = f(x_j) + \tau_j \Big|_{\text{Exact Solution into F.D. Equation}}$$

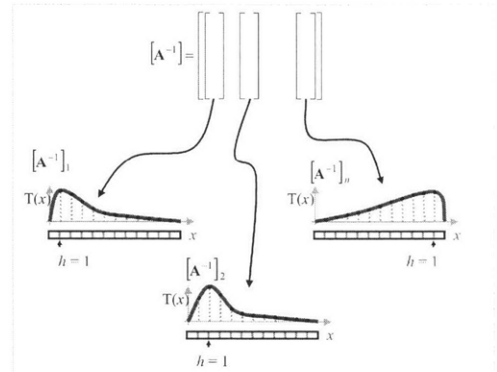
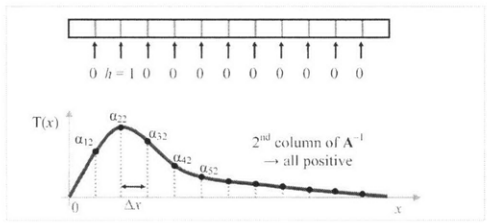
$$\frac{\hat{u}_{j+1} - 2\hat{u}_j + \hat{u}_{j-1}}{\Delta x^2} = f(x_j) \Big|_{\text{F.D. Solution solves F.D. Equation}}$$

$$\frac{-e_{j+1} - 2e_j + e_{j-1}}{\Delta x^2} = \tau_j \Big|_{\text{Subtract and get a F.D. Equation for the Global Error}}$$

$e_0 = u(x_0) - \hat{u}_0 = 0$ B.C. \rightarrow This generates a linear system that can be solved for $\begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$

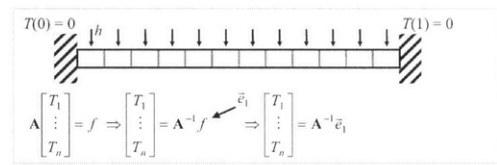
- $Ae = \tau$
- $$e = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}, \quad \tau = \frac{\Delta x^2}{12} \begin{bmatrix} u^{(4)}(x_1 + \theta_1 \Delta x) \\ u^{(4)}(x_2 + \theta_2 \Delta x) \\ \vdots \\ u^{(4)}(x_n + \theta_n \Delta x) \end{bmatrix}$$
- Properties of A^{-1}
 - Let $A^{-1} = \{a_{ij}\}_{1 \leq i, j \leq n}$
 - Non-negativity $a_{ij} \geq 0$, for $1 \leq i, j \leq n$
 - Boundedness $0 \leq \sum_{j=1}^n a_{ij} \leq \frac{1}{8}$, for $1 \leq i \leq n$

$e = \mathbf{A}^{-1}\tau$
 $\mathbf{A}^{-1} = \{\alpha_{ij}\}$ α_{ij} are all positive... WHY?
 How do we visualize the inverse of a matrix?
 $\mathbf{A} \begin{bmatrix} T_1 \\ T_2 \\ \vdots \\ T_n \end{bmatrix} = \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_n \end{bmatrix}$ $\mathbf{A}^{-1} \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_n \end{bmatrix} = \begin{bmatrix} T_1 \\ T_2 \\ \vdots \\ T_n \end{bmatrix}$
 How do we get column j of \mathbf{A}^{-1} ?
 Just pick an appropriate set of heat sources.
 e.g. $\begin{bmatrix} h_1 \\ \vdots \\ h_j = 1 \\ \vdots \\ h_n \end{bmatrix} \leftarrow j \Rightarrow [\mathbf{A}^{-1}]_j = \begin{bmatrix} T_1 \\ T_2 \\ \vdots \\ T_n \end{bmatrix}$



- Qualitative Properties of \hat{u}
 $f \geq 0 \rightarrow \hat{u} \geq 0$
 $\hat{u} = \mathbf{A}^{-1}f$
 If $f_j = f(x_j) \geq 0$, for $1 \leq j \leq n$
 Then $\hat{u}_i = \sum_j \alpha_{ij} f_j \geq 0$, for $1 \leq i \leq n$

\hat{u} is just a linear combination of columns of \mathbf{A}^{-1}
 The columns are all positive if the coefficients are positive
 → the combination is positive.



Discrete Stability
 $\hat{u} = \mathbf{A}^{-1}f$

$$\|\hat{u}\|_\infty = \max_i |\hat{u}_i| = \max_i \left| \sum_j \alpha_{ij} f_j \right| \leq \max_i \left(\sum_j \alpha_{ij} \right) \max_j |f_j| \leq \frac{1}{8} \|f\|_\infty$$

o Convergence

Using the discrete stability estimate on $Ae = \tau$

$$\|e\|_\infty \leq \frac{1}{8} \tau \|e\|_\infty$$

or

$$\max_{1 \leq i \leq n} u(x_i) - \hat{u}_i \leq \frac{\Delta x^2}{96} \max_{0 \leq x \leq 1} |u^{(4)}(x)|$$

A- priori Error Estimate

$Ae = \tau$
 $e = A^{-1}\tau$ ← Global Error
 Local Truncation Error

$\|\tau\|_\infty \leq \frac{\Delta x^2}{12} |u^{(4)}(x)|$ bounded

$\|e\|_\infty = \max_{i \in [1, n]} |e_i| = \|A^{-1}\tau\|_\infty \leq \|A^{-1}\|_\infty \|\tau\|_\infty$
Stability, L.T.E.

$\|A^{-1}\|_\infty = \max_{i \in [1, n]} \sum_j |A^{-1}_{ij}| = \max_{i \in [1, n]} \|A^{-1}x\|_\infty = \max_{\text{all rows } i} \sum_j \alpha_{ij}$ max row sum

$A = \frac{1}{\Delta x^2} \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & & & \\ & & \ddots & \ddots & \\ & & & -1 & 2 \\ & & & & -1 & 2 \end{bmatrix}$

We know A^{-1} has all positive entries
 What x should we pick to maximize $\|A^{-1}x\|_\infty$? $x = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$
 Can look at max column sum since A, A^{-1} symmetric

$T(0)=0$ $\hat{A}^{-1} \tau = \|A^{-1} \tau\|_\infty$ $T(1)=0$

$$\frac{1}{\Delta x^2} \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & & & \\ & & \ddots & \ddots & \\ & & & -1 & 2 \\ & & & & -1 & 2 \end{bmatrix} \begin{bmatrix} T_1 \\ T_2 \\ \vdots \\ T_n \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$$

← h_i

$$\Rightarrow \begin{bmatrix} T_1 \\ T_2 \\ \vdots \\ T_n \end{bmatrix} = \Delta x^2 \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & & & \\ & & \ddots & \ddots & \\ & & & -1 & 2 \\ & & & & -1 & 2 \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$$

⇒ column j^{th} of $\Delta x^2 \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & & & \\ & & \ddots & \ddots & \\ & & & -1 & 2 \\ & & & & -1 & 2 \end{bmatrix}^{-1}$

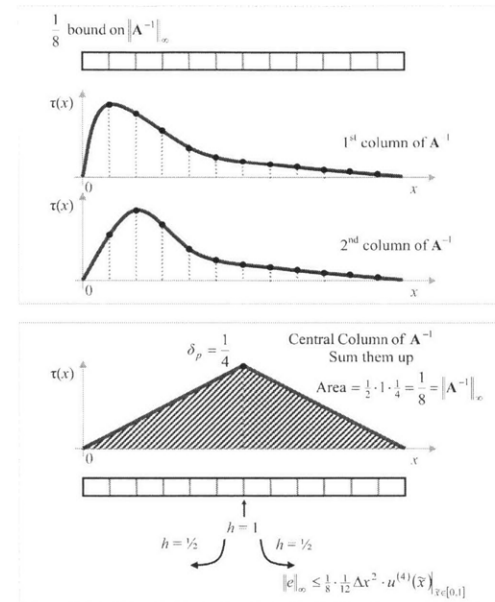
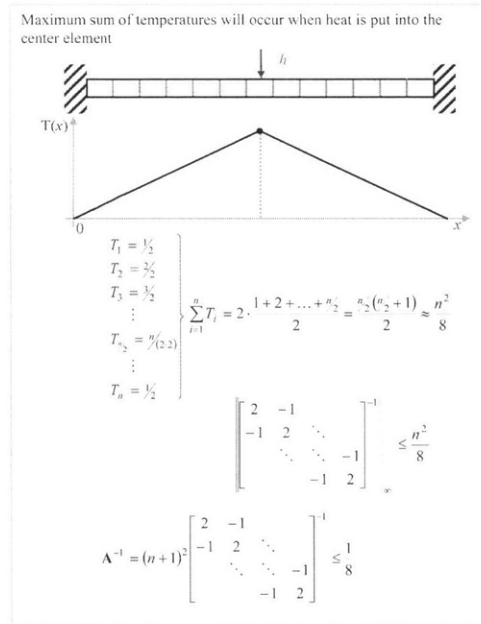
	Lossy	Insulated
LTE	same	same
$\ A^{-1}\ $	smaller	larger

$$A = \frac{1}{\Delta x^2} \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & & & \\ & & \ddots & \ddots & \\ & & & -1 & 2 \\ & & & & -1 & 2 \end{bmatrix} = \frac{1}{(n+1)^2} \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & & & \\ & & \ddots & \ddots & \\ & & & -1 & 2 \\ & & & & -1 & 2 \end{bmatrix}$$

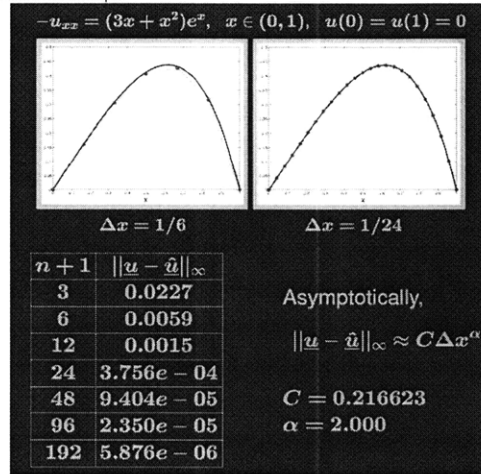
$$A^{-1} = (n+1)^2 \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & & & \\ & & \ddots & \ddots & \\ & & & -1 & 2 \\ & & & & -1 & 2 \end{bmatrix}^{-1} \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & & & \\ & & \ddots & \ddots & \\ & & & -1 & 2 \\ & & & & -1 & 2 \end{bmatrix}^{-1}$$

$\|A^{-1}\|_\infty = \max_{i \in [1, n]} \sum_j |A^{-1}_{ij}| = \max_{\text{all rows } i} \sum_j \alpha_{ij}$ max row sum

$x = \begin{bmatrix} \pm 1 \\ \pm 1 \\ \vdots \\ \pm 1 \end{bmatrix}$



o Numerical Example



o Summary

- For a simple model problem we can produce numerical approximations of **arbitrary accuracy**.
- An **a-priori error estimate** gives the asymptotic dependence of the solution error on the discretization size Δx .

FORMAL GENERALIZATION TO ANY PDE

o Definitions

Consider a linear elliptic differential equation $\mathcal{L}u = f$ and a difference scheme $\hat{\mathcal{L}}\hat{u} = \hat{f}$

e.g. Different PDE

$\mathcal{L}u = f \quad - \frac{\partial^2 u}{\partial x^2} + \alpha u(x) = f(x) \quad \text{Helmholtz}$

$\hat{\mathcal{L}}\hat{u} = \hat{f} \quad - \frac{\hat{u}_{j-1} - 2\hat{u}_j + \hat{u}_{j+1}}{\Delta x^2} + \alpha \hat{u}_j = f_j \quad \text{F.D. eq on F.D. solution } \hat{u}$

o Consistency + Stability yields Convergence

o Consistency

The difference scheme is **consistent** with the differential equation if:

For all smooth functions v

$(\hat{\mathcal{L}}v - \hat{f})_j \rightarrow 0, \quad \text{for } j = 1, \dots, n \quad \text{when } \Delta x \rightarrow 0.$

$(\hat{\mathcal{L}}v - \hat{f})_j \rightarrow 0, \quad \text{for } j = 1, \dots, n \quad \text{when } \Delta x \rightarrow 0.$

$(\hat{\mathcal{L}}v - \hat{f})_j = \mathcal{O}(\Delta x^p)$ for all $j \Rightarrow p$ is **order of accuracy**

Consistency LTE $\tau \rightarrow 0$ if $\Delta x \rightarrow 0$

$\hat{\mathcal{L}}u - \hat{f} = \tau$ Exact solution u into F.D. equation

$$-\frac{u(x_{j-1}) - 2u(x_j) + u(x_{j+1}))}{\Delta x^2} + \alpha u(x_j) - f(x_j) = \tau_j$$

In this case note $\tau_j = \frac{1}{12} \Delta x^2 u^{(4)}(x_j) + \dots$

As for Poisson Eq. $p = 2$

o Truncation Error

$(\hat{\mathcal{L}}u - \hat{f})_j = \tau_j, \quad \text{for } j = 1, \dots, n \quad \text{or, } \hat{\mathcal{L}}u - \hat{f} = \tau.$

The truncation error results from inserting the exact solution into the difference scheme.

Consistency $\Rightarrow \|\tau\|_\infty = \mathcal{O}(\Delta x^p)$

o Error Equation

Original Scheme $\hat{\mathcal{L}}\hat{u} = \hat{f}$

Consistency $\hat{\mathcal{L}}u = \hat{f} + \tau$

The error $e = u - \hat{u}$ satisfies $\hat{\mathcal{L}}e = \tau$

$\hat{f}e = \tau$	$e = u - \hat{u}$	
$-\frac{u(x_{j-1}) - 2u(x_j) + u(x_{j+1}))}{\Delta x^2} + \alpha u(x_j) - f(x_j) = \tau_j$		Exact Solution into F.D. Equation
$-\frac{\hat{u}_{j-1} - 2\hat{u}_j + \hat{u}_{j+1}}{\Delta x^2} + \alpha \hat{u}_j - f_j = 0$		F.D. Solution \hat{u} into F.D. Equation
<hr/>		
$-\frac{e_{j-1} - 2e_j + e_{j+1}}{\Delta x^2} + \alpha e_j = \tau_j$		Subtract

INTRODUCTION TO NUMERICAL SIMULATION

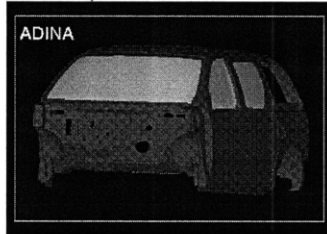
LECTURE 19.
Boundary Value Problems – Solving 3-D Finite-Difference Problems

TODAY'S OUTLINE:

- ❖ Finite Difference Matrices in 1-D, 2-D, and 3-D
- ❖ Gaussian Elimination Cost:
 - Bandlimited GE
 - Sparse GE
- ❖ Krylov Method Cost
 - Counting iterations: Communication Lower Bound

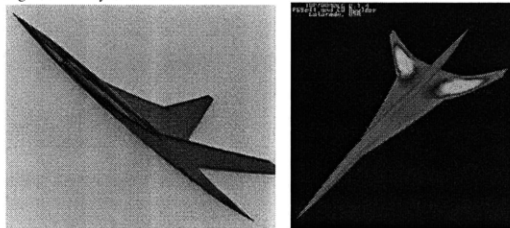
FINITE DIFFERENCE MATRICES IN 1-D, 2-D, AND 3-D

- ❖ Structural Analysis of Automobiles



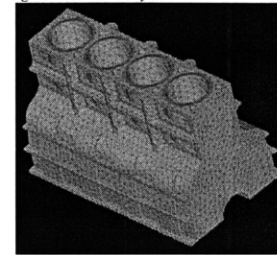
- Equations
 - Force-displacement relationships for mechanical elements (beams, plates, shells) and sum of forces = 0.
 - Partial Differential Equations of Continuum Mechanics

- ❖ Drag Force Analysis of Aircraft

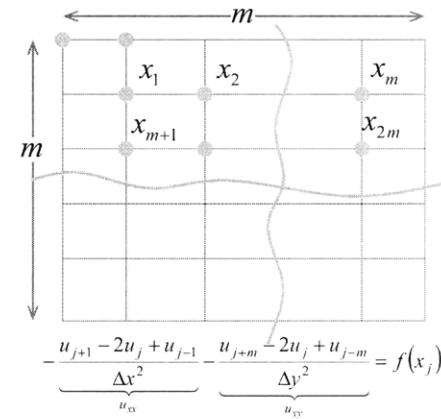


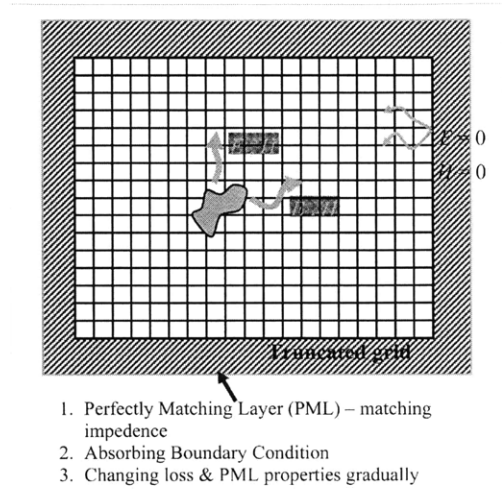
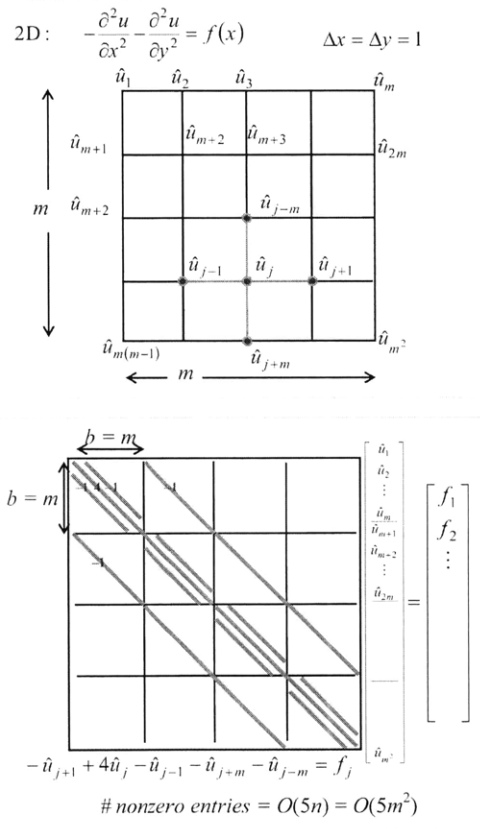
- Equations
 - Navier-Stokes Partial Differential Equations.

- ❖ Engine Thermal Analysis

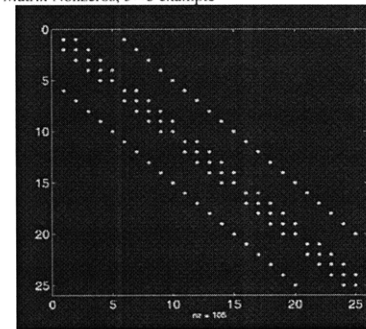


- Equations
 - The Poisson Partial Differential Equation
- ❖ FD Matrix Properties
 - 2-D Discretized Problem
 - Discretized Poisson





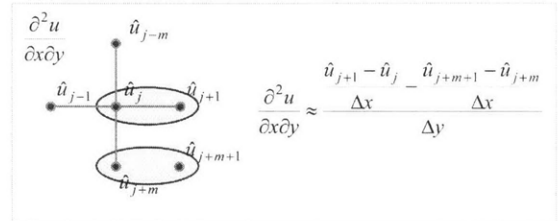
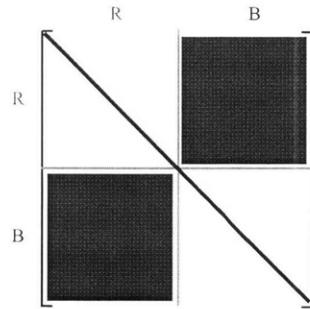
– Matrix Nonzeros, 5x5 example



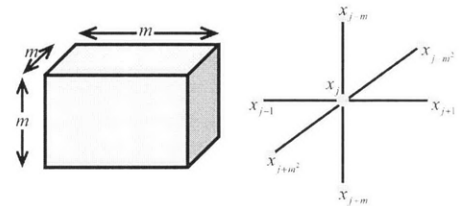
Red – Black Ordering
versus “Natural” Ordering

- 1 ● 33 ● 2 ● 34 ● 3 ● 35 ● 4 ● 36 ●
- 37 ● 5 ● 38 ● 6 ● 39 ● 7 ● 40 ● 8 ●
- 9 ● 41 ● 10 ● 42 ● 11 ● 43 ● 12 ● 44 ●
- 45 ● 13 ● 46 ● 14 ● 47 ● 15 ● 48 ● 16 ●
- 17 ● 49 ● 18 ● 50 ● 19 ● 51 ● 20 ● 52 ●
- 53 ● 21 ● 54 ● 22 ● 55 ● 23 ● 56 ● 24 ●
- 25 ● 57 ● 26 ● 58 ● 27 ● 59 ● 28 ● 60 ●
- 61 ● 29 ● 62 ● 30 ● 63 ● 31 ● 64 ● 32 ●

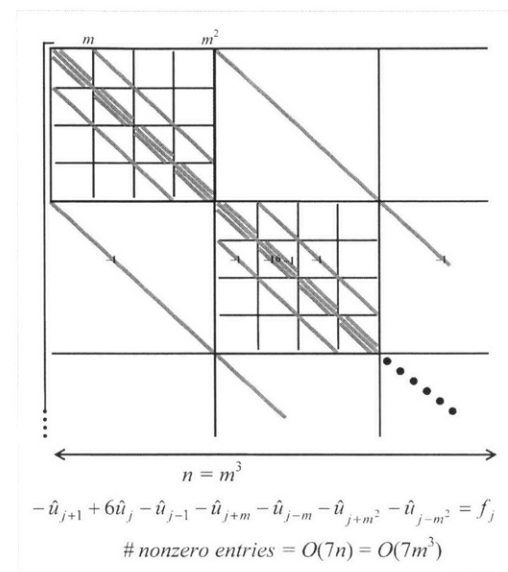
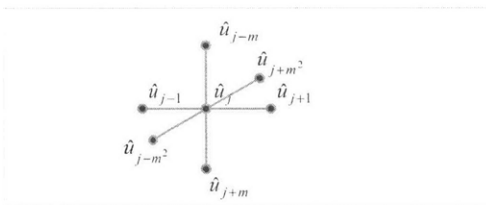
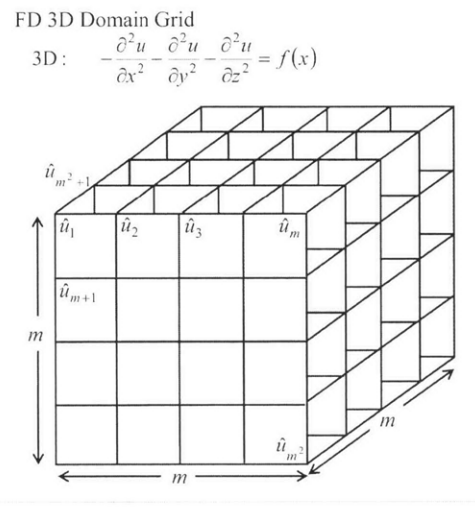
Red – Black Ordering



- o 3-D Discretized Problem
- Discretized Poisson

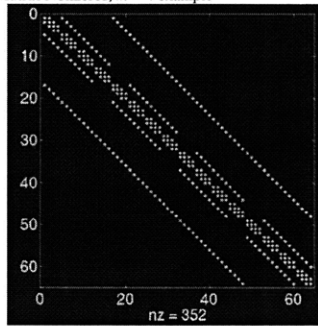


$$\underbrace{\frac{\hat{u}_{j+1} - 2\hat{u}_j + \hat{u}_{j-1}}{\Delta x^2}}_{u_{xx}} - \underbrace{\frac{\hat{u}_{j+m} - 2\hat{u}_j + \hat{u}_{j-m}}{\Delta y^2}}_{u_{yy}} - \underbrace{\frac{\hat{u}_{j+m^2} - 2\hat{u}_j + \hat{u}_{j-m^2}}{\Delta z^2}}_{u_{zz}} = f(x_j)$$



- 1D: $b = 1 \quad n = m$
- 2D: $b = m \quad n = m^2$
- 3D: $b = m^2 \quad n = m^3$

- Matrix Nonzeros, $m = 4$ example



o Summary

- Numerical Properties

Matrix is irreducibly diagonally dominant

$$\left(|A_{ii}| \geq \sum_{j \neq i} |A_{ij}| \right)$$

Matrix is symmetric positive definite

Assuming uniform discretization, diagonal is

$$1-D \rightarrow \frac{2}{\Delta^2}, \quad 2-D \rightarrow \frac{4}{\Delta^2}, \quad 3-D \rightarrow \frac{6}{\Delta^2}$$

- Structural Properties

Matrices in 3D are LARGE

$$1-D \rightarrow m \times m, \quad 2-D \rightarrow m^2 \times m^2, \quad 3-D \rightarrow m^3 \times m^3$$

$0 \times 100 \times 100$ grid in 3D = 1 million \times 1 million matrix

Matrices are very sparse

Nonzeros per row 1-D: 3, 2-D: 5, 3-D: 7

Matrices are banded

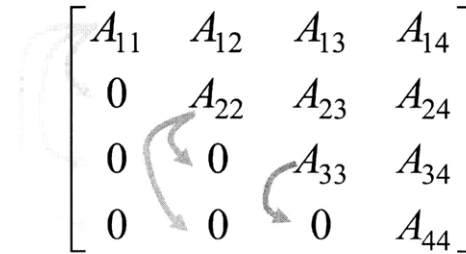
$$1-D \quad A_{ij} = 0 \quad |i - j| > 1 = b$$

$$2-D \quad A_{ij} = 0 \quad |i - j| > m = b$$

$$3-D \quad A_{ij} = 0 \quad |i - j| > m^2 = b$$

GAUSSIAN ELIMINATION COST

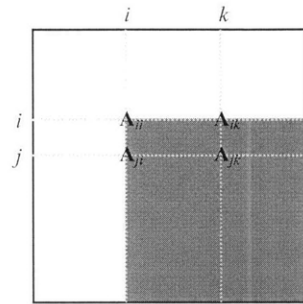
- ❖ Bandlimited GE
 - o Dense GE
 - Picture



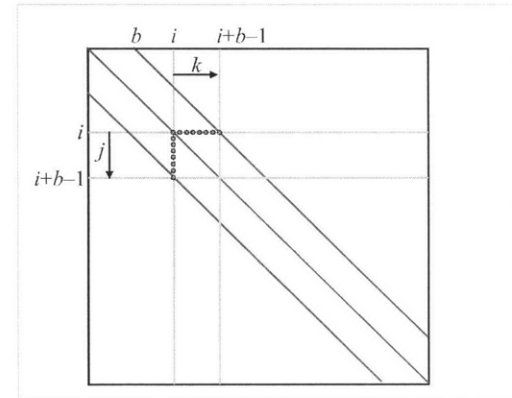
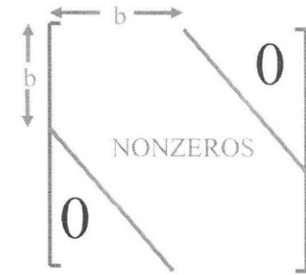
- Algorithm

```

For i = 1 to n - 1 {
  "For each row"
  For j = i + 1 to n {
    "For each row below pivot"
    Form n-1 reciprocals (pivots)
     $\tilde{A}_{ji} \leftarrow \frac{A_{ji}}{A_{ii}}$ 
    Form  $\sum_{i=1}^{n-1} (n-i) = \frac{n^2}{2}$  multipliers
  }
  For k = i + 1 to n {
    "For each element beyond Pivot"
     $\tilde{A}_{jk} \leftarrow A_{jk} - \tilde{A}_{ji} A_{ik}$ 
    Perform  $\sum_{i=1}^{n-1} (n-i)^2 \approx \frac{2}{3} n^3$  multiply-adds
  }
}
    
```



- Complexity
 - 1-D $O(n^3) = O(m^3) \leftarrow 100 \text{ pt grid}$ $O(10^6)$ ops 10 msec
 - 2-D $O(n^3) = O(m^6) \leftarrow 100 \times 100 \text{ grid}$ $O(10^9)$ ops 2h 46min
 - 3-D $O(n^3) = O(m^9) \leftarrow 100 \times 100 \times 100 \text{ grid}$ $O(10^{27})$ ops 317 years
 - e.g. on a 1 GFlops computer using $m = 100$ and a constant $c = 10$
 - For 2-D and 3-D problems – need a faster solver!
 - o Banded GE
 - Triangularizing Algorithm
 - For $i = 1$ to $n - 1$ {
 - For $j = i + 1$ to $i + b - 1$ {
 - For $k = i + 1$ to $i + b - 1$ {
 - $A_{jk} \leftarrow A_{jk} - \frac{A_{ji}}{A_{ii}} A_{ik}$
- Perform $\sum_{i=1}^{n-1} [\min(b-1, n-i)]^2 \approx O(b^2 n)$ Multiply-Adds

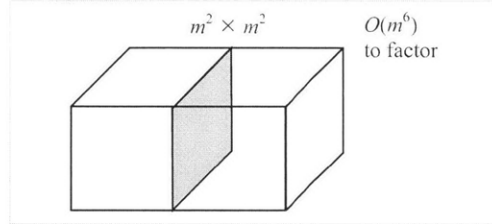


- Complexity

	matrix size	band	
1-D	$n = m$	$b = 1$	domain size = m
2-D	$n = m^2$	$b = m$	
3-D	$n = m^3$	$b = m^2$	

 - 1-D $O(b^2 n) = O(m)$ $\leftarrow 100 \text{ pt grid}$ $O(100)$ ops 1 μ sec
 - 2-D $O(b^2 n) = O(m^3) \leftarrow 100 \times 100 \text{ grid}$ $O(10^6)$ ops 1 sec

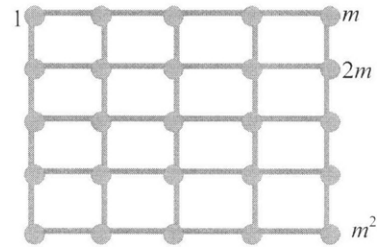
3-D $O(b^3n) = O(m^3) \leftarrow 100 \times 100 \times 100$ grid $O(10^{14})$ ops 11d 13h
 e.g. on a 1 GFlops computer using $m = 100$ and a constant $c = 10$
 For 3-D problems – still need a faster solver!



❖ Sparse GE

○ Matrix Graphs – Grid Example

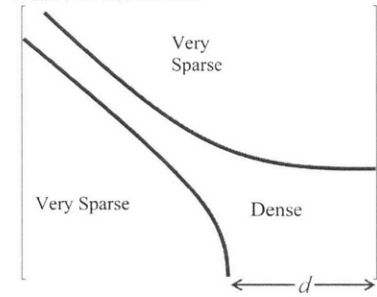
How long does it take to factor an $m \times m$ grid



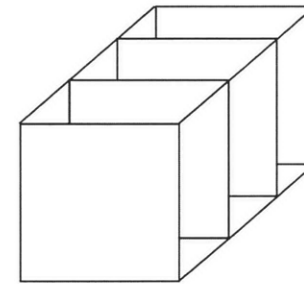
Suppose the center column is eliminated last?

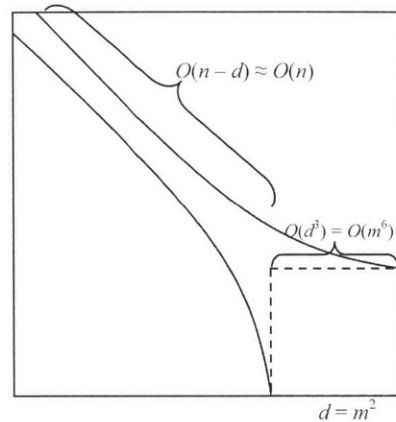
○ Fill-In

Pattern of a Filled-in Matrix



2-D: $d = m$
 3-D: $d = m^2$





G.E. Complexity
 $O(n) + O(d^3)$

- Complexity of Sparse GE

	matrix size	dense block	
1-D	$n = m$	$d = 1$	domain size = m
2-D	$n = m^2$	$d = m$	
3-D	$n = m^3$	$d = m^2$	
1-D	$O(n) + O(d^3) = O(m)$		← 100 pt grid $O(100)$ ops 1 μsec
2-D	$O(n) + O(d^3) = O(m^3)$		← 100 × 100 grid $O(10^6)$ ops 10 msec
3-D	$O(n) + O(d^3) = O(m^6)$		← 100 × 100 × 100 grid $O(10^{12})$ ops 2h 46m

e.g. on a 1 GFlops computer using $m = 100$ and a constant $c = 10$
 For 3-D problems – can we do any better with GCR?

KRYLOV METHOD COST

- ❖ The Generalized Conjugate Residual Algorithm
 - Algorithm Complexity

Compute $\mathbf{A}r^{k-1}$ Sparse Matrix-vector product costs $O(n)$

$$\tilde{p}_{k-1} \leftarrow r^{k-1} - \sum_{j=0}^{k-1} (\mathbf{A}p_j)^T (\mathbf{A}r^{k-1}) p_j$$

inner product : $O(n)$

$$p_{k-1} \leftarrow \frac{\tilde{p}_{k-1}}{\|\mathbf{A}\tilde{p}_{k-1}\|}$$

inner product : $O(n)$

$$y_{k-1} \leftarrow (r^{k-1})^T \mathbf{A}p_{k-1}$$

$$x^k \leftarrow x^{k-1} + y_{k-1} p_{k-1}$$

multiplications : $O(n)$

$$r^k \leftarrow r^{k-1} - y_{k-1} \mathbf{A}p_{k-1}$$

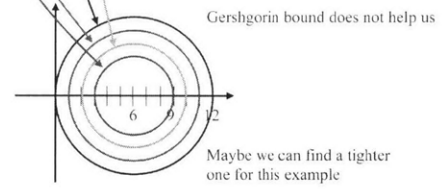
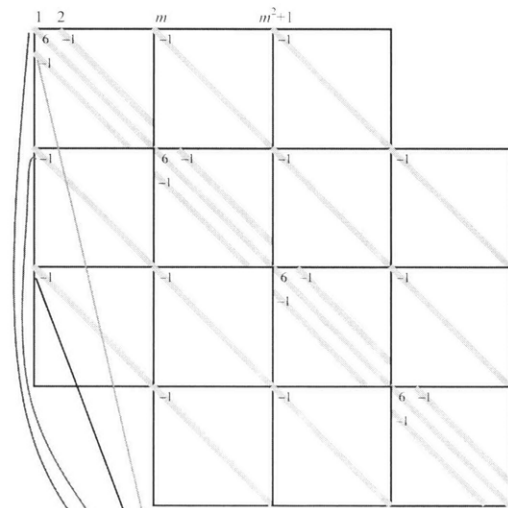
Algorithm is $O(kn)$ where k = number of iterations

- Complexity of sparse GCR (general worst case iterations)

	matrix size	worst case iterations ($k = n$)	
1-D	$n = m$	$k = m$	domain size = m
2-D	$n = m^2$	$k = m^2$	
3-D	$n = m^3$	$k = m^3$	
1-D	$O(kn) = O(m^2)$		← 100 pt grid $O(10^4)$ ops 0.1 ms
2-D	$O(kn) = O(m^4)$		← 100 × 100 grid $O(10^8)$ ops 1 sec
3-D	$O(kn) = O(m^6)$		← 100 × 100 × 100 grid $O(10^{12})$ ops 2h 46m

e.g. on a 1 GFlops computer using $m = 100$ and a constant $c = 10$
 But how many iterations k does it really take?

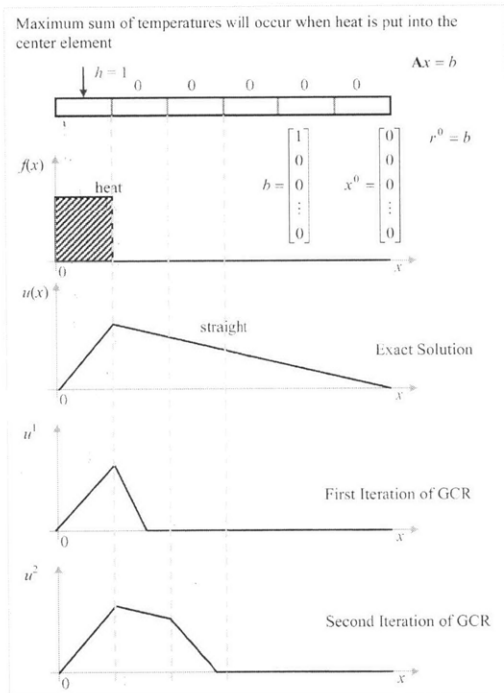
Let's estimate k # of iterations



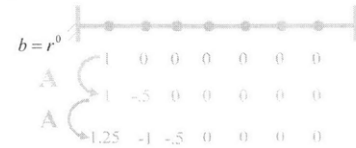
- ❖ Counting iterations
 - 1-D Case

1-D Discretized PDE

$$\begin{array}{c}
 b = r^0 \\
 \begin{array}{cccccccc}
 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
 A & 1 & -5 & 0 & 0 & 0 & 0 & 0 \\
 A & -1.25 & -1 & -5 & 0 & 0 & 0 & 0
 \end{array} \\
 \\
 \underbrace{\begin{bmatrix} 1 & -\frac{1}{2} & & & & & & \\ -\frac{1}{2} & 1 & -\frac{1}{2} & & & & & \\ & -\frac{1}{2} & 1 & -\frac{1}{2} & & & & \\ & & \ddots & \ddots & \ddots & & & \\ & & & -\frac{1}{2} & 1 & & & \\ & & & & -\frac{1}{2} & 1 & & \\ & & & & & & -\frac{1}{2} & 1 \end{bmatrix}}_A \begin{bmatrix} 1 \\ -\frac{1}{2} \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} 1.25 \\ -1 \\ -\frac{1}{2} \\ 0 \end{bmatrix}
 \end{array}$$



o Communication Lower Bound

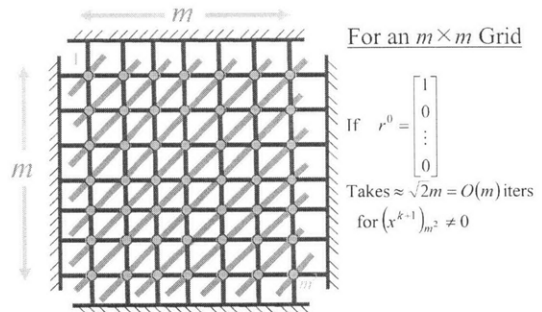


Communication Lower Bound for m gridpoints

$A^k r^0$ is nonzero in m^{th} entry after $k = m$ iterations

Need at least m iterations for $(x^{k+1})_m = \left(x^0 + \sum_{j=1}^k a_j r^j\right)_m \neq 0$

o 2-D Case



o Eigenanalysis

3D Condition number: $\kappa \equiv \frac{\lambda_{\max}}{\lambda_{\min}} = \left(1 - \cos\left(\frac{m\pi}{m+1}\right)\right) \left(1 - \cos\left(\frac{\pi}{m+1}\right)\right)^{-1}$

Number of Iterations for GCR to achieve convergence

$\left|\frac{r^k}{r^0}\right| \leq 2 \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^k \leq \gamma$ (convergence tolerance)

$$k = \frac{\log_2^7}{\log\left(\frac{\sqrt{k}-1}{\sqrt{k}+1}\right)^{m-2x}} \cong O(m)$$

GCR achieves communication lower bound $O(m)$!

- Complexity of sparse GCR (no precondition) for 3D FD methods

	matrix size	iterations ($k = n$)			
1-D	$n = m$	$k = m$	domain size = m		
2-D	$n = m^2$	$k = m$			
3-D	$n = m^3$	$k = m$			
1-D	$O(kn) = O(m^2)$	\leftarrow 100 pt grid	$O(10^4)$ ops	0.1 ms	
2-D	$O(kn) = O(m^3)$	\leftarrow 100x100 grid	$O(10^6)$ ops	10 ms	
3-D	$O(kn) = O(m^4)$	\leftarrow 100x100x100 grid	$O(10^8)$ ops	1 sec	

e.g. on a 1 GFlops computer using $m = 100$ and a constant $c = 10$
 But how many iterations k does it really take?

- Work for Banded Gaussian Elimination, Relaxation and GCR

Dimension	Banded GE	Sparse GE	GCR
1	$O(m)$	$O(m)$	$O(m^2)$
2	$O(m^4)$	$O(m^3)$	$O(m^3)$
3	$O(m^7)$	$O(m^6)$	$O(m^4)$

GCR faster than banded in GE in 2 and 3 dimensions
 Could be faster, 3-D matrix only m^3 nonzeros.
 Must defeat the communication lower bound!

$m = 100$	
3-D	SPG - $O(10^{13})$ ops \leftarrow 10,000 sec
	GCR - $O(10^8)$ ops \leftarrow 1 sec

Dimension	Banded GE	Sparse GE	GCR (no preconditioner)
1	1 μ s	1 μ s	0.1 ms
2	1 sec	10 ms	10 ms
3	11 days 6h	2h 46m	1 sec

e.g. on a 1 GFlops computer using $m = 100$, and constant $c = 10$

GCR faster than banded GE in 2 and 3 dimensions
 Could be faster, 3-D matrix only m^3 nonzeros.
 Must defeat the communication lower bound!

Dimension	Banded GE	Sparse GE	GCR (no preconditioner)
1	10 μ s	10 μ s	10 ms
2	2h 46m	10 sec	10 sec
3	317,000 years	317 years	2h 46m

e.g. on a 1 GFlops computer using $m = 1000$, and constant $c = 10$

Dimension	Banded GE	Sparse GE	GCR	GCR + preconditioner
1	$O(m)$	$O(m)$	$O(m^2)$	$O(km)$
2	$O(m^4)$	$O(m^3)$	$O(m^3)$	$O(km^2)$
3	$O(m^7)$	$O(m^6)$	$O(m^4)$	$O(km^3)$

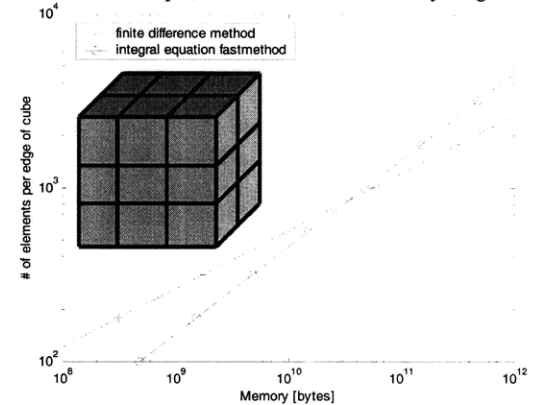
GCR+good preconditioner (i.e. $k = 5-10$) is always the best choice for 2D & 3D problems!

Dimension	Banded GE	Sparse GE	GCR	GCR + preconditioner
1	10 μ s	10 μ s	10 ms	0.1 ms
2	2h 46m	10 sec	10 sec	0.1 sec
3	317,000 years	317 years	2h 46m	1m 20s

e.g. on a 1 GFlops computer using $m = 1000$, $k = 10$, and constant $c = 10$

- Cube Example, Interior Problem - Memory Usage

Cube example, Interior Problem - Memory usage



INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 21.

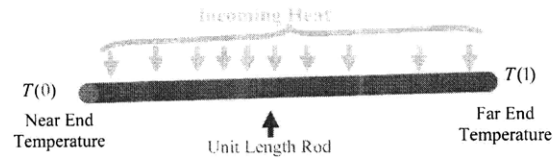
Finite Element Methods & GCR Preconditioners

TODAY'S OUTLINE:

- ❖ Finite Element Methods
 - The FEM Basis Functions
 - The Galerkin scheme
 - The FEM linear system
 - An FEM Example
 - Energy minimization view of FEM
- ❖ GCR Preconditioners for FD and FEM
 - Diagonal Preconditioners
 - Block Diagonal Preconditioners
 - Incomplete Factorization Preconditioners
 - Communication Improving Preconditioners: Gauss-Seidel
 - Examples of GCR + Preconditioners for 3D FD

Finite Element Methods

HEAT FLOW 1-D EXAMPLE



Question: What is the temperature distribution along the bar



❖ Normalized Poisson Equation

$$\frac{\partial}{\partial x} \kappa \frac{\partial T(x)}{\partial x} = -h_s \Rightarrow -\frac{\partial^2 u(x)}{\partial x^2} = f(x)$$

$$-u_{xx}(x) = f(x)$$

407

THE FEM BASIS FUNCTIONS

- ❖ Residual Equation
- Partial Differential Equation Form

$$-\frac{\partial^2 u}{\partial x^2} = f \quad u(0) = 0 \quad u(1) = 0$$

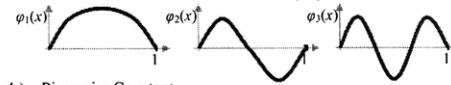
Step 1: Choose Basis Functions to represent the solution

$$u(x) \approx u_h(x) = \sum_{i=1}^n \omega_i \underbrace{\varphi_i(x)}_{\text{Basis Functions}}$$

$$-\frac{\partial^2 u}{\partial x^2} = f(x) \quad u(x) = \sum_{j=1}^n w_j \varphi_j(x) \quad u(0) = 0 \quad u(1) = 0$$

What are good basis functions?

- 1.) $\varphi(x) = a_0 + a_1 x \rightarrow$ No, second derivative is always zero
 - 2.) $\varphi(x) = a_0 + a_1 x + a_2 x^2 \rightarrow$ Could work
 - 3.) I can try to build into the basis function the boundary conditions
e.g. $\varphi_1(x) = \sin 2\pi i x$ [and second derivative is nonzero]
- Good choice. This is the same used by Spectral Methods. Fourier Series.



- 4.) Piecewise Constant.
Good. This leads to Finite Difference Methods



- 5.) Piecewise Linear.

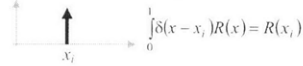


Step 2: Generate equations for the basis functions weights setting residual orthogonal to some test functions

$$\int_0^1 \underbrace{\varphi_j(x)}_{\text{functions}} R(x) dx = 0 \quad R(x) = \sum_{i=1}^n \omega_i \frac{d^2 \varphi_i(x)}{dx^2} + f(x)$$

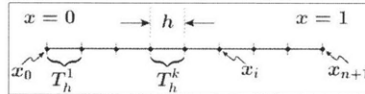
Test Functions $-\frac{\partial^2}{\partial x^2} \left(\sum_{i=1}^n w_i \varphi_i(x) \right) - f(x) = R(x) = 0$

- 1.) Collocation $R(x_i) = 0 \quad i = 1, \dots, n \quad \Psi_i(x) = \delta(x - x_i)$
- 2.) Galerkin $\varphi_i \perp R(x) \quad \int_0^1 \varphi_i(x) R(x) dx = 0 \quad \Psi_i(x) = \varphi_i(x)$
- 3.) Other Test Functions $\Psi_i \perp R(x) \quad \int_0^1 \Psi_i(x) R(x) dx = 0$



❖ Rayleigh-Ritz Approach – Approximation

○ Mesh

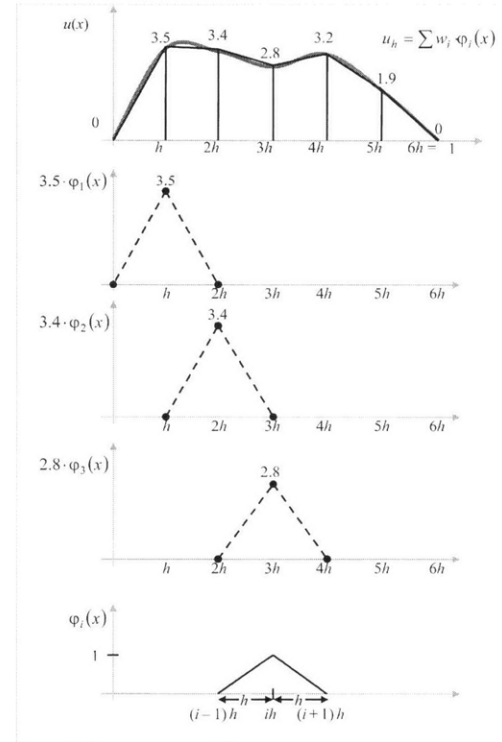
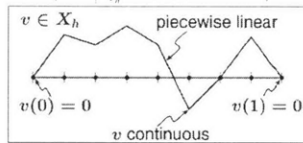


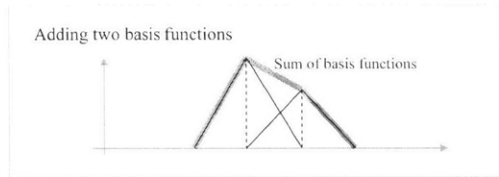
$$\Omega = \bigcup_{k=1}^K T_h^k \quad T_h^k, k = 1, \dots, K = n+1; \text{ elements}$$

$$x_j, j = 0, \dots, n+1; \text{ nodes}$$

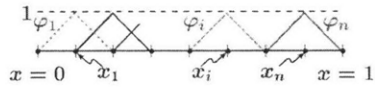
○ Space $X_h \subset X$

$$X_h = \{v \in X \mid v|_{T_h^k} \in \mathbb{P}_1(T_h^k), k = 1, \dots, K\}$$





- o Basis
- Nodal basis for X_h : $\varphi_j, j = 1, \dots, n = \dim(X_h)$



φ_j nonzero only on $T_h^j \cup T_h^{j+1}$

❖ Example Basis Functions

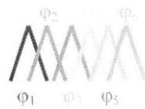
Introduce basis representation $u(x) \approx u_h(x) = \sum_{i=1}^n \omega_i \underbrace{\varphi_i(x)}_{\text{Basis Functions}}$

$\Rightarrow u_h(x)$ is a weighted sum of basis functions
The basis functions define a space

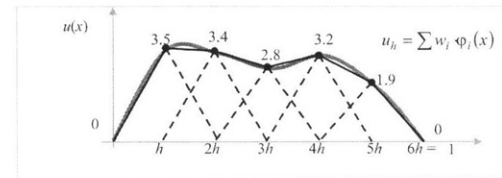
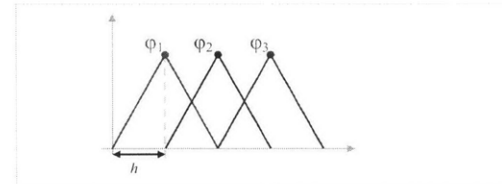
$$X_h = \left\{ v \in X_h \mid v = \sum_{i=1}^n \beta_i \varphi_i \text{ for some } \beta_i \text{'s} \right\}$$

Example

“Hat” basis functions



Piecewise linear Space



THE COLLOCATION SCHEME

Assume the initial condition is satisfied

$$\sum_{i=1}^n \omega_i \varphi_i(0) = x_0$$

Define the Residual

$$R(t) = \sum_{i=1}^n \omega_i \frac{d\varphi_i(t)}{dt} - A \sum_{i=1}^n \omega_i \varphi_i(t) - u(t)$$

Select weights to “minimize” the residual

Force the residual to be zero at n test points

$$R(t_l) = \sum_{i=1}^n \omega_i \frac{d\varphi_i(t_l)}{dt} - A \sum_{i=1}^n \omega_i \varphi_i(t_l) - u(t_l) = 0 \quad l \in \{1, \dots, n\}$$

Generates n equations in n unknowns

THE GALERKIN SCHEME

Force the residual to be “orthogonal” to the basis functions

$$\int_0^1 \varphi_j(x) R(x) dx = 0$$

Generates n equations in n unknowns

$$\int_0^1 \varphi_l(x) \left[\sum_{i=1}^n w_i \frac{d^2 \varphi_i(x)}{dx^2} + f(x) \right] dx = 0 \quad l \in \{1, \dots, n\}$$

$\frac{\partial^2 \phi_i}{\partial x^2} = ?$

delta function area = $\frac{1}{h}$

We don't want to work with these!!!

$$\int_0^1 \varphi_l(x) \left[\frac{\partial^2}{\partial x^2} \sum_{j=1}^n w_j \varphi_j(x) + f \right] dx = 0$$

$$\sum_{j=1}^n w_j \underbrace{\int_0^1 \varphi_l(x) \frac{\partial^2}{\partial x^2} \varphi_j(x) dx}_{A_{lj}} = - \int_0^1 \varphi_l(x) f(x) dx$$

$$\mathbf{Ax} = \mathbf{b}$$

$$\begin{bmatrix} A_{11} \\ \vdots \\ A_{ln} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} = \begin{bmatrix} \int_0^1 \varphi_1(x) f(x) dx \\ \vdots \\ \int_0^1 \varphi_n(x) f(x) dx \end{bmatrix}$$

Integrate by parts to avoid having to deal with second derivatives of hat basis functions.

$$\left[\varphi_l(x) \sum_{i=1}^n w_i \frac{d\varphi_i}{dx} \right]_0^1 - \int_0^1 \frac{d\varphi_l}{dx} \sum_{i=1}^n w_i \frac{d\varphi_i}{dx} dx + \int_0^1 \varphi_l(x) f(x) dx = 0 \quad l = \{1, 2, \dots, n\}$$

$$\underbrace{\varphi_l(1) \sum_{i=1}^n w_i \frac{d\varphi_i}{dx} \Big|_{x=1} - \varphi_l(0) \sum_{i=1}^n w_i \frac{d\varphi_i}{dx} \Big|_{x=0}}_{\text{Impose the Boundary Conditions}} - \int_0^1 \frac{d\varphi_l}{dx} \sum_{i=1}^n w_i \frac{d\varphi_i}{dx} dx + \int_0^1 \varphi_l(x) f(x) dx = 0$$

Note: Do NOT need 2nd derivatives $l = \{1, 2, \dots, n\}$

Integration by Parts

$$\begin{aligned} \sum w_i \int_0^1 \varphi_i(x) \frac{\partial^2 \varphi_i(x)}{\partial x^2} dx &= \sum w_i \left[\varphi_i(x) \frac{\partial \varphi_i(x)}{\partial x} \right]_0^1 - \sum w_i \int_0^1 \frac{\partial \varphi_i(x)}{\partial x} \frac{\partial \varphi_i(x)}{\partial x} dx \\ &= \sum w_i \left[\varphi_i(1) \frac{\partial \varphi_i(1)}{\partial x} \right] - \sum w_i \left[\varphi_i(0) \frac{\partial \varphi_i(0)}{\partial x} \right] \\ &\quad - \sum w_i \int_0^1 \frac{\partial \varphi_i(x)}{\partial x} \frac{\partial \varphi_i(x)}{\partial x} dx \\ &= - \sum w_i \int_0^1 \frac{\partial \varphi_i(x)}{\partial x} \frac{\partial \varphi_i(x)}{\partial x} dx \end{aligned}$$

$\mathbf{Ax} = \mathbf{b}$

$$\left[\begin{array}{c} A_{1j} \\ \vdots \\ A_{ij} \\ \vdots \\ A_{nj} \end{array} \right] = - \int_0^1 \frac{\partial \varphi_i(x)}{\partial x} \frac{\partial \varphi_j(x)}{\partial x} dx \quad \left[\begin{array}{c} w_1 \\ \vdots \\ w_i \\ \vdots \\ w_n \end{array} \right] = \left[\begin{array}{c} \int_0^1 \varphi_1(x) f(x) dx \\ \vdots \\ \int_0^1 \varphi_i(x) f(x) dx \\ \vdots \\ \int_0^1 \varphi_n(x) f(x) dx \end{array} \right]$$

Only first derivatives of basis functions

$$\int_0^1 \frac{d\varphi_l(x)}{dx} \frac{d}{dx} \sum_{i=1}^n \omega_i \varphi_i(x) dx - \int_0^1 \varphi_l(x) f(x) dx = 0 \quad l \in \{1, \dots, n\}$$

$$\begin{aligned} \sum_i \omega_i \int_0^1 \frac{d\varphi_i}{dx} \cdot \frac{d\varphi_i}{dx} dx &= \int_0^1 \varphi_1(x) f(x) dx \\ \omega_1 \frac{1}{(\Delta x)^2} [2 \cdot \Delta x] + \omega_2 \frac{1}{(\Delta x)^2} [-1 \cdot \Delta x] &= \int_0^1 \varphi_1(x) f(x) dx \\ \omega_1 \frac{1}{(\Delta x)^2} [-1 \cdot \Delta x] + \omega_2 \frac{1}{(\Delta x)^2} [2 \cdot \Delta x] - \omega_3 \frac{1}{(\Delta x)^2} [-1 \cdot \Delta x] &= \int_0^1 \varphi_2(x) f(x) dx \\ \frac{1}{\Delta x} \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & -1 & 2 \\ & & & & & -1 & 2 \end{bmatrix} \begin{bmatrix} \omega_1 \\ \omega_2 \\ \vdots \\ \omega_n \end{bmatrix} &= \begin{bmatrix} \int \varphi_1 f(x) dx \\ \int \varphi_2 f(x) dx \\ \vdots \\ \int \varphi_n f(x) dx \end{bmatrix} \end{aligned}$$

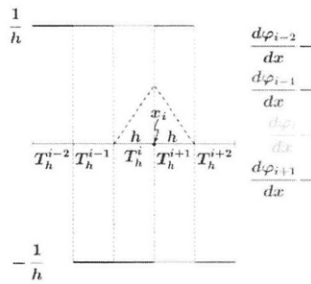
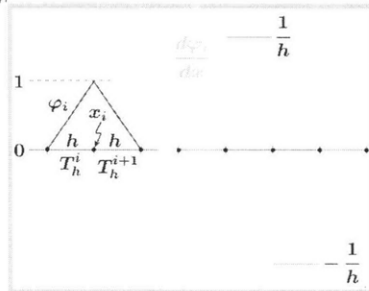
THE FEM LINEAR SYSTEM

$$\sum_{i=1}^n \omega_i \underbrace{\int_0^1 \frac{d\varphi_i}{dx} \frac{d\varphi_i}{dx} dx}_{A_{i,j}} = \underbrace{\int_0^1 \varphi_i(x) f(x) dx}_{F_i}$$

$$\left[\begin{array}{c} \mathbf{A} \\ \vdots \\ \mathbf{A} \end{array} \right] \begin{bmatrix} \omega_1 \\ \omega_2 \\ \vdots \\ \omega_n \end{bmatrix} = \begin{bmatrix} F_1 \\ F_2 \\ \vdots \\ F_n \end{bmatrix}$$

DISCRETE EQUATIONS

❖ ϕ_i



❖ Typical Row

$$A_{h,i,j} = \int_{\Omega} \frac{d\phi_i}{dx} \frac{d\phi_j}{dx} dx = \int_{T_h^{i-1}} \frac{d\phi_i}{dx} \frac{d\phi_j}{dx} dx + \int_{T_h^i} \frac{d\phi_i}{dx} \frac{d\phi_j}{dx} dx$$

is nonzero only for $i = j-1, j, j+1$

$$A_{h,i,i} = \frac{1}{h^2}(h) + \frac{1}{h^2}(h) = \frac{2}{h}$$

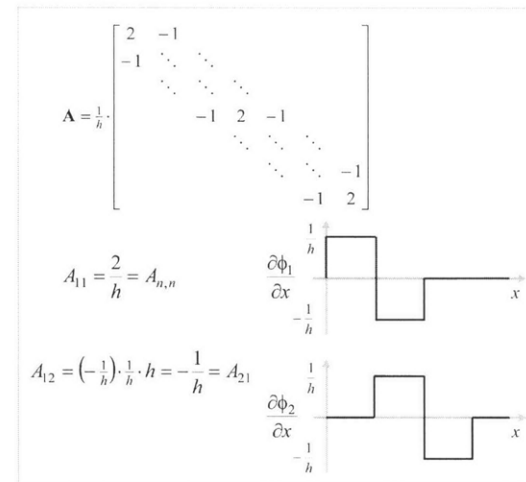
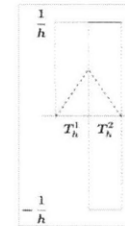
$$A_{h,i,i-1} = \frac{1}{h} \left(-\frac{1}{h} \right) (h) = -\frac{1}{h}$$

$$A_{h,i,i} = \left(-\frac{1}{h} \right) \frac{1}{h} (h) = -\frac{1}{h}$$

❖ Boundary Rows

$$A_{h,1,1} = \frac{2}{h} \quad A_{h,1,2} = -\frac{1}{h}$$

$$A_{h,n,n} = \frac{2}{h} \quad A_{h,n,n-1} = -\frac{1}{h}$$



F.D. F.E.M.

$$\Delta x \cdot f(x_j) \approx \int_0^1 \phi_j(x) f(x) dx$$

If $f(x)$ is constant:

$$\int_0^1 \phi_j(x) f(x) dx = \frac{1}{2} \cdot 2 \cdot \Delta x \cdot f(x_j)$$

→ same as F.D.

If $f(x)$ is not constant the RHS is still an approximation

$$\int_0^1 \phi_j(x) f(x) dx \approx \frac{1}{2} \Delta x \cdot f(x_j)$$

F.E.M.

$$\frac{1}{\Delta x} \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & & & \\ & & \ddots & & \\ & & & -1 & 2 \\ & & & & -1 & 2 \end{bmatrix} \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix} = \Delta x \begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix}$$

Same as F.D. system for 1D hat basis
However F.E.M. gives more flexibility (can choose other basis)

F.E.M. solution vector $\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$

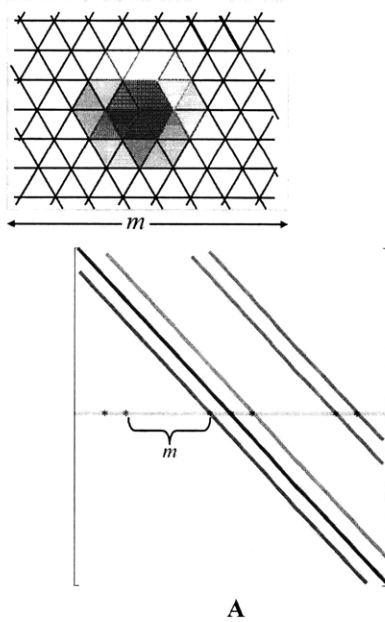
F.D. solution vector $\begin{bmatrix} \hat{u}_1 \\ \hat{u}_2 \\ \vdots \\ \hat{u}_n \end{bmatrix}$

AN FEM EXAMPLE

1-D

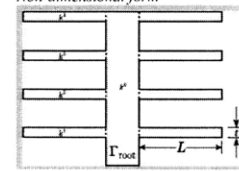
2-D

Top View

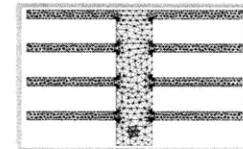


Heat Transfer Problem

Non-dimensional form



k^i : thermal conductivity for $\Omega_i, i = 0, 1, \dots, 4$
 B_i : heat transfer coefficient
 t, L : geometric parameters



$X_h = \text{span}\{\phi_1, \dots, \phi_n\}$
 $\phi_i(x)$: Nodal basis functions
 - first order elements $\dim(X_h) = n$

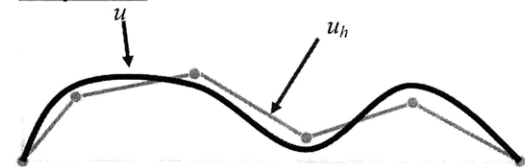
Possible Solutions



CONVERGENCE ANALYSIS OF FEM

❖ Convergence Analysis

The question is



How does $\underbrace{\|u - u_h\|}_{\text{error}}$ decrease with refinement?

❖ Review of FEM

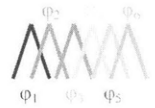
Introduce basis representation $u(x) \approx u_h(x) = \sum_{i=1}^n \beta_i \underbrace{\varphi_i(x)}_{\substack{\text{Basis} \\ \text{Functions}}}$

$\Rightarrow u_h(x)$ is a weighted sum of basis functions
The basis functions define a space

$$X_h = \left\{ v \in X_h \mid v = \sum_{i=1}^n \beta_i \varphi_i \text{ for some } \beta_i \text{'s} \right\}$$

Example

“Hat” basis functions



Piecewise linear Space



Partial Differential Equation Form

$$-\frac{\partial^2 u}{\partial x^2} = f \quad u(0) = 0 \quad u(1) = 0$$

“Nearly” Equivalent Weak Form

$$\int_{\Omega} \frac{\partial u}{\partial x} \frac{\partial v}{\partial x} dx = \int_{\Omega} f v dx \quad \text{for all } v$$

Introduced an abstract notation for the equation u must satisfy
 $a(u, v) = l(v)$ for all v

$a(u, u)$ defines a norm $a(u, u) = \|u\|^2$

u is restricted to be 0 at 0 and 1!!!

Key Idea:

Using the norm properties, it is possible to show

If $a(u_h, \varphi_i) = l(\varphi_i)$ for all $\varphi_i \in \{\varphi_1, \varphi_2, \dots, \varphi_n\}$

Note: this is basically Galerkin

$$\text{Then } \underbrace{\|u - u_h\|}_{\text{Error}} = \min_{v_h \in X_h} \underbrace{\|u - v_h\|}_{\text{Projection Error}}$$

Or in other words: given some basics, Galerkin finds the best coefficients

417

$u(x) \in \mathcal{R}^\infty$ Here $u(x) \in \mathcal{R}^3$ for visualization

Pick some basis functions e.g. $\varphi_1(x), \varphi_2(x)$
 $\|u(x) - u_h(x)\|$

The space of all their linear combinations $X_h = \text{span}\{\varphi_1(x), \varphi_2(x)\}$ is a plane
Here, for visualization, assume the $\{\hat{i}, \hat{j}\}$ plane

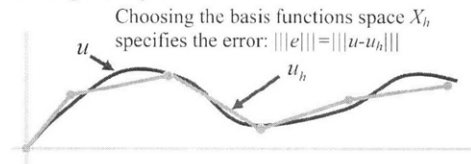
The key idea says that the error $\|u(x) - u_h(x)\| = \|e_h\| = \sqrt{a(e, e)} = \sqrt{a(u(x) - u_h(x), u(x) - u_h(x))}$
 $= \left[\int_0^1 \frac{\partial(u - u_h)}{\partial x} \frac{\partial(u - u_h)}{\partial x} dx \right]^{1/2}$ in 1D example

is minimized by choosing $u_h(x)$ such that
 $\begin{cases} a(u_h, \varphi_1) - l(\varphi_1) = 0 \\ a(u_h, \varphi_2) - l(\varphi_2) = 0 \end{cases}$ the solution is obtained by projection

\Rightarrow for a given set of basis, Galerkin gives the projection and hence gives the smallest possible error (best coefficients)
 \Rightarrow Note: to do better than that I need to change X_h the set of basis

For instance hat basis
As I decrease $h = \frac{1}{n}$ the error decreases linearly
 $\|e_h\| = \|u - u_h\| = \mathcal{O}\left(\frac{1}{n}\right) \xrightarrow{n \rightarrow \infty} 0$

- ❖ Error Analysis – Energy Norm
In words: even if you *knew* u , you could not find a w_h in X_h more accurate than u_h *in the energy norm.*
- ❖ FEM Convergence Analysis

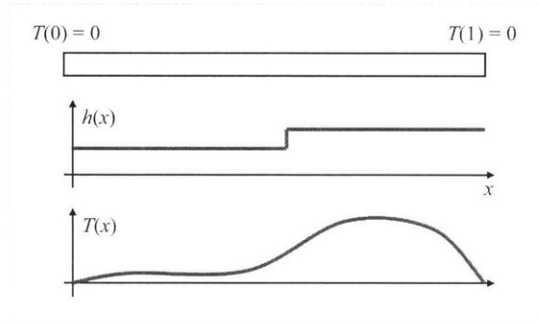


Question: How does that error change when I change basis?

For piecewise linear: $\|u - u_h\| = O\left(\frac{1}{n}\right)$

ENERGY MINIMIZATION VIEW OF FEM

- ❖ Problem of Interest
 - Helmholtz Equation in 1D
Boundary Value Problem (BVP) – Strong Form
 $-u''(x) + \alpha u(x) = f(x) \quad \alpha \geq 0$
 $x \in (0,1) \quad u(0) = u(1) = 0$
 Describes many physical phenomena (e.g.):
 - Temperature distribution in a bar
 - Deformation of an elastic bar
 - Deformation of a string under tension
 - Solution Properties
 - The solution $u(x)$ always *exists*
 - $u(x)$ is always smoother than the data $f(x)$
 - given $f(x)$ the solution $u(x)$ is *unique*



- ❖ Minimization Principle Statement

Find

$$u = \arg \min_{w \in X} J(w)$$

where

$$X = \{v \text{ sufficiently smooth } v(0) = v(1) = 0\}$$

and

$$J(w) = \frac{1}{2} \int_0^1 (w_x w_x + \alpha w w) dx - \int_0^1 f w dx$$

$$J_{\text{Poisson}}(w) = \frac{1}{2} \int_0^1 w_x w_x dx - \int_0^1 f w dx$$

$$\int_0^1 \underbrace{R(x)}_{0 \text{ or } -f(x)} \varphi dx = \int_0^1 \frac{\partial \varphi}{\partial x} \frac{\partial \varphi}{\partial x} dx - \int_0^1 f \varphi dx$$

$$\int_0^1 (u'' - f(x))v = 0 \quad \forall v$$

In words:

Over all functions w in X ,

u that satisfies

$$-u_{xx} + \alpha u = f \text{ in } \Omega$$

$$u(0) = u(1) = 0$$

makes $J(w)$ as small as possible

o Proof

Let $w = u + v$

Then

$$J\left(\frac{u+v}{e^{\frac{w}{X}}}\right) = \frac{1}{2} \int_0^1 (u+v)_x (u+v)_x dx + \frac{\alpha}{2} \int_0^1 (u+v)(u+v) dx - \int_0^1 f(u+v) dx$$

$$J(u+v) = \frac{1}{2} \int_0^1 (u_x u_x + \alpha u u) dx - \int_0^1 f u dx = J(u)$$

$$+ \int_0^1 (u_x v_x + \alpha u v) dx - \int_0^1 f v dx \quad \delta J_v(u)$$

$$+ \frac{1}{2} \int_0^1 (v_x v_x + \alpha v v) dx > 0 \text{ for } v \neq 0$$

$$\delta J_{\text{Poisson}}(w) = \int_0^1 u_x v_x dx - \int_0^1 f v dx$$

If $v_x = \varphi_i$ then Galerkin

$$\delta J_v(u) = \int_0^1 (u_x v_x + \alpha u v) dx - \int_0^1 f v dx$$

$$= v(0)u_x(0) - v(1)u_x(1) - \int_0^1 u_{xx} v dx + \alpha \int_0^1 u v dx - \int_0^1 f v dx$$

$$J(u+v) = J(u) + \underbrace{\int_0^1 (v_x v_x + \alpha v v) dx}_{>0 \text{ unless } v=0}, \quad \forall v \in X$$

$$\Rightarrow \begin{cases} J(w) > J(u), \quad \forall w \in X, w \neq u \\ \Downarrow \\ u \text{ is THE minimizer of } J(w) \end{cases}$$

❖ Rayleigh-Ritz Approach – “Projection”

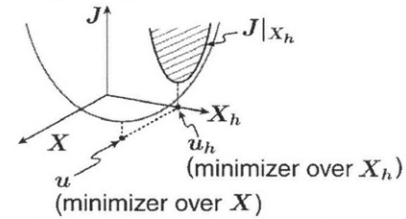
Let

$$u_h \in X_h = \sum_{j=1}^n u_{hj} \varphi_j(x)$$

RR-FE approximation

Set $u_{hj} = w_j$ that minimize

$$J\left(\sum_{j=1}^n w_j \varphi_j\right)$$



$$J\left(\sum_{j=1}^n w_j \varphi_j\right) = \frac{1}{2} \int_0^1 \frac{d}{dx} \left(\sum_{i=1}^n w_i \varphi_i \right) \frac{d}{dx} \left(\sum_{j=1}^n w_j \varphi_j \right) dx + \dots$$

$$\dots + \frac{\alpha}{2} \int_0^1 \sum_{i=1}^n (w_i \varphi_i) \sum_{j=1}^n (w_j \varphi_j) dx - \int_0^1 f \sum_{j=1}^n w_j \varphi_j dx$$

$$= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n w_i w_j \int_0^1 \left(\frac{d\varphi_i}{dx} \frac{d\varphi_j}{dx} + \alpha \varphi_i \varphi_j \right) dx - \sum_{j=1}^n w_j \int_0^1 f \varphi_j dx$$

By bilinearity and linearity.

$$\begin{aligned}
 u_h(x) &= \sum w_j \phi_j(x) \\
 -u_{xx} + \alpha u - f(x) &= 0 \\
 -\frac{\partial^2}{\partial x^2} (\underbrace{\sum w_j \phi_j(x)}_{R(x)}) + \alpha \sum w_j \phi_j(x) - f(x) &\equiv 0 \\
 \int_0^1 \phi_l(x) R(x) dx &= 0 \quad \forall l \\
 -\int_0^1 \phi_l(x) \sum_{i=1}^n w_i \frac{\partial^2}{\partial x^2} \phi_i(x) dx + \int_0^1 \phi_l(x) \alpha \sum_{i=1}^n w_i \phi_i(x) dx - \int_0^1 \phi_l(x) f(x) dx &= 0
 \end{aligned}$$

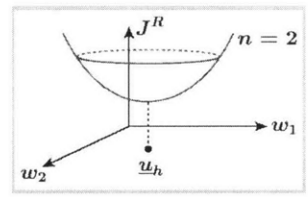
$$\begin{aligned}
 J^R(\underline{w} \in \mathfrak{R}^n) &\equiv J\left(\sum_{j=1}^n w_j \phi_j\right) \\
 &= \frac{1}{2} \underline{w}^T \underline{A}_h \underline{w} - \underline{w}^T \underline{F}_h
 \end{aligned}$$

$$\underline{F}_h \in \mathfrak{R}^n : F_{hi} = \int_0^1 f \phi_i dx$$

$$\underline{A}_h \in \mathfrak{R}^{n \times n} : A_{hij} = \int_0^1 \left(\frac{d\phi_i}{dx} \frac{d\phi_j}{dx} + \alpha \phi_i \phi_j \right) dx$$

$$\begin{aligned}
 J(\underline{\tilde{w}}) &= \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}^T \left[\dots \int_0^1 \left(\frac{\partial \phi_i}{\partial x} \frac{\partial \phi_j}{\partial x} + \alpha \phi_i \phi_j \right) dx \dots \right] \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}^T \begin{bmatrix} F_1 \\ \vdots \\ F_n \end{bmatrix} \\
 \text{For Poisson} & \\
 J(\underline{\tilde{w}}) &= \frac{1}{2} \sum \sum w_i w_j \int_0^1 \frac{\partial \phi_i}{\partial x} \frac{\partial \phi_j}{\partial x} dx - \sum w_j \int f \phi_j dx \\
 &= \frac{1}{2} \underline{\tilde{w}}^T \underline{A} \underline{\tilde{w}} - \underline{\tilde{w}}^T \underline{\tilde{F}}_h \\
 \nabla_{\underline{w}} J(\underline{\tilde{w}}) &= 0 \quad \text{to minimize} \\
 &\Rightarrow \underline{A} \underline{\tilde{w}} - \underline{\tilde{F}}_h = 0
 \end{aligned}$$

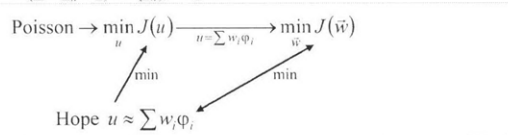
If **A** is SPD $\min_{\underline{\tilde{w}}} \frac{1}{2} \underline{\tilde{w}}^T \underline{A} \underline{\tilde{w}} - \underline{\tilde{w}}^T \underline{\tilde{F}}_h$ $\underline{\tilde{w}}$ satisfies $\underline{A} \underline{\tilde{w}} = \underline{\tilde{F}}_h$



$u_h = \arg \min_{\underline{w} \in \mathfrak{R}^n} J^R(\underline{w})$
 Expand $J(\underline{w} = \underline{u}_h + \underline{v})$;
 require $J(\underline{w}) > J(\underline{u}_h)$
 unless $\underline{v} = 0$.

$$\begin{aligned}
 J^R(u_h + v) &= \frac{1}{2}(u_h + v)^T A_h(u_h + v) - (u_h + v)^T F_h \\
 &= \frac{1}{2}u_h^T A_h u_h - u_h^T F_h + \frac{1}{2}v^T A_h v + \frac{1}{2}u_h^T A_h v - v^T F_h + \frac{1}{2}v^T A_h v \\
 J^R(u_h + v) &= J(u) \\
 &\quad + \underbrace{(A_h u_h - F_h)^T}_{\nabla J^R(u_h)} v \underbrace{\delta J^R(u_h)}_{\text{SPD}} \\
 &\quad + \frac{1}{2}v^T \underbrace{A_h}_{>0, \forall v \neq 0} v \quad \text{SPD}
 \end{aligned}$$

If (and only if)
 $\delta J^R_v(u_h) = 0, \forall v \in \mathfrak{R}^n$
 \Downarrow
 $\nabla J^R(u_h) = A_h u_h - F_h = 0$
 then
 $J(w = u_h + v) > J(u_h), \forall v \neq 0$

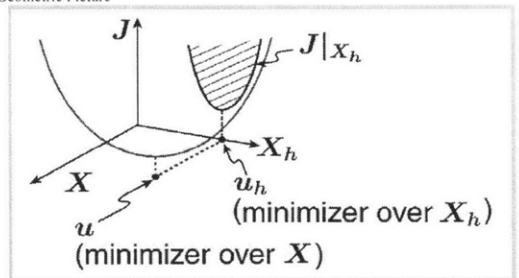


Find $u_h \in \mathfrak{R}^n$ such that
 $A_h u_h = F_h \Rightarrow u_h(x) = \sum_{j=1}^n u_{hj} \phi_j(x)$
 SPD \Leftrightarrow existence and uniqueness

❖ Error Analysis - Energy Norm
 $J(u_h) < J(w_h), \forall w_h \in X_h, w_h \neq u_h$
 If $e = u - u_h$
 $\| \underbrace{u - u_h}_e \| < \| u - w_h \|, \forall w_h \in X_h, w_h \neq u_h$
 and
 $\| e \| = \inf_{w_h \in X_h} \| u - w_h \|$

where the energy norm is defined as: $\| e \| = \sqrt{\int_0^1 (e_x^2 + a e^2) dx}$

❖ Geometric Picture



GCR Preconditioners

$Ax = b$
 $PAx = Pb$ preconditioner
 $P \approx A^{-1}$ Choose \tilde{A} s.t.

- 1.) $\tilde{A} \approx A$ e.g. $\tilde{A} = \text{diag}(A) = D$
- 2.) $P = \tilde{A}^{-1}$ e.g. $\tilde{A} = \text{tridiag}(A)$

easy to solve
 $\tilde{A}^{-1} Ax = \tilde{A}^{-1} b$
 $\tilde{A}^{-1} b = \tilde{b} \Leftrightarrow b = \tilde{A} \tilde{b}$
 $\tilde{A}^{-1} Ax = \tilde{A}^{-1} b = \tilde{b}$
 $\tilde{A}^{-1} A r_{k-1} = r_k$ Solve $\tilde{A} r_k = A r_{k-1}$

$P = A^{-1} \Leftrightarrow A^{-1} A x = A^{-1} b \Leftrightarrow I x = A^{-1} b$ ← GCR converges in 1 iteration

DIAGONAL PRECONDITIONERS

Let $A = D + A_{nd}$

Apply GCR to $(D^{-1}A)\tilde{x} = (I + D^{-1}A_{nd})\tilde{x} = D^{-1}b$

- The inverse of a diagonal is cheap to compute
- Usually improves convergence

Krylov Methods – Convergence Analysis

GCR Optimality Property:

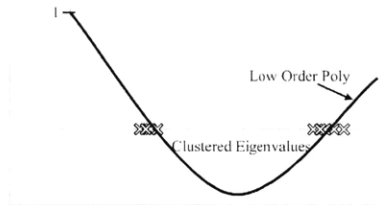
$$\|r^k\| = \|v_k(M)r^0\| \leq |\tilde{\varphi}_k(M)r^0| \leq |\tilde{\varphi}_k(M)| \|r^0\|$$

$\tilde{\varphi}_k$ is any k^{th} order poly such that $\tilde{\varphi}_k(0) = 1$

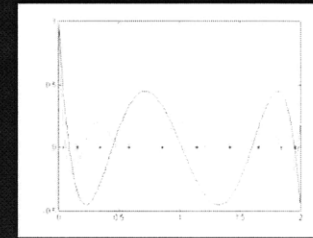
Therefore

Any polynomial which satisfies the constraint can be used to get

an upper bound on $\|r^k\| \leq |\tilde{\varphi}_k(M)|$

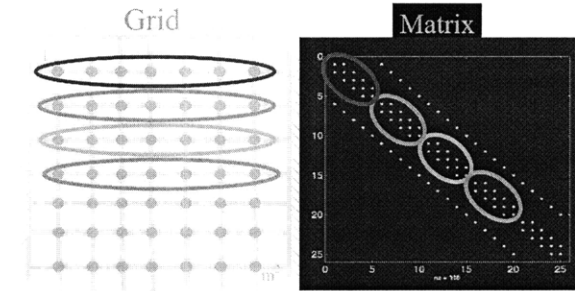


**Residual Poly Picture for Heat Conducting Bar Matrix
No loss to air (n=10)**



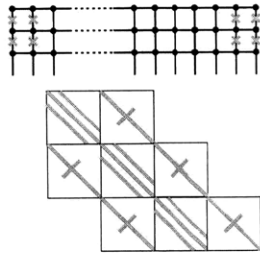
Keep $|\varphi_k(\lambda_i)|$ as small as possible:
Easier if eigenvalues are clustered!

BLOCK DIAGONAL PRECONDITIONERS



Tridiagonal Matrices factor quickly

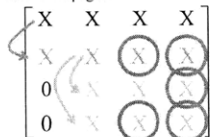
Block Diagonal Preconditioners



$\tilde{A}_B = A$ (blocks)
 \tilde{A}_B is a different problem
 (no vertical interaction)
 \tilde{A}_B is easy to solve (tridiagonal)
 $P = \tilde{A}_B^{-1}$

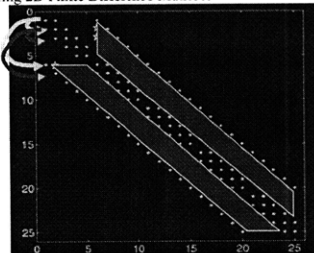
INCOMPLETE FACTORIZATION PRECONDITIONERS

- ❖ Sparse Matrices – Fill-in
 Fill-ins Propagate



Fill-ins from Step 1 results in
 Fill-ins in Step 2

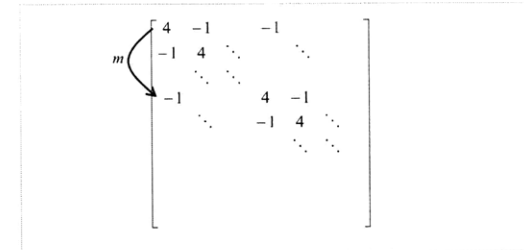
- ❖ Factoring 2D Finite Difference Matrices



Generated Fill-in
 Makes
 Factorization
 Expensive

- ❖ Incomplete Factorization Schemes – Key Idea
 Throw away all fill-ins
 Throw away only fill-ins with small values

Throw away fill-ins produced by other fill-ins
 Throw away fill-ins produced by fill-ins of other fill-ins, etc.



$\tilde{A}_L = L + D$

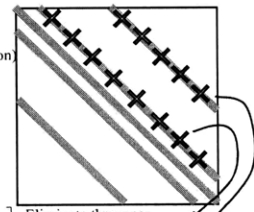
Lower triangular plus diagonal

\tilde{A}_L is easy to solve (Forward Elimination)

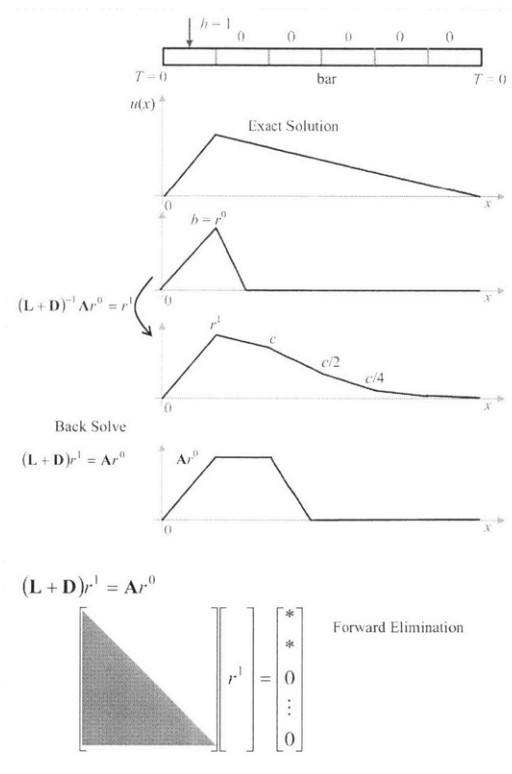
$P = \tilde{A}_L^{-1} = (L + D)^{-1}$

$A = \begin{bmatrix} 2 & -1 & & \\ -1 & 2 & -1 & \\ & -1 & 2 & \\ & & & \ddots \end{bmatrix} = L + D + U$

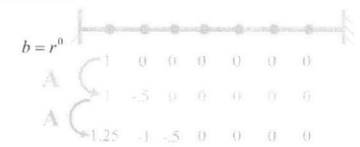
$\rightarrow L + D = \begin{bmatrix} 2 & & & \\ -1 & 2 & & \\ & -1 & 2 & \\ & & & \ddots \end{bmatrix}$



Eliminate the upper
 triangular part of the matrix



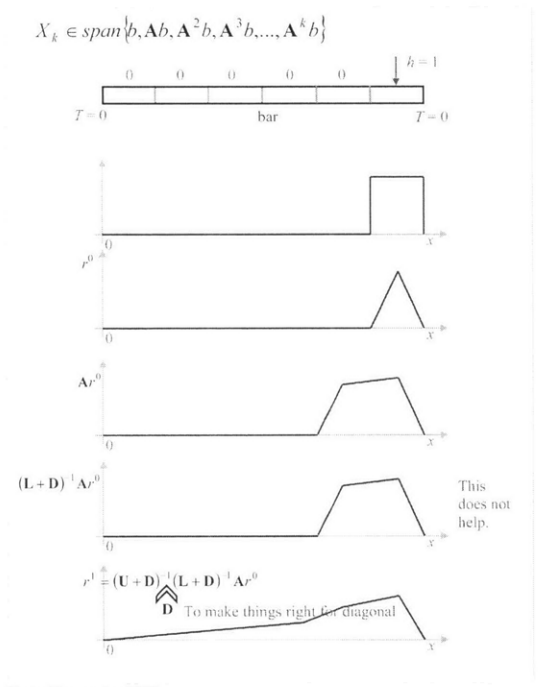
COMMUNICATION IMPROVING PRECONDITIONERS: GAUSS-SEIDEL
 ❖ Lower Bound



Communication Lower Bound for m gridpoints

$A^k r^0$ is nonzero in m^{th} entry after $k = m$ iterations

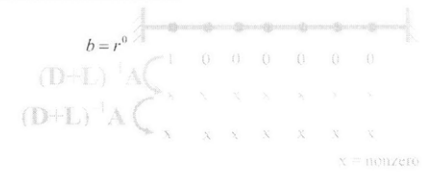
Need at least m iterations for $(x^{k+1})_m = \left(x^0 + \sum_{j=1}^k a_j r^j \right)_m \neq 0$



- ❖ Faster Converging GCR
 - GCR already achieves Communication Lower bound
 - Preconditioner must accelerate communication
 - Multiplying by Preconditioner must move values by more than one grid point.

- ❖ Gauss-Seidel Preconditioning
 - o Krylov Vectors

1-D Discretized PDE



$Ax = b$

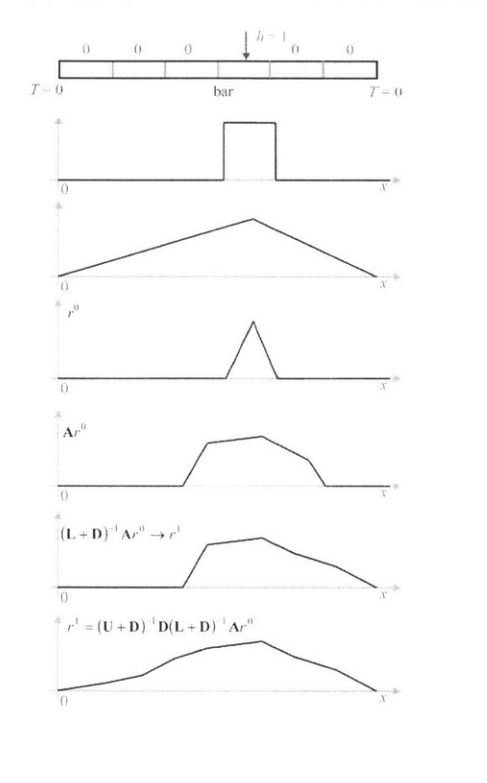
$(L+D)^{-1} Ax = (L+D)^{-1} b$

$A = L + D + U$

1.) $\tilde{A} \approx A = L + D + U$

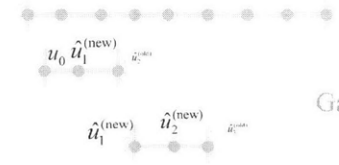
$\tilde{A} = L + D \quad \tilde{A} = \begin{bmatrix} 2 & & & & \\ -1 & 2 & & & \\ & \ddots & \ddots & \ddots & \\ & & & -1 & 2 \end{bmatrix}$

2.) Easy to solve for $\tilde{A} = L + D$
 Forward Elimination



Physical View

1-D Discretized PDE



Gauss Seidel

Each iteration of Gauss-Seidel Relaxation moves data across the grid

$Ax = [L + D + U]x = b$
 $(L + D)^{-1}Ax = (L + D)^{-1}b$

Lower Triangular
 Apply with back solve

$Ax = b$

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

Lower Triangular
 Apply with back solve

$a_{11}u_1^{new} + a_{12}u_2^{old} + a_{13}u_3^{old} = b_1$
 $a_{21}u_1^{new} + a_{22}u_2^{new} + a_{23}u_3^{old} = b_2$
 $a_{31}u_1^{new} + a_{32}u_2^{new} + a_{33}u_3^{new} = b_3$

$\Rightarrow (L + D)u^{new} + Uu^{old} = b$

$(L + D)^{-1}Ax = (L + D)^{-1}b$

o Krylov Vectors

1-D Discretized PDE

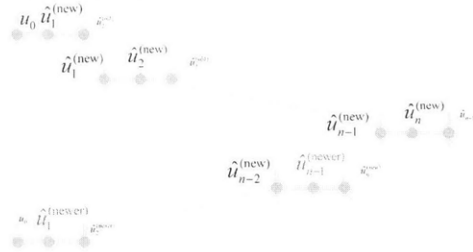


Gauss-Seidel communicates quickly in only one direction.

$$\frac{(\mathbf{L} + \mathbf{D})^{-1} \mathbf{A} x = (\mathbf{L} + \mathbf{D})^{-1} b}{M}$$

$$\frac{p}{(\mathbf{U} + \mathbf{D})^{-1} \mathbf{D} (\mathbf{L} + \mathbf{D})^{-1} \mathbf{A} x = (\mathbf{U} + \mathbf{D})^{-1} \mathbf{D} (\mathbf{L} + \mathbf{D})^{-1} b}$$

o Symmetric Gauss-Seidel



This symmetric Gauss-Seidel preconditioner communicates both directions.

Derivation of the SGS Iteration Equations

Forward Sweep (half step): $(\mathbf{D} + \mathbf{L})x^{k+1/2} + \mathbf{U}x^k = b$
 Backward Sweep (half step): $(\mathbf{D} + \mathbf{U})x^k + \mathbf{L}x^{k+1/2} = b$

$$\Rightarrow x^{k+1} = (\mathbf{D} + \mathbf{U})^{-1} \mathbf{L} (\mathbf{D} + \mathbf{L})^{-1} \mathbf{U} x^k + (\mathbf{D} + \mathbf{U})^{-1} b \dots$$

$$- (\mathbf{D} + \mathbf{U})^{-1} \mathbf{L} (\mathbf{D} + \mathbf{L})^{-1} b$$

$$\Rightarrow x^{k+1} = x^k - (\mathbf{D} + \mathbf{U})^{-1} \mathbf{D} (\mathbf{D} + \mathbf{L})^{-1} \mathbf{A} x^k + (\mathbf{D} + \mathbf{U})^{-1} \mathbf{D} (\mathbf{D} + \mathbf{L})^{-1} b$$

Step 1.) Pick $\tilde{\mathbf{A}}$ s.t.
 a.) $\tilde{\mathbf{A}} \approx \mathbf{A}$
 b.) $\tilde{\mathbf{A}}$ easy to factor (or solve) $\tilde{\mathbf{A}} = \tilde{\mathbf{L}}\tilde{\mathbf{U}}$
 $\Leftrightarrow \mathbf{P} = \tilde{\mathbf{A}}^{-1}$
 e.g. incomplete factorization

Step 2.) $\mathbf{A}x = b$ becomes $\tilde{\mathbf{A}}^{-1} \mathbf{A} x = \tilde{\mathbf{A}}^{-1} b$
 Solve $\tilde{\mathbf{A}} \tilde{b} = b$ use $\tilde{\mathbf{L}}\tilde{\mathbf{U}}\tilde{b} = b$

Step 3.) Do not calculate $\tilde{\mathbf{A}}^{-1} \mathbf{A}$
 Instead at each iteration of GCR method
 $[\tilde{\mathbf{A}}^{-1} \mathbf{A}]_k \rightarrow r_{k+1}$
 Set-up system $\tilde{\mathbf{A}} r_{k+1} = \mathbf{A} r_k$ Use $\tilde{\mathbf{L}}\tilde{\mathbf{U}} r_{k+1} = \mathbf{A} r_k$
 get r_{k+1}

EXAMPLES OF GCR + PRECONDITIONERS FOR 3D FD

- Complexity of Sparse GCR + Preconditioner for 3D FD Methods

	matrix size	dense block	
1-D	$n = m$	$k = 5-10$	domain size = m
2-D	$n = m^2$	$k = 5-10$	
3-D	$n = m^3$	$k = 5-10$	
1-D	$O(kn) = 10 O(m)$	\leftarrow 100 pt grid	$O(10^3)$ ops 1 μ s
2-D	$O(kn) = 10 O(m^2)$	\leftarrow 100x100 grid	$O(10^5)$ ops 0.1 ms
3-D	$O(kn) = 10 O(m^3)$	\leftarrow 100x100x100 grid	$O(10^7)$ ops 10 ms

 e.g. on a 1 GFlops computer using $m = 100$ and a constant $c = 10$
- Krylov - Work for Banded Gaussian Elimination, Relaxation, and GCR

Dimension	Banded GE	Sparse GE	GCR	GCR + preconditioner
1	$O(m)$	$O(m)$	$O(m^2)$	$O(km)$
2	$O(m^2)$	$O(m^2)$	$O(m^3)$	$O(km^2)$
3	$O(m^3)$	$O(m^3)$	$O(m^4)$	$O(km^3)$

 GCR+good preconditioner is always the best choice for 2D & 3D problems!

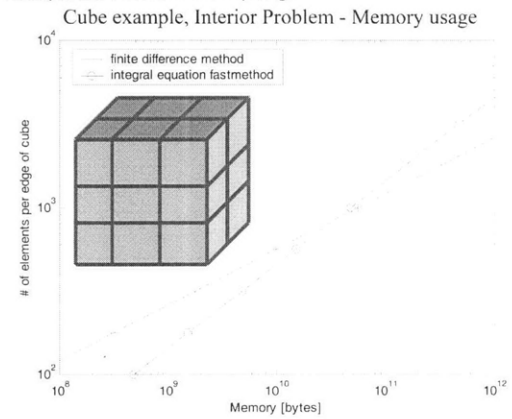
Dimension	Banded GE	Sparse GE	GCR	GCR + preconditioner
1	1 μ s	1 μ s	0.1 ms	1 μ s
2	1 sec	10 ms	10 ms	0.1 sec
3	11d 6h	2h 46m	1 sec	10 ms

e.g. on a 1 GFlops computer using $m=1000$, $k=10$, and constant $c=10$
 GCR+good preconditioner is always the best choice for 2D & 3D problems!

Dimension	Banded GE	Sparse GE	GCR	GCR + preconditioner
1	10 μ s	10 μ s	10 ms	0.1 ms
2	2h 46m	10 sec	10 sec	0.1 sec
3	317,000 years	317 years	2h 46m	1m 20s

e.g. on a 1 GFlops computer using $m=1000$, $k=10$, and constant $c=10$
 GCR+good preconditioner (i.e. $k=5-10$) is always the best choice for 2D & 3D problems!

❖ Cube Example, Interior Problem – Memory Usage

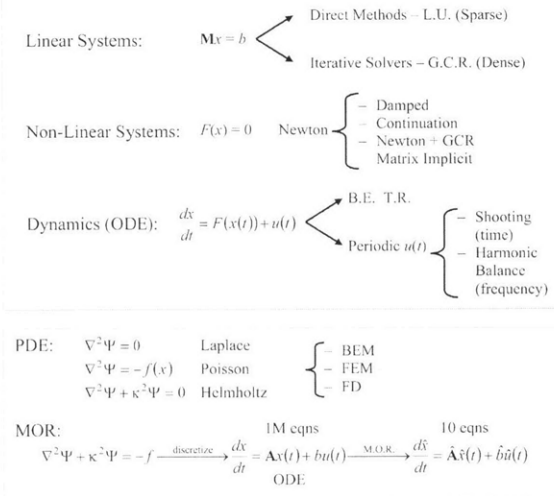


INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 22.
PDE – BEM Integral Equation Method I

- TODAY'S OUTLINE:**
- ❖ Poisson and Laplace Equations: examples
 - ❖ Exterior Versus Interior Problems
 - ❖ Green's Function
 - ❖ Basis Functions
 - ❖ Collocation Method
 - ❖ Galerkin Method

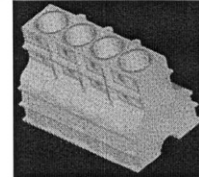
OVERVIEW OF COURSE



429

POISSON AND LAPLACE EQUATIONS: EXAMPLES

- ❖ Thermal Conduction Equations



www.adina.com

- o Thermal Conduction Equations
 - The Poisson Partial Differential Equation.

$$\nabla^2 \Psi(x) = \nabla \cdot \nabla \Psi(x) = -f(x) \quad f(x) = \begin{cases} \neq 0 & \text{Poisson} \\ = 0 & \text{Laplace} \end{cases}$$

1D

$$\frac{d^2 \Psi}{dx^2} = -f(x)$$

3D

$$\frac{\partial^2 \Psi(x, y, z)}{\partial x^2} + \frac{\partial^2 \Psi(x, y, z)}{\partial y^2} + \frac{\partial^2 \Psi(x, y, z)}{\partial z^2} = -f(x, y, z)$$

Heat Flow (conduction)

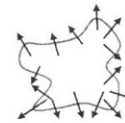
$\Psi(x) \leftrightarrow$ temperature

$\nabla \Psi(x) \propto$ heat flow

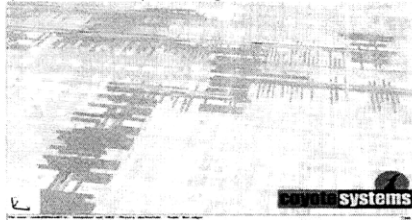
$\nabla \cdot h = \nabla \cdot \nabla \Psi = 0 \leftrightarrow$ Conservation Law

$\Leftrightarrow \sum h = 0$

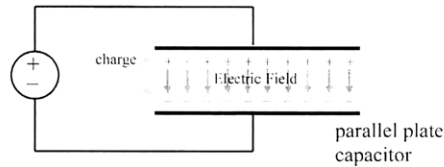
$$\int_V \nabla \cdot h dV = \int_S h \cdot n dS$$



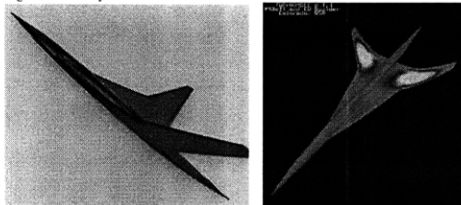
❖ Capacitance on a microprocessor signal line



- o Electrostatic Analysis
 - The Laplace Partial Differential Equation.
- o Electrostatic Analysis
 - $\Psi(x) \leftrightarrow$ electrostatic potential
 - $\nabla\Psi(x) \propto$ electric field



❖ Drag Force Analysis of Aircraft

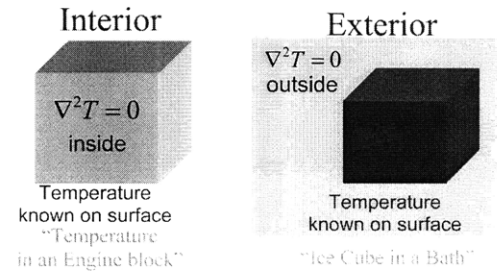


- o Potential Flow Equations

• Poisson Partial Differential Equations

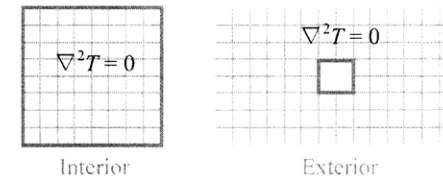
Potential Flow
 $\nabla\Psi =$ air velocity

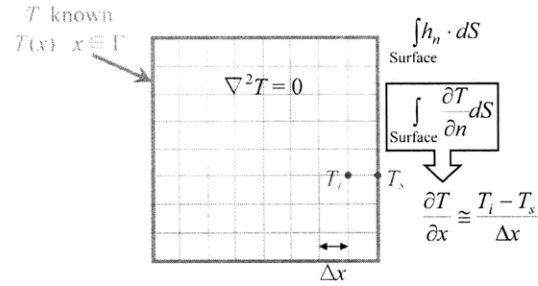
EXTERIOR VERSUS INTERIOR PROBLEMS



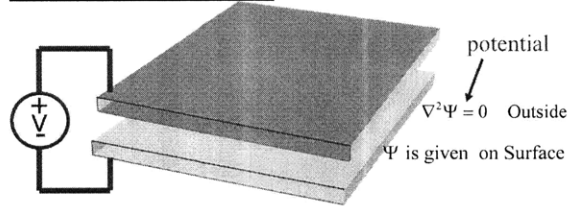
What is the heat flow?
 Heat Flow = Thermal conductivity $\int_{\text{surface}} \frac{\partial T}{\partial n}$

Finite Difference/Finite Element Method Meshing Problem





Exterior Problem in Electrostatics



What is the capacitance?

Capacitance = Dielectric Permittivity $\int_{\text{surface}} \frac{\partial \Psi}{\partial n}$

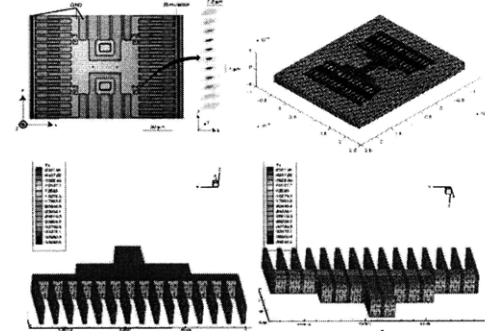
$\Psi(x) \leftrightarrow$ electrostatic potential

$\nabla \Psi(x) = \frac{\partial \Psi}{\partial n} \leftrightarrow$ electric field \propto charge on surface $\leftrightarrow Q$

capacitance = $\frac{Q}{V} = \frac{\text{charge}}{\text{voltage}} = Q_{i=1} = \int_{\text{surface}} \frac{d\Psi}{dn} dS$

$\nabla^2 \Psi = -\frac{\rho}{\epsilon}$

Drag Force in a Microresonator



Comb Drive

- Exterior Problem
- Complicated fluid flow problem
 - Interested in Quality Factor at Resonance
- Used as
- sensors (e.g. gyroscope)
 - RF resonator
- Stokes Equation – Linear
- High viscosity
 - Size of structure

Electromagnetic Coupling in a Package

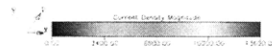
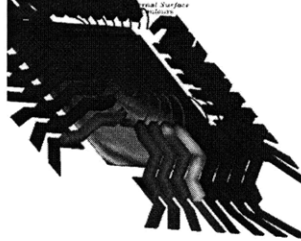


Image Thanks to Conventor

Full Wave
 $\nabla^2 \Psi + \kappa^2 \Psi = 0$ Helmholtz Equation
 $\nabla^2 \Psi = 0$ Laplace (or Poisson) MQS or EMQS
 Ψ = magnetic vector potential

Capacitance on a Microprocessor Signal Line



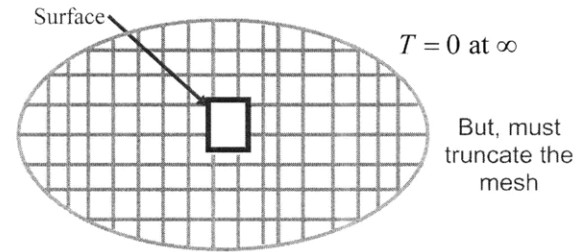
What is common about these problems

- ❖ Exterior Problems
 - Drag Force in MEMS device – fluid (air) creates drag.
 - Coupling in a Package – Field sin exterior create coupling
 - Capacitance of a Signal Line – Fields in exterior.
- ❖ Quantities of Interest are on the Surface
 - MEMS device – just want surface traction force
 - Package – just want coupling between conductors
 - Signal Line – just want surface charge
- ❖ Exterior Problem is Linear and Space-Invariant
 - MEMS – Exterior Stokes Flow equation (linear)
 - Package – Maxwell's equation in free space (linear)
 - Signal Line – Laplace's equation in free space (linear)

But problems are geometrically very complex!!

Exterior Problems

- ❖ Why not use Finite-Difference or FEM methods?
 2D Heat Flow Example



Only need $\frac{\partial T}{\partial n}$ on the surface, but T is computed everywhere
 Must truncate the mesh $\Rightarrow T(\infty) = 0$ becomes $T(R) = 0$

GREEN'S FUNCTION

Laplace's Equation

- ❖ Green's Function
 - Definition: Green's Function $G(x, x')$ is the solution at test location x produced by a unit point source at location x'
 - e.g. for Laplace operator:
 $\nabla^2 G(x, x') = \delta(x - x')$
 - 2-D

$$\Psi(x, y) = \log \sqrt{(x-x_0)^2 + (y-y_0)^2}$$

$$\frac{\partial^2 \Psi}{\partial x^2} + \frac{\partial^2 \Psi}{\partial y^2} = 0 \quad \text{for all } (x, y) \neq (x_0, y_0)$$

$$G(x, x') = \frac{1}{2\pi} \log|x-x'|$$

3-D

$$\Psi(x, y, z) = \frac{1}{\sqrt{(x-x_0)^2 + (y-y_0)^2 + (z-z_0)^2}}$$

$$\frac{\partial^2 \Psi}{\partial x^2} + \frac{\partial^2 \Psi}{\partial y^2} + \frac{\partial^2 \Psi}{\partial z^2} = 0 \quad \text{for all } (x, y, z) \neq (x_0, y_0, z_0)$$

$$G(x, x') = -\frac{1}{4\pi} \frac{1}{|x-x'|}$$

Proof: Just differentiate and see!

$$\phi(x) = \frac{1}{4\pi\epsilon} \frac{1}{|x|}$$

$$G(x, x') = -\frac{1}{4\pi} \frac{1}{|x-x'|}$$

$$\nabla^2 G = \delta(x-x')$$

$$\phi(x) = \frac{1}{4\pi\epsilon} \frac{1}{|x-x'|}$$

$$\nabla^2 \phi = -\frac{\delta(x-x')}{\epsilon}$$

$\nabla^2 \Psi = \delta(x)$

$\nabla^2 \Psi = 0 \quad x \neq 0$

Green Function

$$\Psi(x) = \begin{cases} \log|x| & \text{2D} \\ \frac{1}{|x|} & \text{3D} \end{cases}$$

Think Electrostatics

$$\nabla^2 \phi = \delta(x)$$

$$\phi(x) = \frac{1}{4\pi\epsilon|x|}$$

$$\nabla^2 \phi = \delta(x-x_0)$$

$$\phi(x) = \frac{1}{4\pi\epsilon|x-x_0|} = \frac{1}{4\pi\epsilon} \frac{1}{\sqrt{(x-x_0)^2 + (y-y_0)^2 + (z-z_0)^2}}$$

Simple Idea (in 2-D)

One Point $\Psi(x, y)$ is given on surface

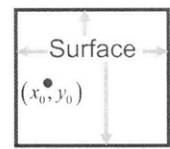
$$\frac{\partial^2 \Psi}{\partial x^2} + \frac{\partial^2 \Psi}{\partial y^2} = 0 \quad \text{outside}$$

Let $\Psi(x, y) = \log \sqrt{(x-x_0)^2 + (y-y_0)^2}$

$$\frac{\partial^2 \Psi}{\partial x^2} + \frac{\partial^2 \Psi}{\partial y^2} = 0 \quad \text{outside}$$

~~Problem Solved~~

Does not match boundary conditions!



SPHERE

Point source
 $R = 10$
 $T = 1$
 $(0,0,0)$
 $\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} + \frac{\partial^2 u}{\partial z^2} = 0$
 $u = \frac{10}{\sqrt{x^2 + y^2 + z^2}}$

CIRCLE

Can you think of a problem we can solve with just one point source?

Point source
 (x_0, y_0)
 $R = 10$
 $\Psi = 1$
 Outside:
 $\nabla^2 \Psi = 0$
 $\Leftrightarrow \frac{\partial^2 \Psi}{\partial x^2} + \frac{\partial^2 \Psi}{\partial y^2} = 0$

$\Psi = A \log \sqrt{(x-x_0)^2 + (y-y_0)^2}$

$A = \frac{1}{\log R}$ Weighted point source

Note: with one point source I can match only one B.C. For more complicated B.C. I could use more point sources, how many point sources do I need? What do I put them?

o "More Points"
 $\Psi(x, y)$ is given on surface

$\frac{\partial^2 \Psi}{\partial x^2} + \frac{\partial^2 \Psi}{\partial y^2} = 0$ outside

Let $\Psi(x, y) = \sum_{j=1}^n \omega_j \log \sqrt{(x-x_j)^2 + (y-y_j)^2} = \sum_{j=1}^n \omega_j G(x-x_j, y-y_j)$

Pick the ω_j 's to match the boundary conditions!

Source Strengths selected to give correct potential at test points.


$$\begin{bmatrix} G(x_1 - x_1, y_1 - y_1) & \dots & G(x_1 - x_n, y_1 - y_n) \\ \vdots & & \vdots \\ G(x_c - x_1, y_c - y_1) & \dots & G(x_c - x_n, y_c - y_n) \end{bmatrix} \begin{bmatrix} \omega_1 \\ \vdots \\ \omega_n \end{bmatrix} = \begin{bmatrix} \Psi(x_1, y_1) \\ \vdots \\ \Psi(x_c, y_c) \end{bmatrix}$$

$\Psi(x_i, y_i) = w_1 G(x_i - x_1, y_i - y_1) + w_2 G(x_i - x_2, y_i - y_2) + \dots$

SOURCE POINT LOCATIONS

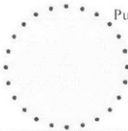
If Cluster of Charges in Center:

- Smaller oscillations
- Less change to B.C.
- Matrix is ill conditioned (if all on top of one another – singular matrix)

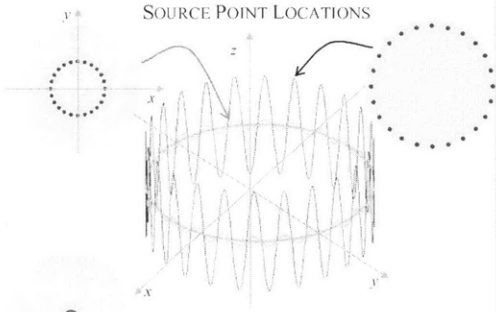


Push Charges to Surface:

- To avoid ill conditioning
- Use many points to reduce oscillation amplitude

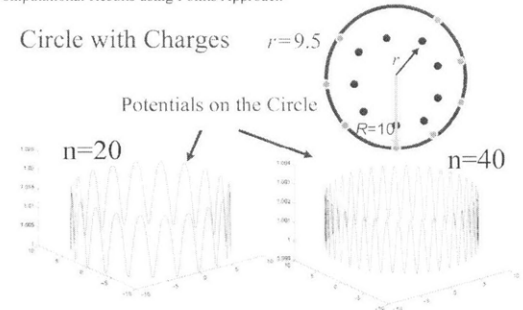


SOURCE POINT LOCATIONS



More singular weight matrix.

o Computational Results using Points Approach



$$\Psi(x) = \sum_i \overset{\text{unknown}}{w_i} G(x, x_i)$$

$$\Psi(x)_{\text{surface}} = \int_{\text{surface}} \overset{\text{unknown}}{\sigma(x')} G(x, x') dx' \Leftrightarrow \nabla^2 \Psi = 0 + \underset{\Psi(x)=\Psi(x)_{\text{surface}}}{\text{B.C}}$$

given unknown

❖ Charge Density
Want to smear point charges to the surface

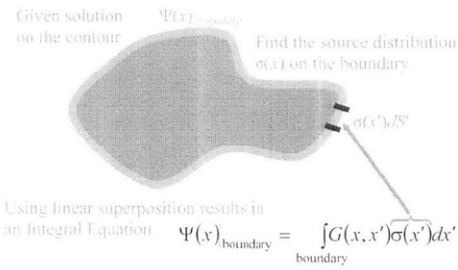


Results in an Integral Equation

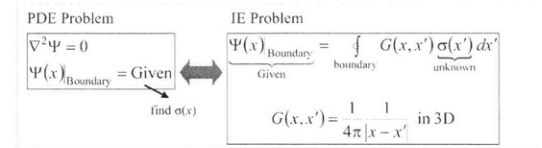
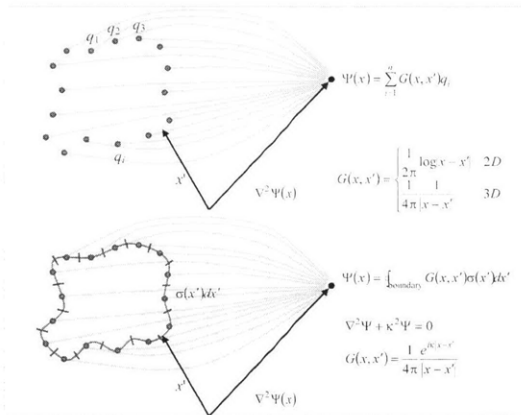
$$\Psi(x) = \int_{\text{surface}} G(x, x') \sigma(x') dS'$$

How do we solve the integral equation?

❖ Integral Formulation



How do we solve the integral equation?



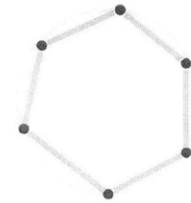
BASIS FUNCTIONS

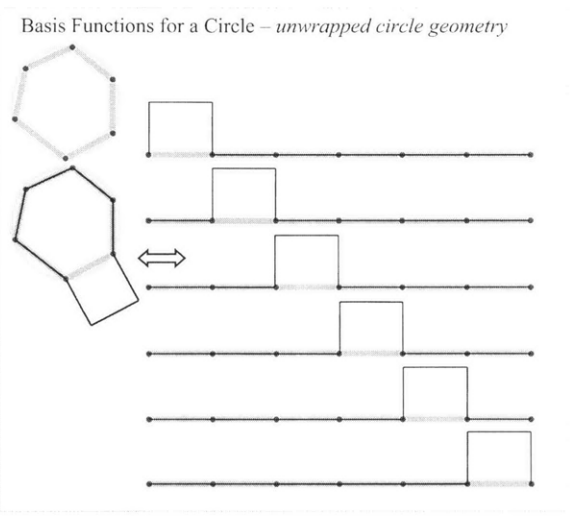
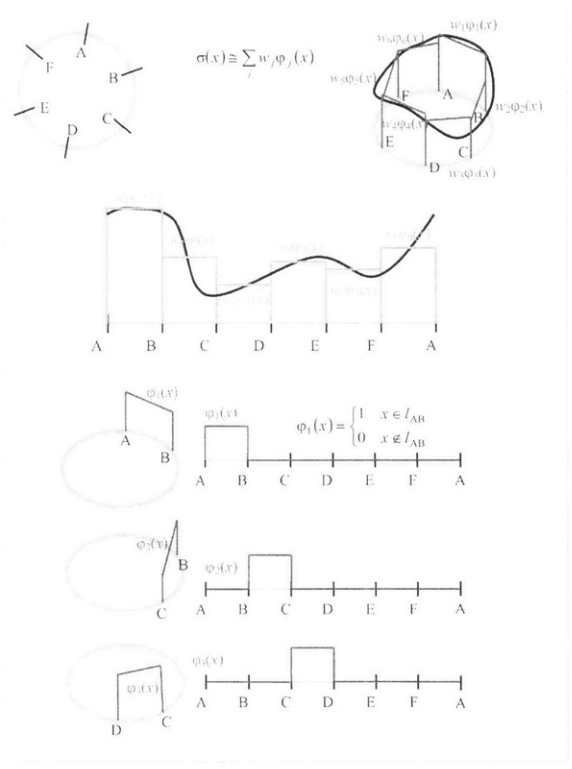
Basic Idea

Represent $\sigma(x) = \sum_{i=1}^n \omega_i \underbrace{\phi_i(x)}_{\text{Basis Functions}}$

Example Basis

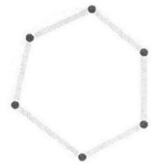
Represent a circle with straight lines
 Assume σ is constant along each line
 The basis functions are "on" the surface
 Can be used to approximate the density
 May also approximate the geometry



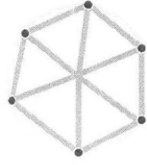


Geometric Approximation

❖ Not a New Idea



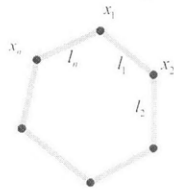
Piecewise Straight surface basis functions approximate the circle



Triangles for 2-D FEM approximate the circle too!

$$\Psi(x) = \int_{\text{approx. surface}} G(x, x') \sum_{i=1}^n \omega_i \phi_i(x') dS'$$

❖ Piecewise Constant Straight Sections Example



Represent $\sigma(x) = \sum_{i=1}^n \omega_i \phi_i(x)$
Basis Functions

- 1.) Pick a set of n points on the surface
- 2.) Define a new surface by connecting points with n lines
- 3.) Define $\phi_j(x) = 1$ if x is on line l_j otherwise, $\phi_j(x) = 0$

$$Q = \int_I \sigma(x) dx = \int \sum_{i=1}^n \omega_i \phi_i(x)$$

$$\Psi(x) = \int_I G(x, x') \sum_{i=1}^n \omega_i \phi_i(x') dx'$$

$$\begin{aligned} \Psi(x)_{\text{boundary}} &= \int_{\text{boundary}} G(x, x') \sigma(x') dx' \\ &= \int_{\text{approximate boundary}} G(x, x') \sum_{j=1}^n \omega_j \phi_j(x) dx' \\ &= \sum_{j=1}^n \omega_j \int_{\text{boundary}} G(x, x') \phi_j(x') dx' \\ &= \sum_{j=1}^n \omega_j \int_{\text{line } l_j} G(x, x') dx' \end{aligned}$$

How do we determine the ω_j 's?

$$\nabla^2 \Psi = 0 + \Psi(x)_{\text{surface}} \text{ B.C.}$$

Basis Functions $\sigma(x) \cong \sum_{j=1}^n w_j \phi_j(x)$

Integral Equation $\Psi(x)_{\text{surface}} = \int_{\text{surface}} G(x, x') \sigma(x') dS'$

Substitute Basis Functions $\Psi(x)_{\text{surface}} \cong \int_{\text{surface}} G(x, x') \sum_{j=1}^n w_j \phi_j(x) dS'$
 $\cong \sum_{j=1}^n w_j \int_{\text{surface}} G(x, x') \phi_j(x) dS'$

If basis functions are piecewise constant $\Psi(x)_{\text{surface}} \cong \sum_{j=1}^n w_j \int_{l_j} G(x, x') dS'$

$$\phi_j(x') = \begin{cases} 1 & \text{if } x' \in l_j \\ 0 & \text{if } x' \notin l_j \end{cases}$$

unknowns
 We are going to try and find the weights w_j

Q1. How do I find w_j ?
 Q2. Given w_j , what is the charge on the surface?

$$\sigma(x) = \sum_{j=1}^n w_j \phi_j(x)$$

 Q3. Given w_j , what is the potential anywhere outside?

$$\phi(x) = \sum_{j=1}^n w_j \int_{S_j} G(x, x') dS'$$

 Can evaluate this anywhere outside if we evaluate on the surface we obtain the given B.C.

Residuals

❖ Residual Definition

$$R(x) \equiv \Psi(x)_{\text{boundary}} - \int_{\text{approx surface}} G(x, x') \sum_{i=1}^n \omega_i \phi_i(x') dS'$$

❖ Residual Minimization

We will pick the ω_i 's to make $R(x)$ small.
 General Approach: Pick a set of test functions $\phi_1, \phi_2, \dots, \phi_n$ and force $R(x)$ to be orthogonal to the set $\int \phi_i(x) R(x) dS = 0$ for all i .

$\phi_i = ?$

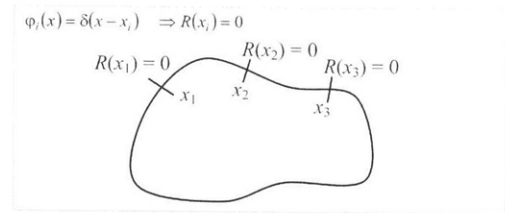
$$\int_{\text{surface}} \phi_i(x) R(x) dS = 0 \rightarrow R(x_i) = 0 \quad i = 1, 2, \dots, N$$

$$\int_{\text{surface}} \delta(x - x_i) R(x) dS = R(x_i)$$

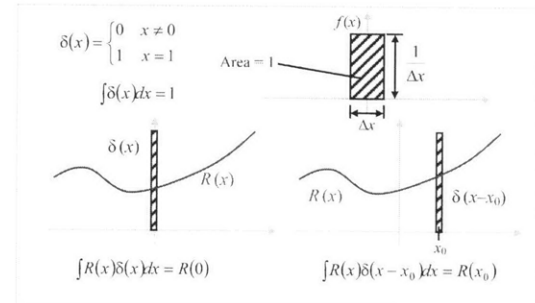
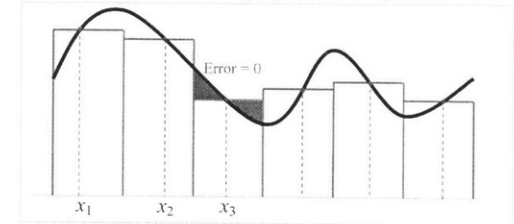
○ Residual Minimization Using Test Functions

$$\int \phi_i(x) R(x) dS = \int \phi_i(x) \Psi(x) dS - \int \int_{\text{approx surface}} \phi_i(x) G(x, x') \sum_{j=1}^n \omega_j \phi_j(x') dS' dS = 0$$

We will generate the different methods by choosing the $\phi_1, \phi_2, \dots, \phi_n$
 Collocation: $\phi_i(x) = \delta(x - x_i)$ (point - matching)



Galerkin Method: $\phi_i(x) = \phi_j(x)$ (basis = test)



COLLOCATION METHOD



$$w_j \quad j = 1, 2, \dots, n$$

$$\varphi(x_i)_{\text{surface}} = \sum_{j=1}^n w_j \underbrace{\int_{I_j} G(x_i, x') dS'}_{A_{ij}}$$

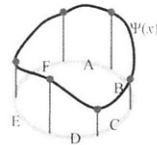
$$[A][w] = [\varphi(x_i)]$$

Collocation: $\phi_i(x) = \delta(x_i)$ (point - matching)

$$\int \delta(x - x_i) R(x) dS = R(x_i) = \Psi(x_i) - \int_{\text{surface}} G(x_i, x') \sum_{j=1}^n \omega_j \varphi_j(x') dS' = 0$$

$$\Rightarrow \Psi(x_i) = \sum_{j=1}^n \omega_j \underbrace{\int_{\text{surface}} G(x_i, x') \varphi_j(x') dS'}_{A_{ij}}$$

$$\begin{bmatrix} A_{1,1} & \dots & A_{1,n} \\ \vdots & \ddots & \vdots \\ A_{n,1} & \dots & A_{n,n} \end{bmatrix} \begin{bmatrix} \omega_1 \\ \vdots \\ \omega_n \end{bmatrix} = \begin{bmatrix} \Psi(x_1) \\ \vdots \\ \Psi(x_n) \end{bmatrix}$$



$$R(x_i) = 0 \quad i = 1, 2, \dots, N$$

$$\Psi(x)_{\text{surface}} = \int_{\text{surface}} G(x, x') \sigma(x') dS'$$

Basis Functions
 $\sigma(x) \cong \sum w_j \varphi_j(x)$

$$\Psi(x)_{\text{surface}} = \sum_{j=1}^n w_j \int_{\text{surface}} G(x, x') \sigma_j(x') dS'$$

Collocation

$$\Psi(x_i) = \sum_{j=1}^n w_j \underbrace{\int_{\text{surface}} G(x_i, x') \sigma_j(x') dS'}_{A_{ij}}$$

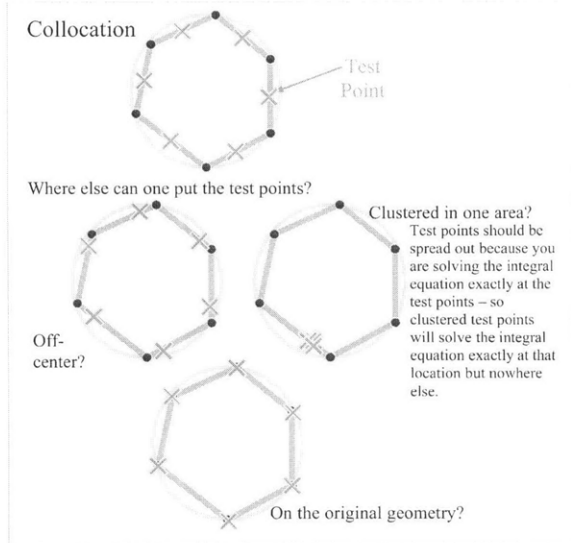
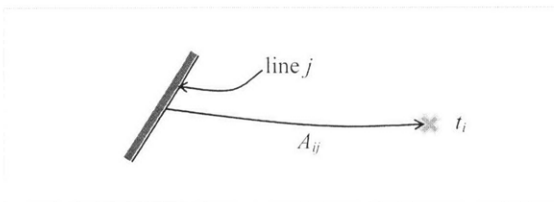
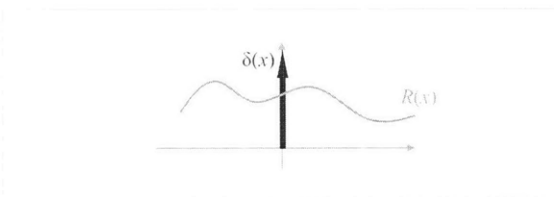
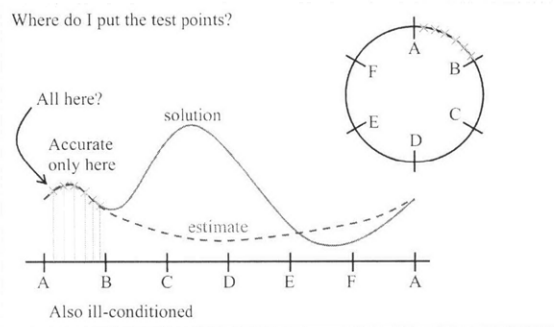
Piecewise constant basis functions

$$\Psi(x_i) = \sum_{j=1}^n w_j \underbrace{\int_{I_j} G(x_i, x') dS'}_{A_{ij}}$$

test source

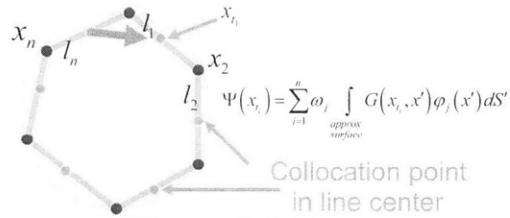
$$\sum_{j=1}^n w_j \underbrace{\int_{I_j} G(x_i, x') dS'}_{A_{ij}} = \underbrace{\Psi(x_i)}_{\text{given}} \quad i = 1, 2, \dots, N$$

$$\begin{bmatrix} A_{11} & \dots & A_{1n} \\ \vdots & \ddots & \vdots \\ A_{n1} & \dots & A_{nn} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} = \begin{bmatrix} \Psi(x_1) \\ \Psi(x_2) \\ \vdots \\ \Psi(x_n) \end{bmatrix}$$



Centroid Collocation

❖ Piecewise Constant Bases



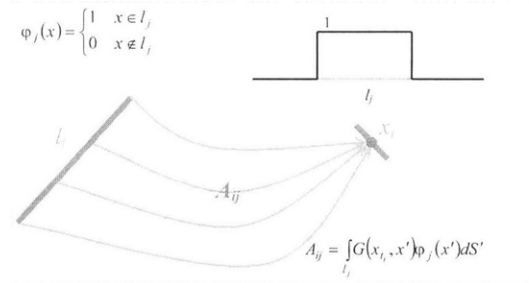
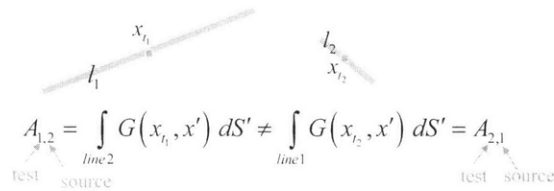
$$\Psi(x_i) = \sum_{j=1}^n \omega_j \int_{\text{approx surface}} G(x_i, x') \varphi_j(x') dS'$$

$$\begin{bmatrix} A_{1,1} & \dots & A_{1,n} \\ \vdots & \ddots & \vdots \\ A_{n,1} & \dots & A_{n,n} \end{bmatrix} \begin{bmatrix} \omega_1 \\ \vdots \\ \omega_n \end{bmatrix} = \begin{bmatrix} \Psi(x_{i_1}) \\ \vdots \\ \Psi(x_{i_n}) \end{bmatrix}$$

$$\Psi(x_i) = \sum_{j=1}^n \omega_j \int_{\text{line } j} G(x_i, x') dS'$$

❖ Generates a Nonsymmetric A

$$\Psi(x_i) = \sum_{j=1}^n \omega_j \int_{\text{line } j} G(x_i, x') dS'$$



GALERKIN METHOD

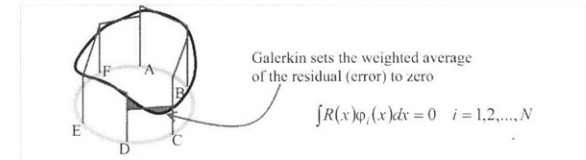
Galerkin: $\varphi_i(x) = \varphi_j(x)$ (test = basis)

$$\int_{\text{approx surface}} \varphi_i(x) R(x) dS = \int_{\text{approx surface}} \varphi_i(x) \Psi(x) dS - \int_{\text{approx surface}} \varphi_i(x) G(x, x') \sum_{j=1}^n \omega_j \varphi_j(x') dS' dS = 0$$

$$\int_{\text{approx surface}} \varphi_i(x) \Psi(x) dS = \sum_{j=1}^n \omega_j \int_{\text{approx surface}} \int_{\text{approx surface}} G(x, x') \varphi_i(x) \varphi_j(x') dS' dS$$

$$\begin{bmatrix} A_{1,1} & \dots & A_{1,n} \\ \vdots & \ddots & \vdots \\ A_{n,1} & \dots & A_{n,n} \end{bmatrix} \begin{bmatrix} \omega_1 \\ \vdots \\ \omega_n \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$

If $G(x, x') = G(x', x)$ then $A_{i,j} = A_{j,i}$ A is symmetric



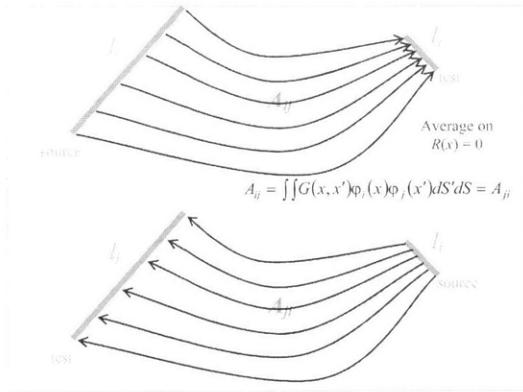
$$\Psi(x)_{surface} = \int_{surface} G(x, x') \sigma(x') dS'$$

$$\Psi(x)_{surface} = \sum_{j=1}^n w_j \int_{surface} G(x, x') \sigma_j(x') dS'$$

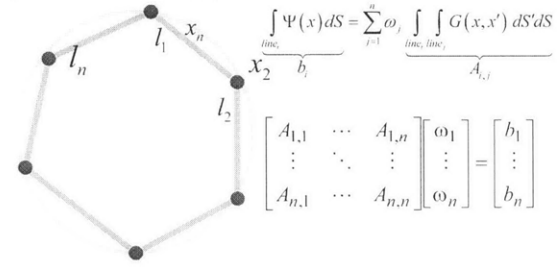
$$\int_s \Psi(x) \phi_i(x) dS = \sum_{j=1}^n w_j \int_s \int_{s'} G(x, x') \phi_i(x) \phi_j(x') dS' dS$$

$$\int_{l_i} \Psi(x) dS = \sum_{j=1}^n w_j \int_{l_i} \int_{l_j} G(x, x') dS' dS$$

Basis Functions
 $\sigma(x') \equiv \sum w_j \phi_j(x')$
 Galerkin
 Piecewise constant basis functions



Piecewise Constant Bases



INTRODUCTION TO NUMERICAL SIMULATION

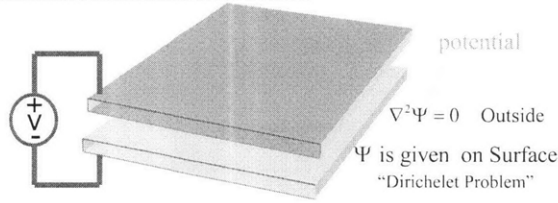
LECTURE 23.

PDE – BEM Integral Equation Method II

TODAY'S OUTLINE:

- ❖ Collocation Method
- ❖ Galerkin Method
- ❖ 3D Panel Integration
- ❖ Solving Discretized Integral Equations

EXTERIOR PROBLEM IN ELECTROSTATICS



❖ First Kind Integral Equation for Charge:

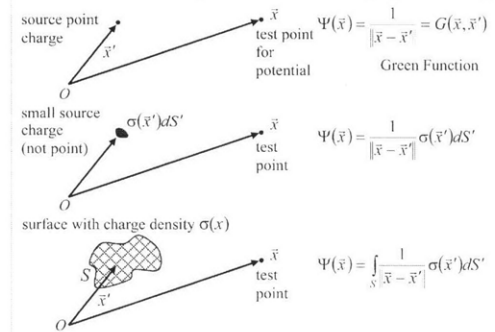
$$\Psi(x) = \int_{\text{surface}} \underbrace{\frac{1}{\|x-x'\|}}_{\text{Green's Function}} \underbrace{\sigma(x')}_{\text{Charge Density}} dS'$$

potential

First Kind Integral Equation for Charge

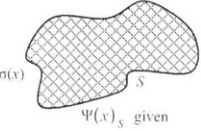
$$\nabla^2 \Psi = 0 \Rightarrow \Psi(\bar{x}) = \int_S \frac{1}{\|\bar{x} - \bar{x}'\|} \sigma(\bar{x}') dS' \quad \text{in 3D}$$

one way to view this is as a superposition integral



$$\nabla^2 \Psi = 0 + \underbrace{\Psi(\bar{x})_S}_{\text{known boundary condition}} = \int_S \frac{\underbrace{\sigma(\bar{x}')}_{\text{unknown}}}{\|\bar{x} - \bar{x}'\|} dS'$$

Problem.



$\sigma(x)$ $\Psi(x)$ given

$\nabla^2 \Psi = 0$ everywhere outside

Find $\sigma(x)$ on S

Solve $\frac{\Psi(\bar{x})}{\text{given}} = \int_S \frac{1}{|\bar{x} - \bar{x}'|} \sigma(\bar{x}') dS'$

Use Basis Functions $\sigma(\bar{x}') \cong \sum w_j \phi_j(\bar{x}')$

$\Psi(\bar{x})_{\bar{x} \in S} = \int_S \frac{1}{|\bar{x} - \bar{x}'|} \sum_{j=1}^N w_j \phi_j(\bar{x}') dS'$

e.g. piecewise constant basis (panels)

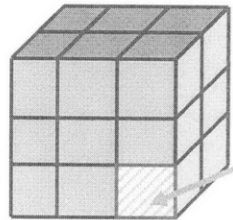
$\Psi(\bar{x})_{\bar{x} \in S} = \sum_{j=1}^N w_j \int_{\text{panel } j} \frac{1}{|\bar{x} - \bar{x}'|} dS'$

$\phi_j(\bar{x}') = \begin{cases} 1 & \text{if } \bar{x}' \in \text{panel } j \\ 0 & \text{otherwise} \end{cases}$

N unknowns

BASIS FUNCTION APPROACH
 ❖ Piecewise Constant Basis Function

Integral Equation: $\Psi(x) = \int_{\text{surface}} \frac{1}{|x - x'|} \sigma(x') dS'$

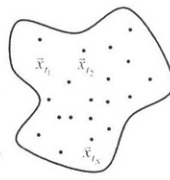


Discretize Surface into Panels

Represent $\sigma(x) \approx \sum_{j=1}^n \omega_j \phi_j(x)$

Panel j
 $\phi_j(x) = 1$ if x is on panel j
 $\phi_j(x) = 0$ otherwise

Collocation.



N equations

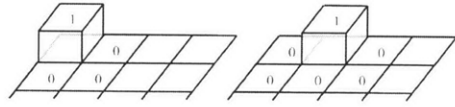
$$\begin{cases} \Psi(\bar{x}_{i_1})_{\bar{x}_{i_1} \in S} = \sum_{j=1}^N w_j \int_{S_j} \frac{1}{|\bar{x}_{i_1} - \bar{x}'|} dS' \\ \Psi(\bar{x}_{i_2})_{\bar{x}_{i_2} \in S} = \sum_{j=1}^N w_j \int_{S_j} \frac{1}{|\bar{x}_{i_2} - \bar{x}'|} dS' \\ \vdots \\ \Psi(\bar{x}_{i_N})_{\bar{x}_{i_N} \in S} = \sum_{j=1}^N w_j \int_{S_j} \frac{1}{|\bar{x}_{i_N} - \bar{x}'|} dS' \end{cases}$$

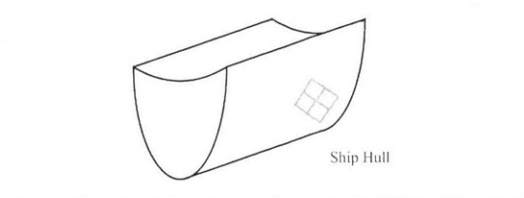
$$\begin{bmatrix} \int_{S_1} \frac{1}{|\bar{x}_{i_1} - \bar{x}'|} dS' & \int_{S_2} \frac{1}{|\bar{x}_{i_1} - \bar{x}'|} dS' & \dots \\ \int_{S_1} \frac{1}{|\bar{x}_{i_2} - \bar{x}'|} dS' & \int_{S_2} \frac{1}{|\bar{x}_{i_2} - \bar{x}'|} dS' & \dots \\ \vdots & \vdots & \ddots \\ \int_{S_1} \frac{1}{|\bar{x}_{i_N} - \bar{x}'|} dS' & \int_{S_2} \frac{1}{|\bar{x}_{i_N} - \bar{x}'|} dS' & \dots \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_N \end{bmatrix} = \begin{bmatrix} \Psi(\bar{x}_{i_1})_{\bar{x}_{i_1} \in S} \\ \vdots \\ \Psi(\bar{x}_{i_N})_{\bar{x}_{i_N} \in S} \end{bmatrix}$$

Piecewise Constant Basis Functions + Collocation.

$$A_{i,j} = \int_{S_j} \frac{1}{|\bar{x}_{i_1} - \bar{x}'|} dS'$$

Basis Functions





Ship Hull

- 1.) Discretize Unknowns $\alpha(x) \rightarrow N$ unknowns using basis functions
- 2.) Discretize Equations \rightarrow Collocation or Galerkin

❖ Residual Definition

$$R(x) \equiv \Psi(x) - \int_{\text{approx surface}} G(x, x') \sum_{i=1}^n \omega_i \phi_i(x') dS'$$

❖ Residual Minimization

We will pick the ω_i 's to make $R(x)$ small.
 General Approach: Pick a set of test functions $\phi_1, \phi_2, \dots, \phi_n$ and force $R(x)$ to be orthogonal to the set $\int \phi_i(x) R(x) dS = 0$ for all i .

○ Residual Minimization Using Test Functions

$$\int \phi_i(x) R(x) dS = \int \phi_i(x) \Psi(x) dS - \int \int_{\text{approx surface}} \phi_i(x) G(x, x') \sum_{j=1}^n \omega_j \phi_j(x') dS' dS = 0$$

We will generate the different methods by choosing the $\phi_1, \phi_2, \dots, \phi_n$
 Collocation: $\phi_i(x) = \delta(x - x_i)$ (point - matching)

Collocation.

$$\bar{x}_{i_1} \quad \bar{x}_{i_2} \quad \dots \quad \bar{x}_{i_N}$$

$$\Psi(\bar{x}_{i_1})_{\bar{x}_{i_1} \in S} = \sum_{j=1}^N w_j \int_{A_j} \frac{1}{|\bar{x}_{i_1} - \bar{x}'|} dS' \quad i = 1, 2, \dots, N$$

$$\Psi(\bar{x}_{i_1})_{\bar{x}_{i_1} \in S} = \sum_{j=1}^N w_j A_{ij} \quad \left[\begin{matrix} \mathbf{A} \\ \mathbf{\bar{w}} \end{matrix} \right] = \left[\begin{matrix} \Psi(\bar{x}_{i_1}) \\ \vdots \\ \Psi(\bar{x}_{i_N}) \end{matrix} \right]$$

Galerkin Method: $\phi_i(x) = \varphi_i(x)$ (basis = test)

COLLOCATION METHOD

❖ Basis Function Approach

Collocation: $\phi_i(x) = \delta(x - x_i)$ (point - matching)

$$\int \delta(x - x_i) R(x) dS = R(x_i) = \Psi(x_i) - \int_{\text{approx surface}} G(x_i, x') \sum_{j=1}^n \omega_j \phi_j(x') dS' = 0$$

$$\Rightarrow \Psi(x_i) = \sum_{j=1}^n \omega_j \int_{\text{approx surface}} G(x_i, x') \phi_j(x') dS'$$

$$\left[\begin{matrix} A_{1,1} & \dots & A_{1,n} \\ \vdots & \ddots & \vdots \\ A_{n,1} & \dots & A_{n,n} \end{matrix} \right] \left[\begin{matrix} \omega_1 \\ \vdots \\ \omega_n \end{matrix} \right] = \left[\begin{matrix} \Psi(x_{i_1}) \\ \vdots \\ \Psi(x_{i_n}) \end{matrix} \right]$$

Collocation.

$$\Psi(x_{t_1}) = \sum_{j=1}^N w_j \int_{\text{panel } j} \frac{1}{|x_{t_1} - x'|} dS'$$

$$\Psi(x_{t_2}) = \sum_{j=1}^N w_j \int_{\text{panel } j} \frac{1}{|x_{t_2} - x'|} dS'$$

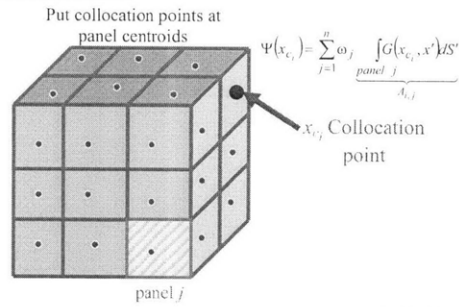
$$\Psi(x_{t_n}) = \sum_{j=1}^N w_j \int_{\text{panel } j} \frac{1}{|x_{t_n} - x'|} dS'$$

$$\Psi(x_{t_i}) = \sum_{j=1}^N w_j A_{ij}$$

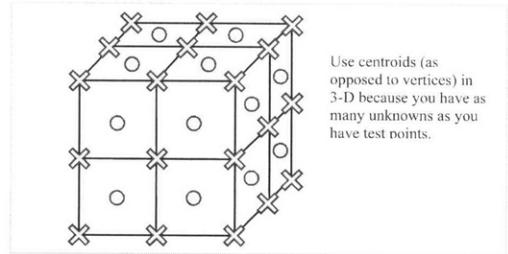
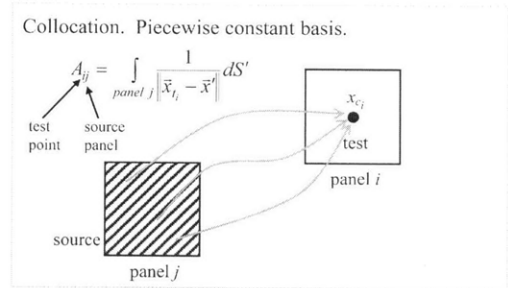
$$A_{ij} = \int_{\text{panel } j} \frac{1}{|x_{t_i} - x'|} dS'$$

$$\begin{bmatrix} A_{11} & A_{12} & \dots & A_{1N} \\ \vdots & \vdots & \ddots & \vdots \\ A_{n1} & A_{n2} & \dots & A_{nN} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_N \end{bmatrix} = \begin{bmatrix} \Psi(x_{t_1}) \\ \Psi(x_{t_2}) \\ \vdots \\ \Psi(x_{t_n}) \end{bmatrix}$$

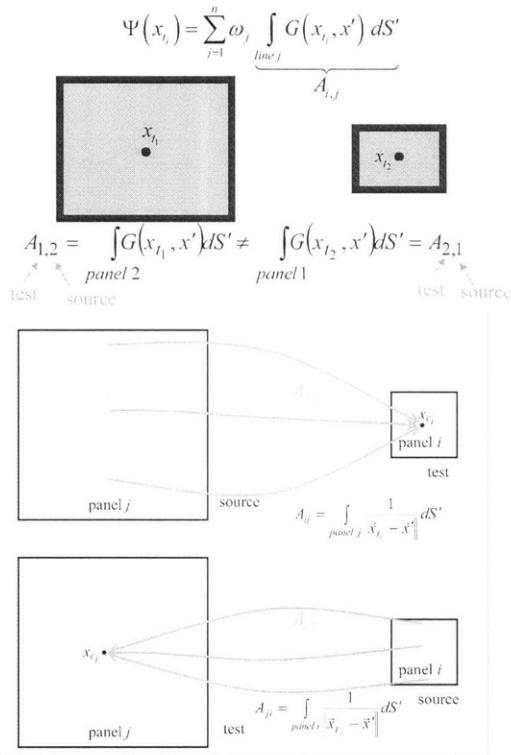
o Centroid Collocation



$$\begin{bmatrix} A_{1,1} & \dots & A_{1,n} \\ \vdots & \ddots & \vdots \\ A_{n,1} & \dots & A_{n,n} \end{bmatrix} \begin{bmatrix} \omega_1 \\ \vdots \\ \omega_n \end{bmatrix} = \begin{bmatrix} \Psi(x_{t_1}) \\ \vdots \\ \Psi(x_{t_n}) \end{bmatrix}$$



o Nonsymmetric A



GALERKIN METHOD

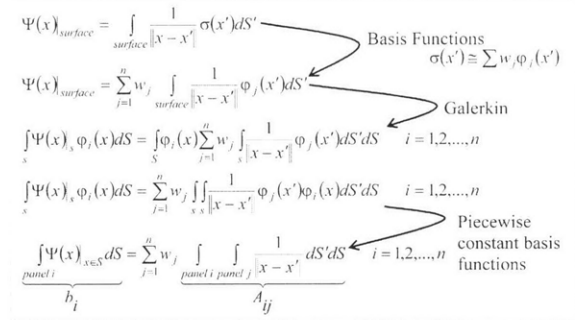
❖ Piecewise Constant Basis Function

Galerkin: $\phi_i(x) = \phi_j(x)$ (test = basis)

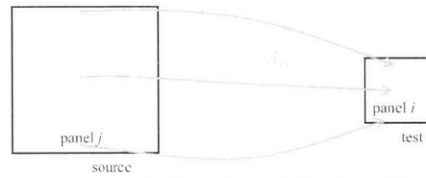
$$\int_{\text{approx surface}} \phi_i(x) R(x) dS = \int_{\text{approx surface}} \phi_i(x) \Psi(x) dS - \int_{\text{approx surface}} \phi_i(x) G(x, x') \sum_{j=1}^n \omega_j \phi_j(x') dS' dS = 0$$

$$\int_{\text{approx surface}} \phi_i(x) \Psi(x) dS = \sum_{j=1}^n \omega_j \int_{\text{approx surface}} \int_{\text{approx surface}} G(x, x') \phi_i(x) \phi_j(x') dS' dS$$

$$\begin{bmatrix} A_{1,1} & \dots & A_{1,n} \\ \vdots & \ddots & \vdots \\ A_{n,1} & \dots & A_{n,n} \end{bmatrix} \begin{bmatrix} \omega_1 \\ \vdots \\ \omega_n \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$



$$\begin{bmatrix} \leftarrow \int_S \int_S \frac{\phi_1(x) \phi_1(x')}{|x-x'|} dx' dx = \int_{S_1, S_1} \frac{1}{|x-x'|} dx' dx \rightarrow \uparrow \\ \leftarrow \int_S \int_S \frac{\phi_2(x) \phi_2(x')}{|x-x'|} dx' dx = \int_{S_2, S_2} \frac{1}{|x-x'|} dx' dx \rightarrow \downarrow \\ \vdots \end{bmatrix} \vec{w} = \begin{bmatrix} \int_{S_1} \Psi(x) dx \\ \int_{S_2} \Psi(x) dx \\ \vdots \end{bmatrix}$$



Galerkin requires that the boundary conditions are satisfied "on averages" over each test panel i

General Galerkin: $A_{ij} = \int_S \int_S \phi_i(x) \phi_j(x') G(x, x') dx' dx$

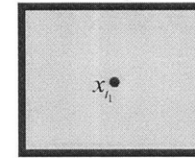
$$\int_{S_i} \phi_i(x) \Psi(x) dS = \sum_{j=1}^n \omega_j \int_{S_i} \int_{S_j} \phi_i(x) G(x, x') \phi_j(x') dS' dS$$

$$\int_{S_i} \phi_i(x) \Psi(x) dS = \sum_{j=1}^n \omega_j \int_{\text{panel } i} \int_{\text{panel } j} \frac{1}{|x-x'|} dS' dS$$

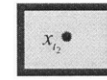
$$\begin{bmatrix} A_{1,1} & \dots & A_{1,n} \\ \vdots & \ddots & \vdots \\ A_{n,1} & \dots & A_{n,n} \end{bmatrix} \begin{bmatrix} \omega_1 \\ \vdots \\ \omega_n \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$

❖ Symmetric A

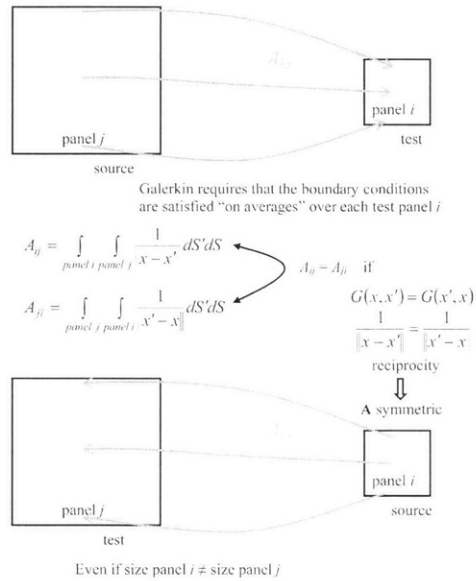
$$\int_{S_i} \Psi(x) dS = \sum_{j=1}^n \alpha_j \int_{\text{panel } i} \int_{\text{panel } j} \frac{1}{|x-x'|} dS' dS$$



If $G(x, x') = G(x', x)$ then $A_{i,j} = A_{j,i}$

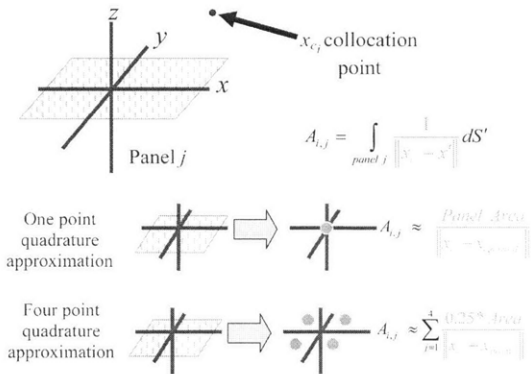


A is symmetric



3D PANEL INTEGRATION

❖ Collocation Approach – Calculating Matrix Elements



Example: Thin Metal Strip

Piecewise Constant Basis Function

Unknowns: w_j charges on each panel
Assume $\Psi(x) = 1$ given on the metal strip

$$\phi_j(x') = \begin{cases} 1 & x' \in \text{panel } j \\ 0 & x' \notin \text{panel } j \end{cases}$$

Toeplitz Matrix

Collocation

$$A_{1,6} = \int_{\text{panel 6}} \frac{1}{\|x_{c1} - x'\|} dS' = \frac{\text{Area Panel 6}}{\|x_{c1} - x_{c6}\|} = \frac{1}{5}$$

$$A_{2,6} = \int_{\text{panel 6}} \frac{1}{\|x_{c2} - x'\|} dS' = \frac{\text{Area Panel 6}}{\|x_{c1} - x_{c2}\|} = \frac{1}{4}$$

$$A_{6,1} = \int_{\text{panel 1}} \frac{1}{\|x_{c6} - x'\|} dS' = \frac{\text{Area Panel 1}}{\|x_{c6} - x_{c1}\|} = \frac{1}{5}$$

$A_{6,1} = A_{1,6}$ Why? Collocation in general gives $A_{6,1} \neq A_{1,6}$
But if $\text{Area}_{\text{panel 1}} = \text{Area}_{\text{panel 6}}$, then $A_{6,1} = A_{1,6}$

Self-term
What about $A_{1,1}$?

$$A_{1,1} = \int_{\text{panel 1}} \frac{1}{\|x_{c1} - x'\|} dS' = \frac{\text{Area Panel 1}}{\|x_{c1} - x_{c1}\|} = \frac{1}{0} \quad \text{No.}$$

Four Point Quadrature

$$A_{1,6} = \int_{\text{panel 6}} \frac{1}{\|x_{c1} - x'\|} dS'$$

$$= \text{Area 6A} + \text{Area 6B} + \text{Area 6C} + \text{Area 6D}$$

Does $A_{6,1} = A_{1,6}$ with four point quadrature?
In general, yes, if shape panel j = shape panel i

❖ Basis Function Approach – Calculating “Self-Term”

$$A_{i,i} = \int_{\text{panel } i} \frac{1}{\|x_{c,i} - x'\|} dS'$$

One point quadrature Approximation

$$A_{i,i} \approx \frac{\text{Panel Area}}{0}$$

$A_{i,i} = \int_{\text{panel } i} \frac{1}{\|x_{c,i} - x'\|} dS'$ is an integrable singularity

$$A_{i,j} = \int_{\text{panel } i} \frac{1}{\|x_j - x'\|} dS'$$

$$A_{i,j} = \int_{\text{disk}} \frac{1}{\|x_j - x'\|} dS' = \int_0^R \int_0^{2\pi} \frac{1}{r} r dr d\theta = 2\pi R$$

Integrate in two pieces

Disk Integral has singularity but has analytic formula

Self-Term

$A_{6,6} = ?$

$$\int_{\text{disk}} \frac{1}{\|x - x'\|} dx' = \int_{\text{disk}} \frac{1}{r} dS' = \int_0^{R} \int_0^{2\pi} \frac{1}{r} r dr d\theta = 2\pi R$$

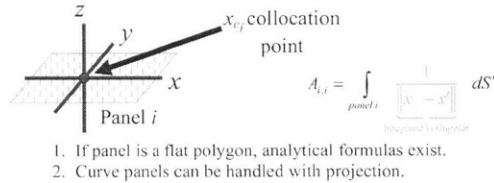
$$A_{6,6} = \int_{\text{Panel 6}} \frac{1}{\|x_{c_6} - x'\|} dS' = \int_{\text{disk}} \frac{1}{r} dS' + \int_{\text{rest}} \frac{1}{r} dS'$$

$$2\pi \cdot \frac{1}{2} = \int_{\text{disk}} \frac{1}{r} dS' \leq A_{6,6} \leq \int_{\text{large disk}} \frac{1}{r} dS' = 2\pi \cdot \frac{\sqrt{2}}{2}$$

≈ 3.8

Toeplitz Matrix

- Not diagonally dominant
- Dense
- All positive entries



SOLVING DISCRETIZED INTEGRAL EQUATIONS

- ❖ Basis Function – problem with dense matrices
- Integral Equation Method Generates Huge Dense Matrices

$$\begin{bmatrix} A_{1,1} & \dots & A_{1,m} \\ \vdots & \ddots & \vdots \\ A_{n,1} & \dots & A_{n,n} \end{bmatrix} \begin{bmatrix} \omega_1 \\ \vdots \\ \omega_n \end{bmatrix} = \begin{bmatrix} \Psi(x_{t_1}) \\ \vdots \\ \Psi(x_{t_n}) \end{bmatrix}$$

Gaussian Elimination Much Too Slow!

$$n = 100,000 = 10^5$$

Storage for an $n \times n$ matrix = $\underbrace{8 \text{ bytes}}_{\text{double precision}} \times \underbrace{10^{10}}_{n^2} = 80 \text{ gigabytes}$

❖ The Generalized Conjugate Residual Algorithm –

- o The k^{th} step of GCR

compute $\mathbf{A}r^{k-1}$

For discretized integral equations, \mathbf{A} is dense

$$\tilde{p}_{k-1} \leftarrow r^{k-1} - \sum_{j=0}^{k-2} (\mathbf{A}p_j)^T (\mathbf{A}r^{k-1}) p_j$$

\mathbf{A} - Orthonormalize search direction

$$p_{k-1} \leftarrow \frac{\tilde{p}_{k-1}}{\|\mathbf{A}\tilde{p}_{k-1}\|}$$

Determine optimal stepsize in k^{th} search direction

$$y_{k-1} \leftarrow (r^{k-1})^T \mathbf{A}p_{k-1}$$

Update the solution and the residual

$$x^k \leftarrow x^{k-1} + y_{k-1} p_{k-1}$$

$$r^k \leftarrow r^{k-1} - y_{k-1} \mathbf{A}p_{k-1}$$

- o Complexity of symmetric GCR (e.g. for Galerkin)

compute $\mathbf{A}r^{k-1}$

Dense matrix - vector product costs $O(n^2)$

$$\tilde{p}_{k-1} \leftarrow r^{k-1} - \sum_{j=0}^{k-2} (\mathbf{A}p_j)^T (\mathbf{A}r^{k-1}) p_j$$

inner product: $O(n)$

$$p_{k-1} \leftarrow \frac{\tilde{p}_{k-1}}{\|\mathbf{A}\tilde{p}_{k-1}\|}$$

inner product: $O(n)$

$$y_{k-1} \leftarrow (r^{k-1})^T \mathbf{A}p_{k-1}$$

$O(n)$ mult.

$$x^k \leftarrow x^{k-1} + y_{k-1} p_{k-1}$$

$$r^k \leftarrow r^{k-1} - y_{k-1} \mathbf{A}p_{k-1}$$

3D Finite Difference Method	100 × 100 × 100 grid
Number of unknowns = 10 ⁶	
Number of nonzeros = 7 · 10 ⁶ → 7 million	$O(n)$
3D Integral Equation Method	
Number of unknowns = 10 ⁶	
Number of nonzeros = 10 ¹² → 8,000 gigabytes	$O(n^2)$

- ❖ Fast Matrix Vector Products
 - Computing exactly A_r^{k-1}
 - Dense matrix-vector product costs $O(n^2)$
 - Computing approximately A_r^{k-1}
 - Reduces matrix-vector product costs to $O(n)$ or $O(n \log n)$



- Given a P.D.E. e.g. $\nabla^2 \Psi = 0$ $\nabla^2 \Psi + \kappa^2 \Psi = 0$
- 1.) Find a Green Function e.g. $G(x, x') = \frac{1}{|x - x'|}$ $G(x, x') = \frac{e^{i\kappa|x-x'|}}{|x - x'|}$
 - 2.) Discretize Surface Boundary Conditions Using Basis Functions

$$\sigma(x) = \sum_{j=1}^d w_j \varphi_j(x) \xrightarrow{\text{Piecewise Constant}} \Psi(x) = \sum_{j=1}^d w_j \int_{S_j} G(x, x') dS'$$

$$\varphi_j(x) = \begin{cases} 1 & x \in S_j \\ 0 & x \notin S_j \end{cases}$$
 - 3.) Impose Boundary Conditions

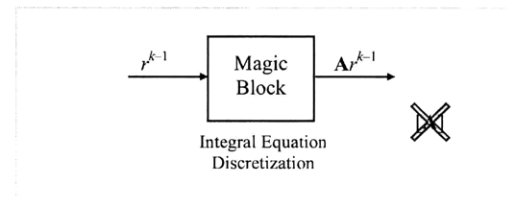
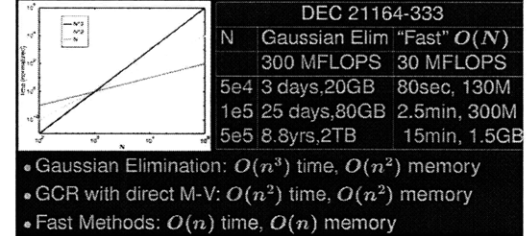
e.g. Collocation $\Psi(x_j) = \sum_{j=1}^N w_j \int_{S_j} G(x, x') dS'$

$$A \begin{bmatrix} w_1 \\ \vdots \\ w_N \end{bmatrix} = b$$

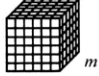
e.g. Galerkin $\int_{S_i} \Psi(x) dS = \sum_{j=1}^N w_j \int_{S_j} \int_{S_i} G(x, x') dS dS'$
 - 4.) Solve system using Krylov iterative methods at each step need A_r^{k-1}
 - 5.) Use Fast Matrix Vector Product for A_r^{k-1} (e.g. Fast Multipole or PFFT)
- } Cost $O(N \log N)$
 $O(N)$ memory

454

○ Computational Costs

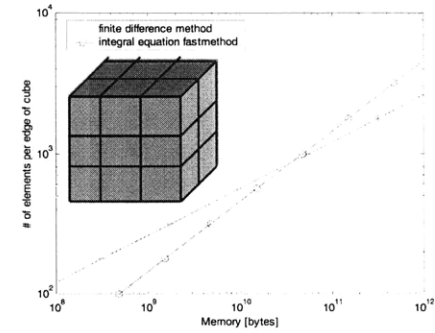


▪ Cube Example, Interior Problem – Memory Usage

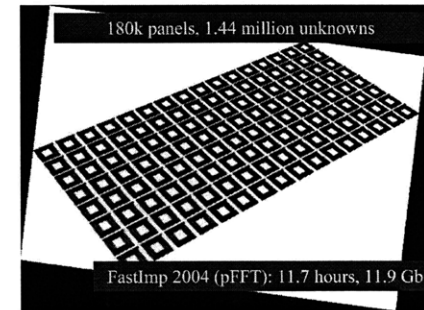
Solving a cube interior problem 

	F.D.			I.E. (BEM)		
	80's	2005	$O(\)$	80's	2005	$O(\)$
# panels on edge m	100	800×8		100	800×8	
# unknowns n	1M	$\times 512$	$O(m^3)$	60,000	$\times 64$	$O(6m^2)$
# nonzero entries	7M	$\times 512$	$O(7n) = O(7m^3)$	3.6 G	$\times 4096$	$O(n^2) = O(36m^4)$
memory	56 MB	$\times 512$ 30 GB	$O(56n) = O(56m^3)$	0.5 GB	$\times 64$ 120 TB	$O(8n^2) = O(290m^4)$ $O(8000n) = O(48000m^2)$

GCR + Fast Methods



▪ 128 3-turn spiral RF inductors array



$$\nabla^2 \Psi = 0 \Rightarrow \Psi(x) = \int G(x, x') \sigma(x') dS \Rightarrow \mathbf{A}\vec{x} = \vec{b}$$

- (1) Choice of basis functions
- (2) Choice of evaluation – Collocation or Galerkin

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 24.A.

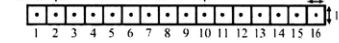
PDE – BEM Integral Equation Method III
Fast Methods for Integral Equation Solvers

TODAY'S OUTLINE:

- ❖ Fast Matrix-Vector Product
 - Multipole Algorithm (for Laplace Green Function)
 - Precorrected FFT (Green Function Independent)

FAST MATRIX-VECTOR PRODUCT

Example: Thin Metal Strip



Piecewise Constant Basis Function

$$\varphi_j(x') = \begin{cases} 1 & x' \in S_j \\ 0 & x' \notin S_j \end{cases}$$

Assume $\Psi(x) = 1$ given everywhere on the metal strip

Find $\sigma(x)$ on strip \leftrightarrow Find weights w_1, w_2, \dots, w_{16}
(charge on each panel) $A_{i,j} = \int_{\text{panel } j} \frac{1}{|x_{c_i} - x'|} dS'$

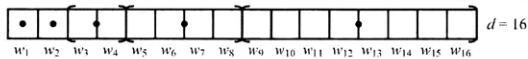
3.8	1	1/2	1/3	1/4	1/5	1/6	...	1/16
1
1/2	...	3.8	1	1/2	1/3	1/4	...	1/16
1/3

$$A_{i,j} \approx \frac{1}{|x_{c_i} - x_{c_j}|} \approx \frac{1}{|i - j|}$$

Area Panel 6 = $\frac{1}{5} = A_{6,1}$

$$\pi \leq A_{1,1} \leq \pi\sqrt{2} \quad A_{1,1} \approx 3.8$$

Structure: Diagonals (Toeplitz)



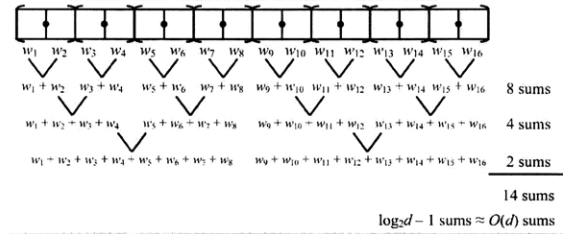
Matrix Vector Product $A \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_{16} \end{bmatrix} = \Psi$

$$\Psi(x_1) = 3.8w_1 + 1w_2 + \frac{1}{2}w_3 + \frac{1}{3}w_4 + \dots + \frac{1}{15}w_{16}$$

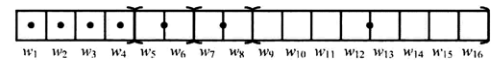
$$\Psi(x_1) = 3.8w_1 + \frac{w_2}{1} + \frac{w_3}{2.5} + \frac{w_4}{5.5} + \frac{w_5}{11.5} + \frac{w_6}{17.5} + \frac{w_7}{23.5} + \frac{w_8}{29.5} + \frac{w_9}{35.5} + \frac{w_{10}}{41.5} + \frac{w_{11}}{47.5} + \frac{w_{12}}{53.5} + \frac{w_{13}}{59.5} + \frac{w_{14}}{65.5} + \frac{w_{15}}{71.5} + \frac{w_{16}}{77.5}$$

4 products
4 sums
= $\log_2 d$

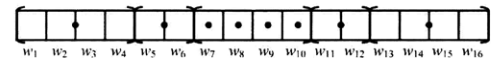
In order to do all the others $d \cdot \log_2 d \rightarrow$ cost of the entire matrix vector product provided j can get the clusters cheaply.



How do we use those sums?



$$\Psi(x_2) \approx \frac{1}{1}w_1 + 3.8w_2 + \frac{1}{1}w_3 + \frac{1}{2}w_4 + \frac{w_5 + w_6}{3.5} + \frac{w_7 + w_8}{5.5} + \frac{w_9 + w_{10} + w_{11} + \dots + w_{16}}{10.5}$$

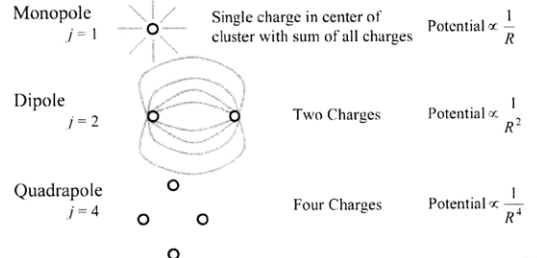


$$\Psi(x_8) \approx \frac{w_1 + \dots + w_4}{5.5} + \frac{w_5 + w_6}{2.5} + w_7 + 3.8w_8 + w_9 + \frac{1}{2}w_{10} + \frac{w_{11} + w_{12}}{3.5} + \frac{w_{13} + \dots + w_{16}}{6.5}$$

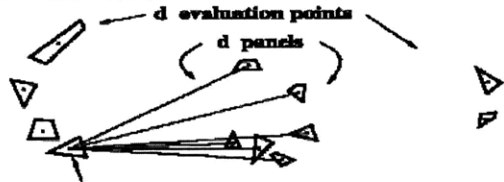
All are $O(\log_2 d) \rightarrow$ total cost $O(d \cdot \log_2 d)$

Multipole Algorithm (for Laplace Green Function)

❖ Basic Multipole Concepts – Multipole Representation

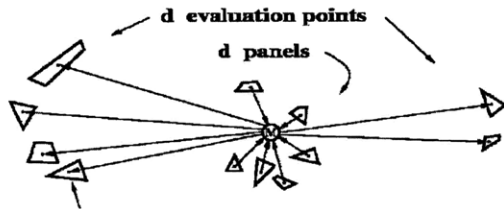


○ Direct Potential Evaluation



- Potential at point i :
$$v_i(r_i, \phi_i, \theta_i) = \sum_{j=1}^d q_j P_{ij}$$
- Complete evaluation at d points costs d^2 operations

○ Multipole Potential Evaluation



• Approximate potential at point i :

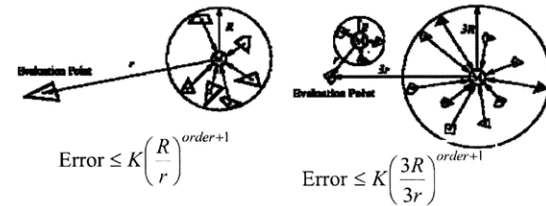
$$v_i(r_i, \phi_i, \theta_i) \approx \sum_{j=0}^{order} \sum_{k=-j}^j \frac{M_j^k}{r_i^{j+1}} Y_j^k(\phi_i, \theta_i)$$

• Multipole coefficients function of panel charges:

$$M_j^k = \sum_{i=1}^d \frac{q_i}{A_j} \int \rho^i Y_j^{-k}(\alpha, \beta) dA$$

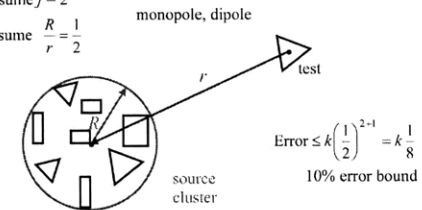
- Computing Multipole expansions costs order d operations.
- Each approximate potential evaluation costs order 1 operations. d potential evaluation due to d panels in order d operations

○ Scale Invariance of Error



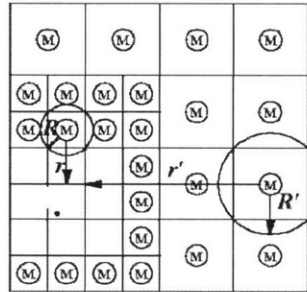
e.g. assume $j = 2$

assume $\frac{R}{r} = \frac{1}{2}$



Key Point: If test is much further away (r larger) then the clusters can be larger (R larger by same amount)

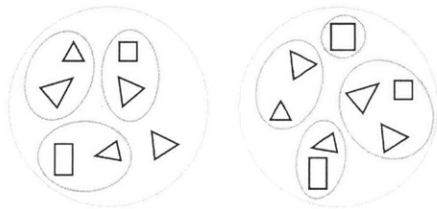
o Multipole Algorithm Hierarchy



Hierarchy guarantees:
 Bounded error:

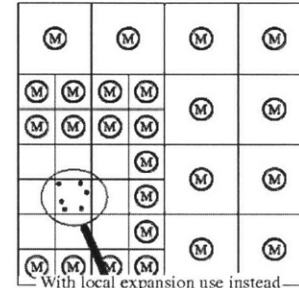
$$\text{Error} \leq K \left(\frac{R}{r}\right)^{\text{order}+1}$$

$$\leq K \left(\frac{1}{2}\right)^{\text{order}+1}$$
 order = 2 yields one percent accuracy.

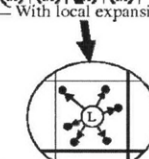


❖ Multipole Optimizations – Local Expansions

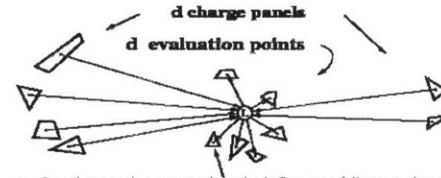
o Cost Reduction



- Construct a local expansion to represent distant charge potentials.
- Evaluate a single local expansion, rather than many multipole expansions, at each evaluation point.



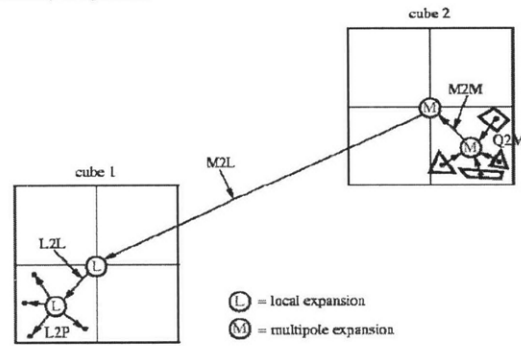
o Clustered Evaluations



- Local expansion summarizes the influence of distance charge or clusters of evaluation points.
- Gives $O(n)$ potential evaluation when combined with coalescing of charge done by multipole expansions.
- Approximate potential at point i :

$$v_i(r_i, \phi_i, \theta_i) \approx \sum_{j=0}^{\text{order}} \sum_{k=-j}^j L_j^k Y_j^k(\phi_i, \theta_i) r_i^j$$

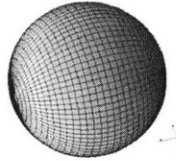
o Summary of Operations



- Multipole and local expansions are built using complementary hierarchies.
- Complete calculation consists of:
 1. Build multipoles (Upward Pass).
 2. Build locals (Downward Pass).
 3. Evaluation local expansions and nearby charge potential (Evaluation Pass).

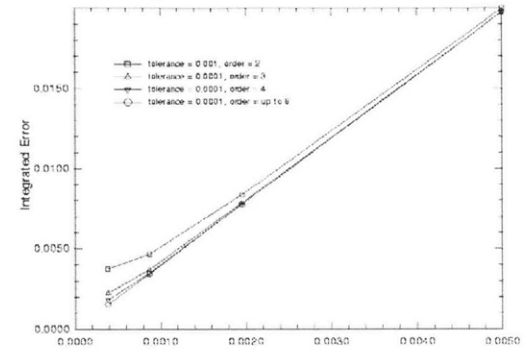
❖ Computational Examples

- o Translating Sphere
 - Potential Distribution

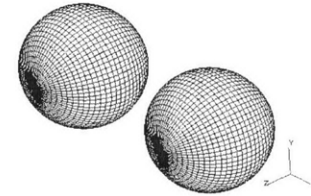


- Potential given by $\psi(x) = -\frac{x_3}{2|x|^3}$
- Charge given by $\sigma(x) = -\frac{3}{8\pi}x_3$

▪ Discretization Convergence

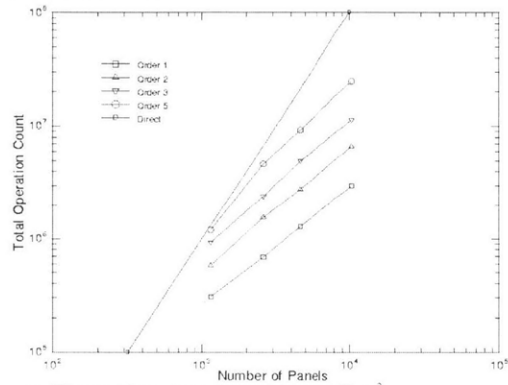


- Error should decay like $1/n$
- Multipole approximations eventually interfere
- Higher-order multipole expansions needed for higher accuracy
- o Two Sphere Example
 - Potential Distribution



- Potential on each sphere: $\psi(x) = -\frac{x_3}{2|x|^3}$
- Does not correspond to a simple physical problem.

Matrix-Vector Product Cost



- Direct matrix-vector product cost increases like n^2 .
- Multipole matrix-vector product cost increases like n .
- The slope for the multipole algorithm depends on accuracy.
- For order 2 expansions, breakpoint is about $n = 400$.

Complexity Summary

For an integral equation discretized with n panels:

- Gaussian elimination: $O(n^3)$
- GCR, direct M-V $O(n^2)$.
- Multipole accelerated GCR $O(mn)$.

$$\Psi(x_i) = \sum_{j=1}^n f(i-j)\sigma_j \quad \leftarrow \text{Discrete Convolution} \quad \text{Computed in } n \log n \text{ using FFT}$$

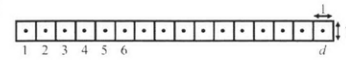
translation invariant

$\nabla^2 \Psi + \kappa^2 \Psi = 0$ Helmholtz Equation (Full Wave Equation)

$$G(x, x') = \frac{e^{i\kappa|x-x'|}}{|x-x'|}$$

Can still use multipole expansions but every time I change the Green function I need to come up with multipole expansions.
Fast Multipole is Green Function dependent!

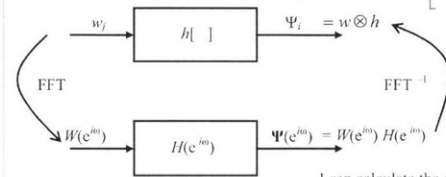
Example: Thin Metal Strip



$$\Psi(x_i) = \sum_{j=1}^d w_j A_{i,j} \quad A_{i,j} = \frac{\text{Area}}{|x_i - x_j|} = \frac{1}{|i-j|} = A_{i-j} = A_{j,i}$$

$$\Psi(x_i) = \sum_{j=1}^d w_j h_{i-j} = w \otimes h$$

Recognize Discrete Convolution Sum
 h is the unit sample response.



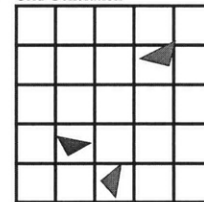
I can calculate the convolution efficiently using 2 FFT [$O(d \log_2 d)$]

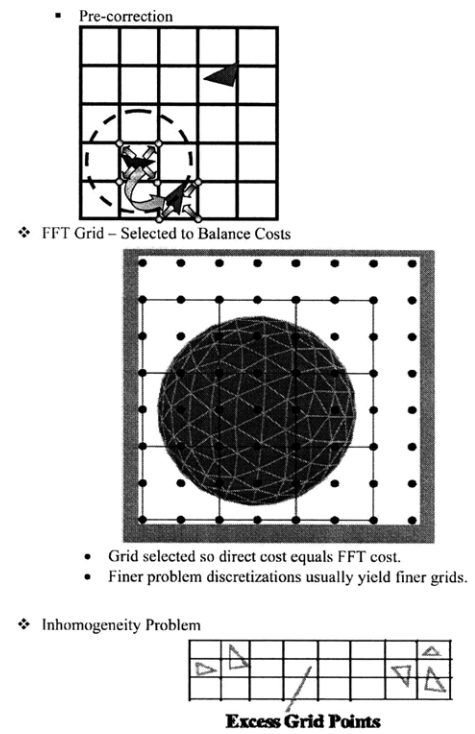
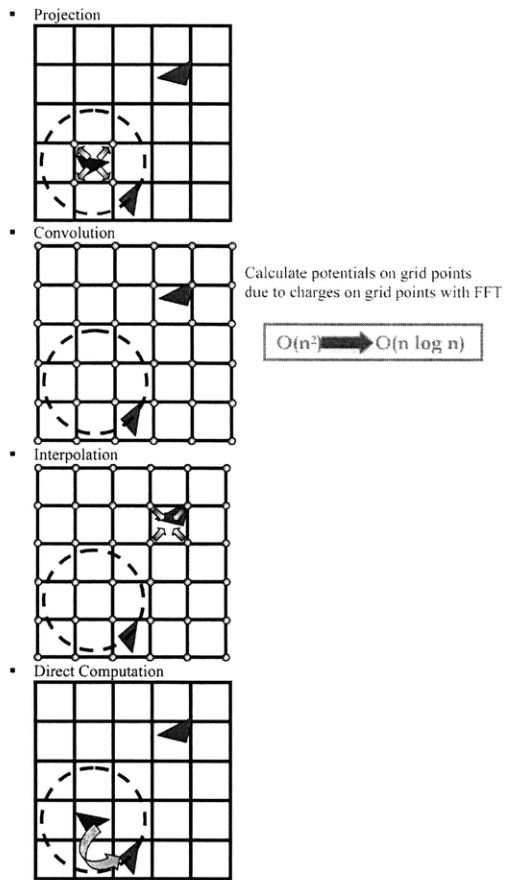
This works only if $G(x, x') = G(|x-x'|)$ translation invariant

Precorrected FFT (Green Function Independent)

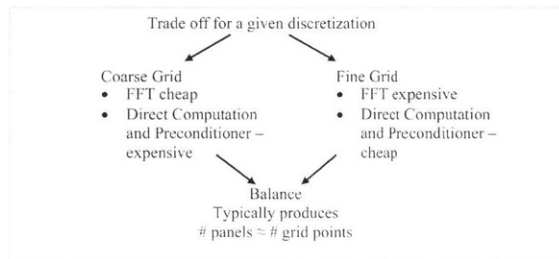
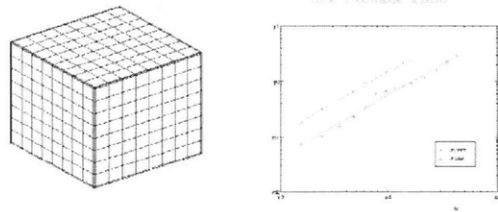
P-FFT Matrix Vector Product

- o Problem: Solve iteratively $\mathbf{A}\bar{\mathbf{w}} = \Psi$
- o At each iteration evaluate matrix-vector products $\mathbf{A}\bar{\mathbf{w}}$ using the following steps:
 - Grid Generation



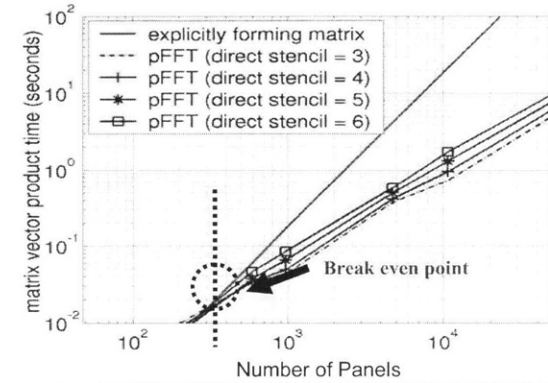


Inhomogeneity – Empty Grid due to FFT - Inefficiency
 Refining Cube Discretization – Worsening Inhomogeneity



462

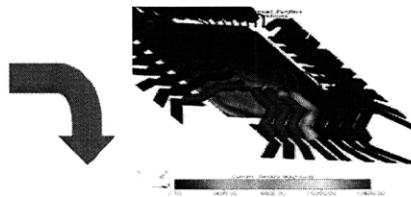
❖ Computational Complexity of PFFT++ is nearly $O(N)$
 PFFT++ is available at www.rle.mit.edu/cpg



EQS	FASTCAP	$\nabla^2\psi=0$	Capacitance
MQS	FASTHENRY	\rightarrow	L, R
<hr/>			
Fullwave	@UIUC		
Fast Multipole			
<hr/>			
pFFT++			
FASTIMP	EMQS + Fullwave impedance extraction		
FASTMAXWELL	EMQS + Fullwave impedance extraction		

❖ Example Maxwell Integral Equation Solver
Kernel Not 1/R Green Functions

$$\begin{aligned} \nabla \times E &= -\mu \frac{dH}{dt} \\ \nabla \times H &= \epsilon \frac{dE}{dt} \\ \nabla \cdot \epsilon E &= 0 \\ \nabla \cdot \mu H &= 0 \end{aligned}$$

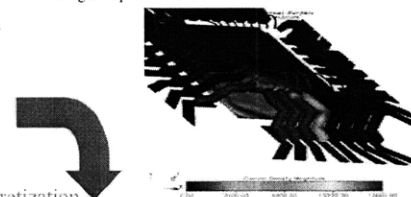


$$\frac{\mathbf{J}(\mathbf{r})}{\sigma} + j\omega \frac{\mu}{4\pi} \int \frac{\mathbf{J}(\mathbf{r}') e^{jk|\mathbf{r}-\mathbf{r}'|}}{|\mathbf{r}-\mathbf{r}'|} d\mathbf{r}' = -\nabla\phi$$

$$\frac{1}{4\pi\epsilon} \int \frac{\rho(\mathbf{r}') e^{jk|\mathbf{r}_s-\mathbf{r}_s'|}}{|\mathbf{r}_s-\mathbf{r}_s'|} d\mathbf{r}_s' = \phi(\mathbf{r}_s)$$

Preconditioner for Maxwell Integral Equation Solver

$$\begin{aligned} \nabla \times E &= -\mu \frac{dH}{dt} \\ \nabla \times H &= \epsilon \frac{dE}{dt} \\ \nabla \cdot \epsilon E &= 0 \\ \nabla \cdot \mu H &= 0 \end{aligned}$$



discretization

$$M \begin{bmatrix} L & 0 \\ 0 & P \end{bmatrix} M^T x = b$$

very sparse

very dense, symmetric diagonally dominant

❖ Preconditioning

- Krylov Methods – Diagonal Preconditioners

Let $A = D + A_{nd}$



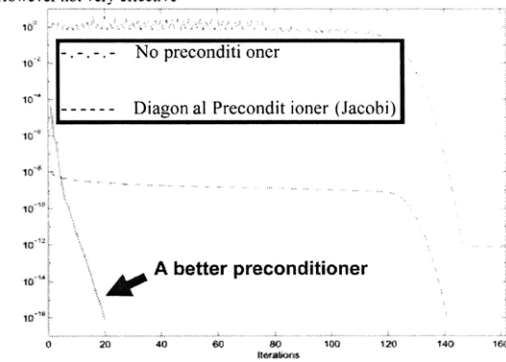
Apply GCR to $(D^{-1}A)x = (I + D^{-1}A_{nd})x = D^{-1}b$

- The inverse of a diagonal is cheap to compute
- Usually improves convergence

- Integral Equation Preconditioning
Diagonal Preconditioner (Jacobi)

$$P = \left(\text{diag} \left(M \begin{bmatrix} L & 0 \\ 0 & P \end{bmatrix} M^T \right) \right)^{-1}$$

Extremely easy to invert
However not very effective



Diagonal Preconditioner (Jacobi)

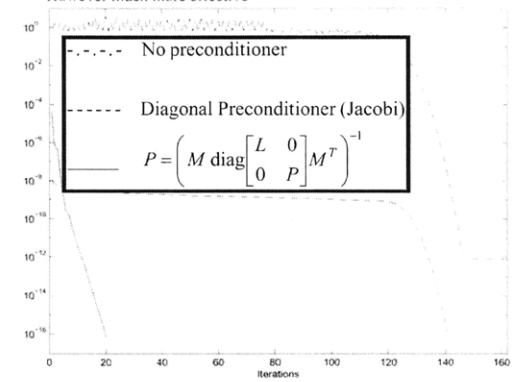
$$P = \left(\text{diag} \left(M \begin{bmatrix} L & 0 \\ 0 & P \end{bmatrix} M^T \right) \right)^{-1}$$

Extremely easy to invert
However not very effective

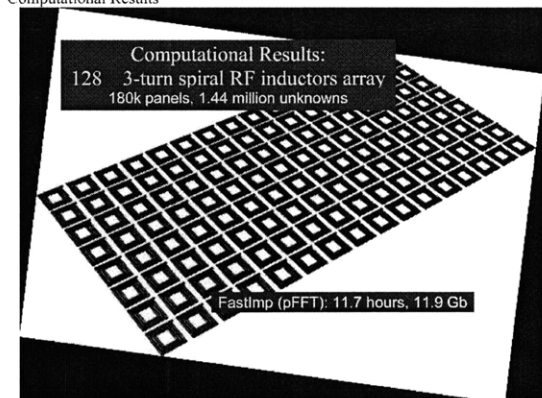
A Better Preconditioner:

$$P = \left(M \text{diag} \begin{bmatrix} L & 0 \\ 0 & P \end{bmatrix} M^T \right)^{-1}$$

A bit more difficult to factor (but still doable since M is sparse)
 However much more effective

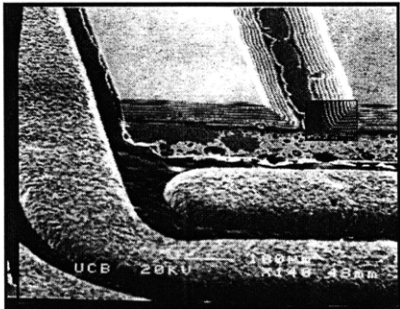
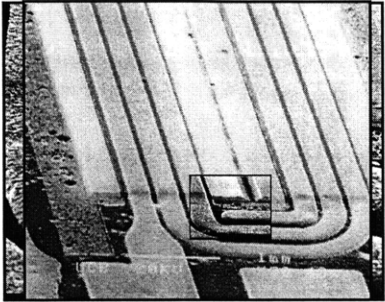


❖ Computational Results



3.24 Model Order Reduction I

2



INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 24.B.

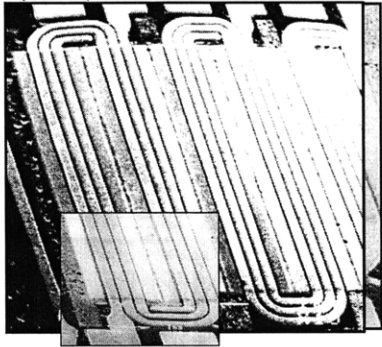
Model-Order Reduction

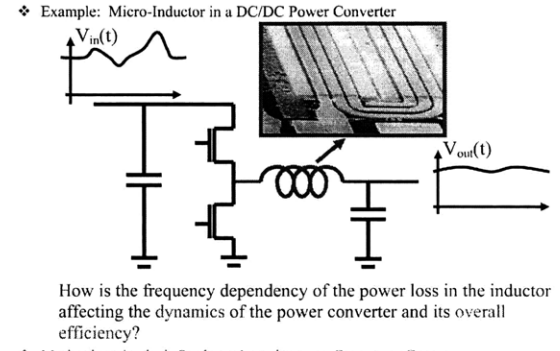
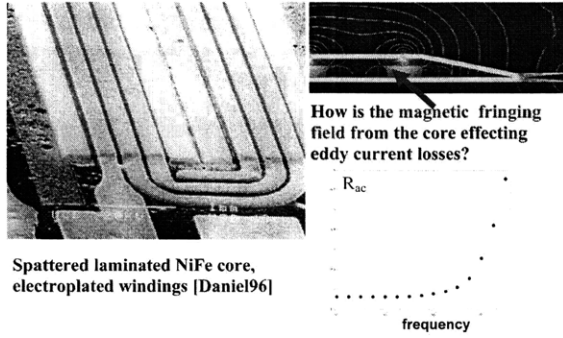
TODAY'S OUTLINE:

- ❖ Motivations and Examples
- ❖ From PDE to a State Space Model
 - Heat Conducting Bar

MOTIVATION AND EXAMPLES

- ❖ Analysis Example: Power Micro-Inductor



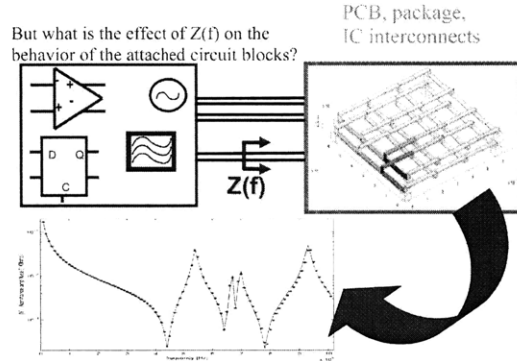


❖ Motivation: Analysis Produces Impedance vs. Frequency Curves

- How are parasitic and the resonances of the power distribution grid affecting the impedance?
- *Analysis tools* can produce for instance impedance vs. frequency curves.

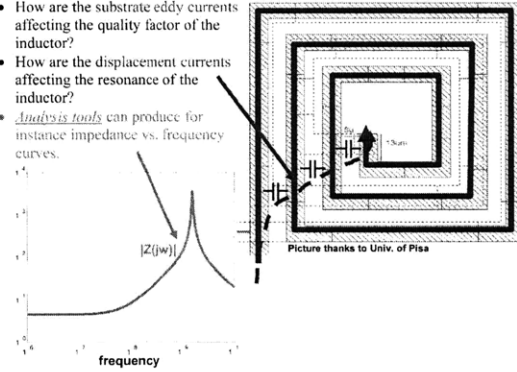
PCB, package, IC interconnects

❖ Example: Power Grid used Feeding Circuit Blocks



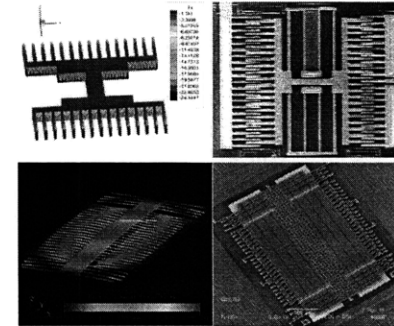
❖ Motivation. Example: RF Micro-Inductor

- How are the substrate eddy currents affecting the quality factor of the inductor?
- How are the displacement currents affecting the resonance of the inductor?
- *Analysis tools* can produce for instance impedance vs. frequency curves.



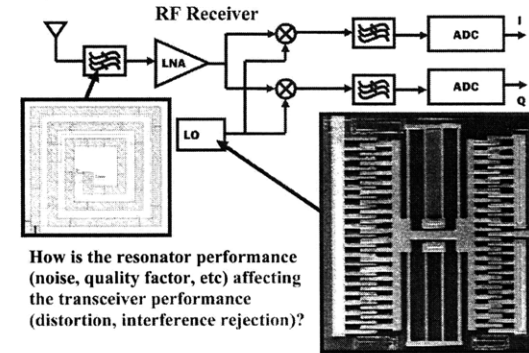
❖ Example: Accelerometer and RF Resonator

- What is the Drag force on the fingers of a resonator or accelerometer?
- How does it affect the quality factor?

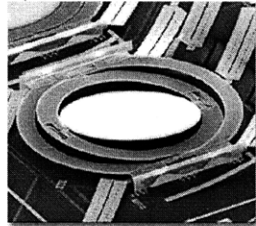


Pictures generated by FastStokes (Thanks to Xin Wang)

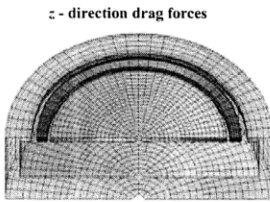
❖ Example: Micro-Inductor and Resonator in a Wireless Transceiver



❖ Example: Micro-Mirror



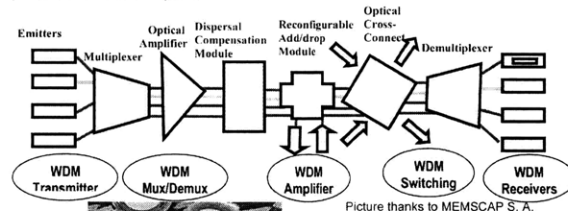
Picture thanks to Lucent



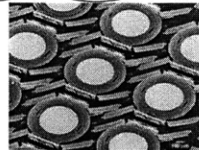
Picture by Xin Wang

What are the forces applied on the mirror, how do they affect the dynamic response of the mirror?

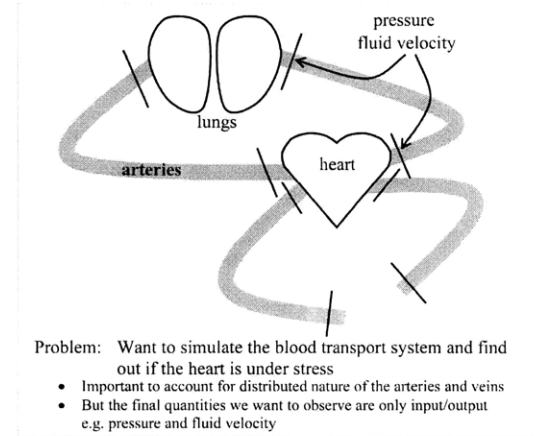
❖ Example: Micromirror Switch in a Dense Wavelength Division Multiplexing Optical Communication System



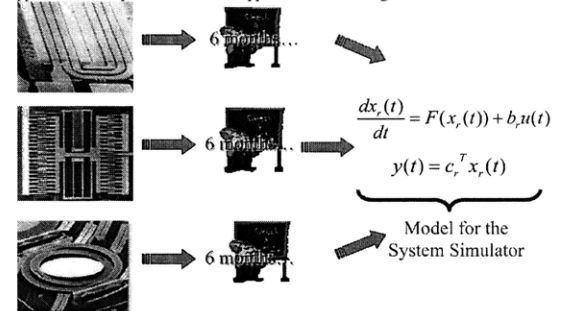
Picture thanks to MEMSCAP S. A.



How is the micromirror performance affecting the communication system functionality?

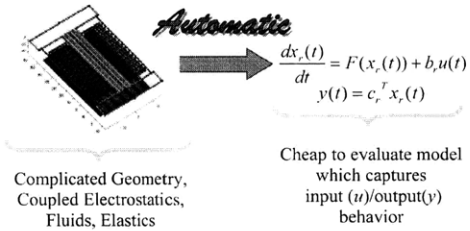


❖ Application Examples – Traditional Approach to Generating Models

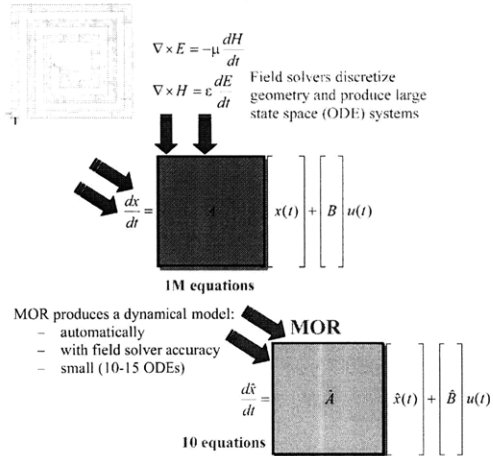


❖ The Numerical Macromodeling or Model Reduction Paradigm

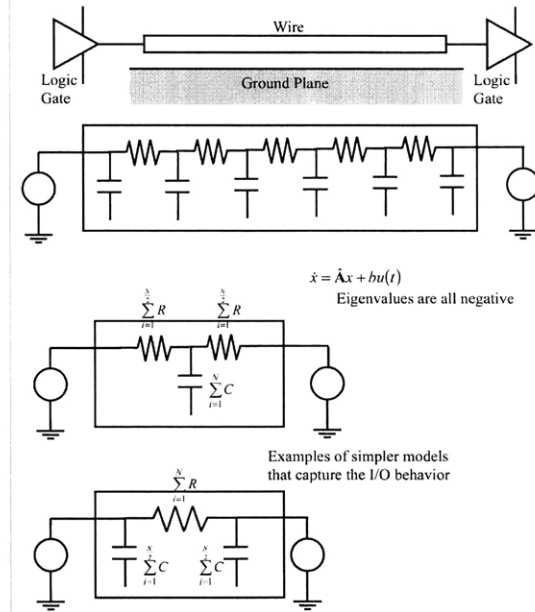
Generate a Reduced-Order Model Directly from 3-D Geometry and Physics



❖ From 3D geometry to small state space systems (MOR)



Signal propagation along a wire



FROM PDE TO A STATE SPACE MODEL

❖ Compare PDE Solvers and Model Order Reduction

- PDE Solvers:
 - Accurate
 - Relatively fast (minutes to hours)
 - Challenges:
 - very large matrices: sparse (FD or FEM), or dense (BEM)
 - Application: verification and characterization of component properties
- Model Order Reduction
 - Preserve PDE solver accuracy
 - Model construction relatively fast (same as PDE solvers minutes to hours)
 - Challenges: same as PDE solvers but only in model construction
 - Only capture input/output behavior (don't show field distributions)
 - But model evaluation in msec (e.g. get dynamical response from any input)
 - Application: analysis of functionality and interaction with other components

❖ Model Order Reduction State of the Art

Model Order Reduction for simple linear systems is well understood (e.g. interconnect, heat diffusion)

$$\frac{dx}{dt} = A x(t) + B u(t)$$

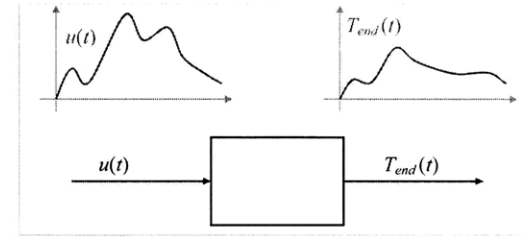
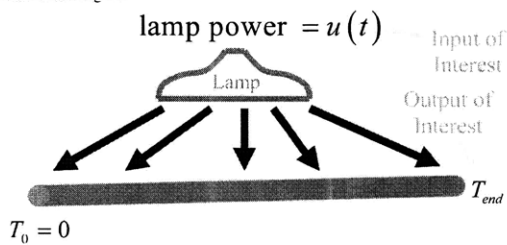
$$y(t) = C x(t)$$

Not many techniques yet for NON-LINEAR systems

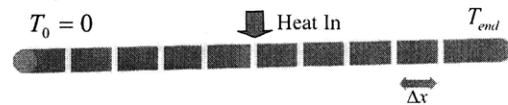
$$\frac{dx}{dt} = F[x(t)] + B u(t)$$

$$y(t) = C x(t)$$

❖ Heat Conducting Bar



o Basic Equations



Temperature Differential Equation

$$\underbrace{\gamma}_{\text{specific heat}} \frac{\partial T(x,t)}{\partial t} - \underbrace{\kappa}_{\text{thermal conductivity}} \frac{\partial^2 T(x,t)}{\partial x^2} = h(x) \underbrace{u(t)}_{\text{scalar input}}$$

Spatial Discretization (except at end)

$$\gamma \frac{d\hat{T}_i}{dt} - \frac{\kappa}{(\Delta x)^2} (\hat{T}_{i+1} - 2\hat{T}_i + \hat{T}_{i-1}) = h(x_i) u(t)$$

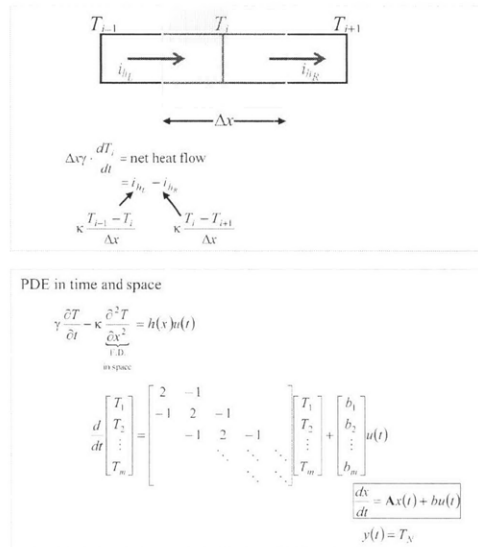
$\downarrow h(x) u(t)$

	$T(x)$	
--	--------	--

In steady-state $\sum h = 0$ Conservation Law $\Rightarrow -\kappa \frac{\partial^2 T}{\partial x^2} - h(x) = 0$

What if $\sum h \neq 0$? e.g. $h(x)u(t)$ increases suddenly

$-\kappa \frac{\partial^2 T}{\partial x^2} - h(x) = -\gamma \frac{\partial T}{\partial t}$ then the temperature will take some time to readjust



o Input-Output Discrete Equations

$$\gamma \frac{d\hat{T}_i}{dt} - \frac{\kappa}{(\Delta x)^2} (\hat{T}_{i+1} - 2\hat{T}_i + \hat{T}_{i-1}) = h(x_i)u(t) \quad i \in [1, \dots, N-1]$$

$$\gamma \frac{d\hat{T}_N}{dt} - \frac{\kappa}{(\Delta x)^2} (\hat{T}_N - \hat{T}_{N-1}) = h(x_N)u(t)$$

$$T_{end} = \hat{T}_N$$

o State-Space Description

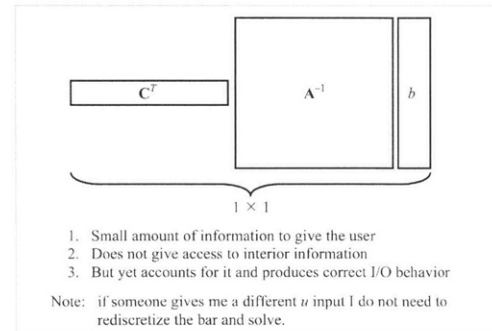
$$\frac{dx(t)}{dt} = \underbrace{A}_{N \times N} x(t) + \underbrace{b}_{N \times 1} \underbrace{u(t)}_{\text{scalar input}} \quad \underbrace{y(t)}_{\text{scalar output}} = \underbrace{c^T}_{1 \times N} x(t)$$

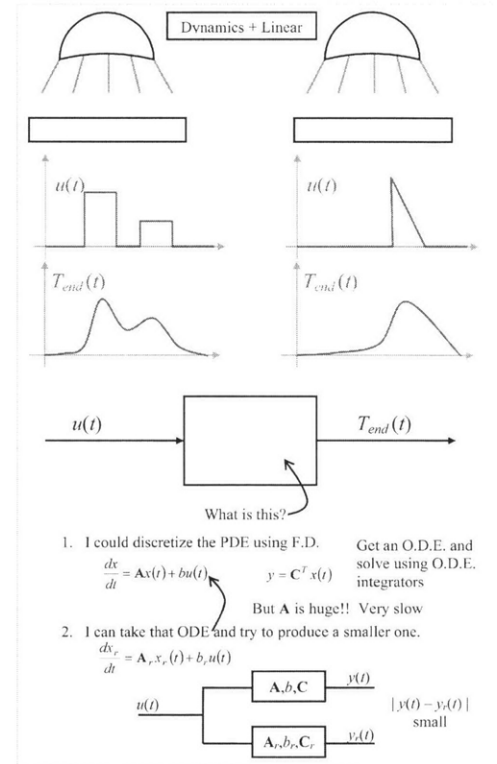
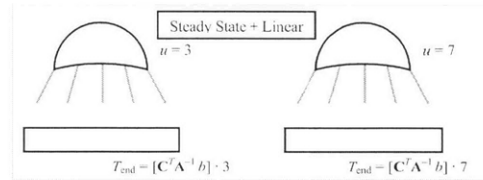
Given the right scaling

$$A = \begin{bmatrix} 2 & -1 & & \\ -1 & 2 & -1 & \\ & -1 & 2 & \ddots \\ & & \ddots & 2 & -1 \\ & & & -1 & 2 \end{bmatrix} \quad b = \begin{bmatrix} h(x_1) \\ h(x_2) \\ \vdots \\ h(x_N) \end{bmatrix} \quad c = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$$

$$\frac{d}{dt} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}_{\text{states}} = \begin{bmatrix} 2 & -1 & & \\ -1 & 2 & -1 & \\ & -1 & 2 & -1 \\ & & \ddots & \ddots \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} \underbrace{u(t)}_{\text{scalar input}}$$

$$\underbrace{y(t)}_{\text{scalar}} = [c_1 \quad c_2 \quad \dots \quad c_n] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \sum_{i=1}^n c_i x_i(t)$$





❖ Dynamic Linear Case – State-Space Description

Original Dynamical System - Single Input/Output

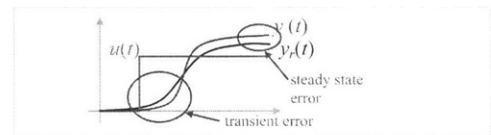
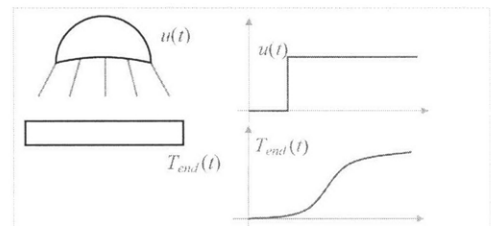
$$\frac{dx(t)}{dt} = \underbrace{A}_{N \times N} x(t) + \underbrace{b}_{N \times 1} \underbrace{u(t)}_{\text{scalar input}} \quad \underbrace{y(t)}_{\text{scalar}} = \underbrace{c^T}_{1 \times N} x(t)$$

Reduced Dynamical System

$$\frac{dx_r(t)}{dt} = \underbrace{A_r}_{q \times q} x_r(t) + \underbrace{b_r}_{q \times 1} \underbrace{u(t)}_{\text{scalar input}} \quad \underbrace{y_r(t)}_{\text{scalar}} = \underbrace{c_r^T}_{1 \times q} x_r(t)$$

$q \ll N$, but input/output behavior preserved

$$\bar{y} = \underbrace{[c_1 \ \dots \ c_N]}_{\text{scalar}} \underbrace{\begin{bmatrix} \mathbf{A}^{-1} \\ b_1 \\ \vdots \\ b_N \end{bmatrix}}_{\text{scalar}} \bar{u}$$



INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 25.

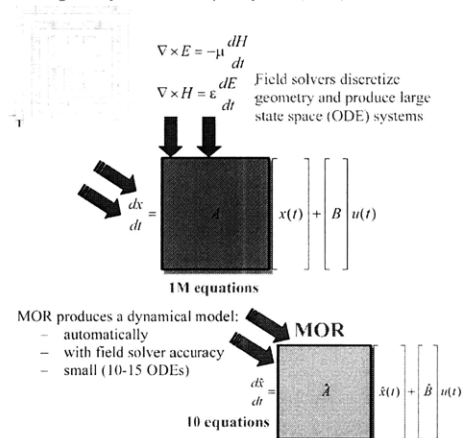
Model-Order Reduction II

TODAY'S OUTLINE:

- ❖ Problem Setup: from PDE to large ODE
- ❖ Reduction via eigenmode truncation method
- ❖ Reduction via transfer function fitting
 - Point Matching
 - Least Square
 - Quasi-convex Optimization Method
- ❖ Reduction via Projection Framework

FROM PDE TO LARGE ODE

- ❖ From 3D geometry to small state space systems (MOR)



❖ Dynamic Linear Case – State-Space Description

- o Original Dynamical System – Single Input/Output

$$\frac{dx(t)}{dt} = \underbrace{A}_{N \times N} x(t) + \underbrace{b}_{N \times 1 \text{ scalar}} u(t) \quad \underbrace{y(t)}_{\text{output}} = \underbrace{c^T}_{N \times 1 \text{ scalar}} x(t)$$

474

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 25.

Model-Order Reduction II

- o Reduced Dynamical System

$$\frac{d\hat{x}(t)}{dt} = \underbrace{\hat{A}}_{q \times q} \hat{x}(t) + \underbrace{\hat{b}}_{q \times 1 \text{ scalar}} u(t) \quad \underbrace{\hat{y}(t)}_{\text{output}} = \underbrace{\hat{c}^T}_{q \times 1 \text{ scalar}} \hat{x}(t)$$

- o $q \ll N$, but input/output behavior preserved

REDUCTION VIA EIGENMODE TRUNCATION METHOD

- ❖ Reminder about Eigenanalysis

Consider an ODE: $\frac{dx(t)}{dt} = Ax(t) + bu(t), \quad x(0) = 0$

Eigendecomposition:

$$A = \begin{bmatrix} \uparrow & \uparrow & \uparrow & \uparrow \\ E_1 & E_2 & \dots & E_N \\ \downarrow & \downarrow & & \downarrow \end{bmatrix} \begin{bmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \lambda_N & \\ & & & \ddots \end{bmatrix} \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ E_1 & E_2 & \dots & E_N \\ \downarrow & \downarrow & & \downarrow \end{bmatrix}^{-1}$$

Change of Variables: $Ew(t) = x(t) \Leftrightarrow w(t) = E^{-1}x(t)$

Substituting: $\frac{dEw(t)}{dt} = AEw(t) + bu(t), \quad Ew(0) = 0$

Multiply by E^{-1} : $\frac{dw(t)}{dt} = E^{-1}AEw(t) + E^{-1}bu(t)$.

$$\hat{E}_1 w_1(t) + \hat{E}_2 w_2(t) + \dots + \hat{E}_N w_N(t) = x(t)$$

Decoupled Equations

$$\frac{d}{dt} \begin{bmatrix} w_1 \\ \vdots \\ w_N \end{bmatrix} = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_N \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_N \end{bmatrix} + \begin{bmatrix} (E^{-1}b)_1 \\ \vdots \\ (E^{-1}b)_N \end{bmatrix} u(t)$$

Model-Order Reduction II

$$\frac{dw_i}{dt} = \lambda_i w_i(t) + \tilde{b}_i u(t)$$

how fast mode i responds

how much the input excites mode i
if $\tilde{b}_i = 0 \Rightarrow$ that mode is not controllable - eliminate it!

$$y(t) = \sum_{i=1}^n \tilde{c}_i w_i(t)$$

how much the mode i is seen at the output
if $\tilde{c}_i = 0 \Rightarrow$ the mode is not observable - eliminate it!

Output Equation

$$y(t) = c^T x(t) = c^T E w(t) = \underbrace{(E^T c^T)}_c w(t)$$

Solving Decoupled Equations

$$w_i(t) = \int_0^t e^{\lambda_i(t-\tau)} \tilde{b}_i u(\tau) d\tau \quad \text{Assuming Zero Initial Conditions}$$

Output Equation

$$y(t) = \sum_{i=1}^n \tilde{c}_i w_i(t)$$

$$\begin{aligned} \frac{d}{dt} w_1 &= \lambda_1 w_1 + \tilde{b}_1 u(t) & y(t) &= \sum \tilde{c}_i w_i(t) \\ &\vdots & & \\ \frac{d}{dt} w_N &= \lambda_N w_N + \tilde{b}_N u(t) & \tilde{c}_i &= 0 \quad i^{\text{th}} \text{ node unobservable!} \\ & & \tilde{b}_i &= 0 \quad i^{\text{th}} \text{ node uncontrollable!} \end{aligned}$$

❖ Dynamic Linear Case – Reduced Models via Mode Truncation

$$\begin{bmatrix} \dot{w}_1 \\ \vdots \\ \dot{w}_q \\ \vdots \\ \dot{w}_N \end{bmatrix} = \begin{bmatrix} \lambda_1 & & & & \\ & \ddots & & & \\ & & \lambda_q & & \\ & & & \ddots & \\ & & & & \lambda_N \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_q \\ \vdots \\ w_N \end{bmatrix} + \begin{bmatrix} \tilde{b}_1 \\ \vdots \\ \tilde{b}_q \\ \vdots \\ \tilde{b}_N \end{bmatrix} u(t)$$

$$\begin{bmatrix} \dot{w}_1 \\ \vdots \\ \dot{w}_q \end{bmatrix} = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_q \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_q \end{bmatrix} + \begin{bmatrix} \tilde{b}_1 \\ \vdots \\ \tilde{b}_q \end{bmatrix} u(t)$$

475

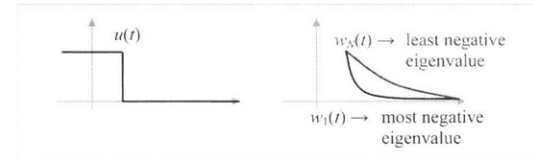
Model-Order Reduction II

$$y(t) = \sum_{i=1}^n \tilde{c}_i w_i(t) \quad \text{Output Equation}$$

- Certain modes are not affected by the input
 $\tilde{b}_{k+1}, \dots, \tilde{b}_N$ are all small
- Certain modes do not affect the output
 $\tilde{c}_{k+1}, \dots, \tilde{c}_N$ are all small
- Keep least negative eigenvalues (slowest modes)
– Look at response to a constant input

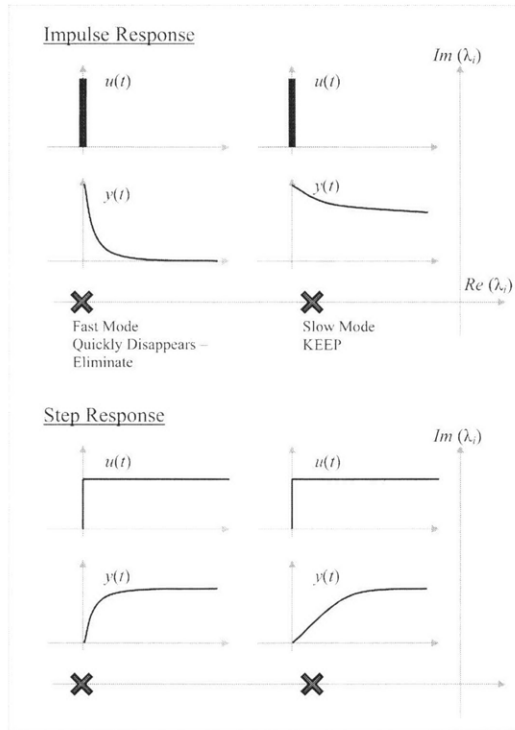
$$w_i(t) = \int_0^t e^{\lambda_i(t-\tau)} \tilde{b}_i u(\tau) d\tau = \frac{1}{\lambda_i} (\tilde{b}_i u - \tilde{b}_i u e^{\lambda_i t})$$

Small if $|\lambda_i|$ large



INTRODUCTION TO NUMERICAL SIMULATION

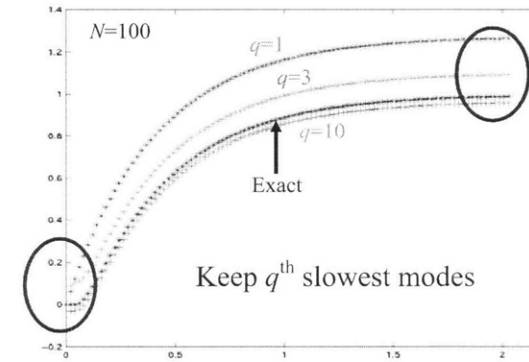
LECTURE 25.
Model-Order Reduction II



476

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 25.
Model-Order Reduction II



Final Value Theorem

For accuracy at steady state of step response need good accuracy for $s \rightarrow 0$:

$$\lim_{s \rightarrow 0} H(s) = \lim_{s \rightarrow 0} \sum_{i=1}^n \frac{\tilde{c}_i \tilde{b}_i}{s - \lambda_i} = \sum_{i=1}^n \frac{\tilde{c}_i \tilde{b}_i}{\lambda_i}$$

eliminate terms where $\frac{\tilde{c}_i \tilde{b}_i}{\lambda_i}$ is small

$$H(s) = \frac{\tilde{c}_1 \tilde{b}_1}{s - \lambda_1} + \dots + \frac{\tilde{c}_N \tilde{b}_N}{s - \lambda_N}$$

❖ Transfer Functions

○ Laplace Transform

Consider an ODE: $\frac{dx(t)}{dt} = Ax(t) + bu(t)$

Bilateral Laplace Transform: $X(s) = \int_{-\infty}^{\infty} x(t)e^{-st} dt$

Key Transform Property: $sX(s) = \int_{-\infty}^{\infty} \frac{dx(t)}{dt} e^{-st} dt$

INTRODUCTION TO NUMERICAL SIMULATION

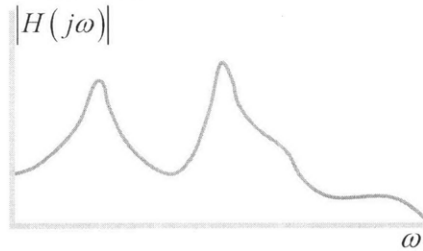
LECTURE 25.

Model-Order Reduction II

Rewrite the ODE in transformed variables
 $sX(s) = \mathbf{A}X(s) + bU(s) \quad Y(s) = c^T X(s)$
 $\Rightarrow Y(s) = \underbrace{c^T (sI - \mathbf{A})^{-1} b}_{H(s)} U(s)$

$$\begin{aligned} sX(s) &= \mathbf{A}X(s) + bU(s) \\ (sI - \mathbf{A})X(s) &= bU(s) \end{aligned}$$

- Meaning of H(s)
 For Stable Systems, $H(j\omega)$ is the frequency response
 If $u(t) = e^{j\omega t}$ ← sinusoid
 then $y(t) = H(j\omega)e^{j\omega t}$ ← sinusoid with shifted phase and amplitude



- EigenAnalysis
 Transfer Function
 $H(s) = c^T (sI - \mathbf{A})^{-1} b$
 Apply Eigendecomposition $\mathbf{A} = \mathbf{E}\mathbf{\Lambda}\mathbf{E}^{-1}$
 $H(s) = c^T \mathbf{E}(sI - \mathbf{A})^{-1} \mathbf{E}^{-1} b$
 $= \tilde{c}^T \begin{bmatrix} \frac{1}{s - \lambda_1} & & \\ & \ddots & \\ & & \frac{1}{s - \lambda_N} \end{bmatrix} \tilde{b} \Rightarrow H(s) = \sum_{i=1}^n \frac{\tilde{c}_i \tilde{b}_i}{s - \lambda_i}$
 Eliminate each mode for which this term is small.

477

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 25.

Model-Order Reduction II

$$H(s) = \frac{\tilde{b}_1 \tilde{c}_1}{s - \lambda_1} + \frac{\tilde{b}_2 \tilde{c}_2}{s - \lambda_2} + \dots + \frac{\tilde{b}_N \tilde{c}_N}{s - \lambda_N}$$

leave out $w_i(t)$ or leave out $\tilde{b}_i \tilde{c}_i$ from $H(s)$

$$\begin{aligned} \mathbf{A} &= \mathbf{E}\mathbf{\Lambda}\mathbf{E}^{-1} \\ H(s) &= \mathbf{C}^T (s\mathbf{I} - \mathbf{E}\mathbf{\Lambda}\mathbf{E}^{-1})^{-1} b \\ &= \mathbf{C}^T (s\mathbf{E}\mathbf{E}^{-1} - \mathbf{E}\mathbf{\Lambda}\mathbf{E}^{-1})^{-1} b \\ &= \mathbf{C}^T \mathbf{E}(s\mathbf{I} - \mathbf{\Lambda})^{-1} \mathbf{E}^{-1} b = \left[\mathbf{E}^T \mathbf{C}^T \right] \begin{bmatrix} \frac{1}{s - \lambda_1} & & \\ & \ddots & \\ & & \frac{1}{s - \lambda_n} \end{bmatrix} \begin{bmatrix} \mathbf{E}^{-1} b \\ \vdots \\ \mathbf{E}^{-1} b \end{bmatrix} \\ &= \sum_{i=1}^n \frac{[\mathbf{C}^T \mathbf{E}] [\mathbf{E}^{-1} b]}{s - \lambda_i} \quad \text{Pole-Residue Form} \end{aligned}$$

$$\mathbf{A} = \mathbf{E} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \mathbf{E}^{-1} \quad \leftarrow \text{Assume } \mathbf{A} \text{ diagonalizable}$$

$$\frac{\mathbf{C}^T \mathbf{E}}{\tilde{c}} \left(sI - \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \right)^{-1} \frac{\mathbf{E}^{-1} b}{\tilde{b}} = \frac{\tilde{b}_1 \tilde{c}_1}{s - \lambda_1} + \frac{\tilde{b}_2 \tilde{c}_2}{s - \lambda_2} + \dots$$

- Model Order Reduction via Eignmode Analysis
 Pole - Residue Form Pole - Zero Form (SISO)
 $H(s) = \sum_{i=1}^n \frac{\tilde{c}_i \tilde{b}_i}{s - \lambda_i} \quad H(s) = \prod_{i=1}^{n-1} (s - \zeta_i)$
 $h(t) = \sum_{i=1}^n \tilde{c}_i \tilde{b}_i e^{\lambda_i t} \quad \prod_{i=1}^n (s - \lambda_i)$
 Ideas for reducing order:
 - Drop terms with small residues $\tilde{c}_i \tilde{b}_i$
 - Drop terms with large negative $\text{Re}(\lambda_i)$ ("fast" modes)
 - Remove pole/zero near-cancellations
 - Cluster poles that are "together"

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 25.

Model-Order Reduction II

REDUCTION VIA TRANSFER FUNCTION FITTING

o Counting Degrees of Freedom

Dynamical System

$$\begin{aligned} \frac{d\hat{x}}{dt} &= \hat{A}\hat{x}(t) + \hat{b}u(t) \\ \hat{y}(t) &= \hat{c}^T \hat{x}(t) \end{aligned}$$

Degrees of Freedom

$$q^2 + 2q \text{ coefficients}$$

Transfer Function

$$\begin{aligned} H(s) &= \frac{r_1}{(s-\lambda_1)} + \dots + \frac{r_n}{(s-\lambda_n)} \\ &= k \frac{(s-\zeta_1) \dots (s-\zeta_{n-1})}{(s-\lambda_1) \dots (s-\lambda_n)} \\ &= \frac{\hat{b}_0 + \hat{b}_1 s + \dots + \hat{b}_{q-1} s^{q-1}}{1 + \hat{a}_1 s + \dots + \hat{a}_q s^q} \end{aligned}$$

Degrees of Freedom

$$2q \text{ coefficients}$$

o Fully Invertible Change of Coordinates

Reduced Model Transfer Function

$$\begin{aligned} \frac{d\hat{x}}{dt} &= \hat{A}\hat{x}(t) + \hat{b}u(t) \Rightarrow H(s) = \hat{c}^T (sI - \hat{A})^{-1} \hat{b} \\ \hat{y}(t) &= \hat{c}^T \hat{x}(t) \end{aligned}$$

Apply any invertible change of coordinates to the state $\hat{x}(t) = U\bar{x}(t)$

$$\begin{aligned} U^{-1}U \frac{d\bar{x}}{dt} &= U^{-1}\hat{A}U\bar{x}(t) + U^{-1}\hat{b}u(t) \Rightarrow \bar{H}(s) = \hat{c}^T (sI - \hat{A})^{-1} \hat{b} = H(s) \\ \bar{y}(t) &= \hat{c}^T U\bar{x}(t) \end{aligned}$$

Many Dynamical System have the same transfer function!!!

$$\begin{aligned} \bar{H}(s) &= C^T U (sI - U^{-1}AU)^{-1} U^{-1}b \\ &= C^T U U^{-1} (sI - A)^{-1} U U^{-1} b \\ &= C^T (sI - A)^{-1} b \end{aligned} \quad U = \begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ u_1 & u_2 & \dots & u_q \\ \downarrow & \downarrow & \dots & \downarrow \end{bmatrix} \times$$

o Model Order Reduction via Rational Transfer Function Fitting

Original System Transfer Function:

$$H(s) = \frac{b_0 + b_1 s + \dots + b_{N-1} s^{N-1}}{1 + a_1 s + \dots + a_N s^N} \leftarrow \text{rational function}$$

478

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 25.

Model-Order Reduction II

Model Reduction = Find a low order ($q \ll N$) rational function matching

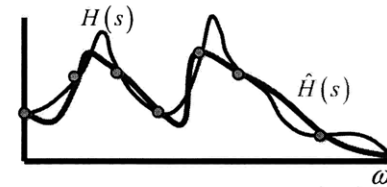
$$\hat{H}(s) = \frac{\hat{b}_0 + \hat{b}_1 s + \dots + \hat{b}_{q-1} s^{q-1}}{1 + \hat{a}_1 s + \dots + \hat{a}_q s^q} \leftarrow \text{reduced order rational function}$$

$$\begin{aligned} (sI - A)X(s) &= bU(s) \\ X(s) &= (sI - A)^{-1} bU(s) \\ Y(s) &= \underbrace{C^T (sI - A)^{-1} b}_{H(s)} U(s) \end{aligned}$$

o Point Matching

o Rational Transfer Function Fitting: via Point Matching

$$\underbrace{H(s_i)}_{\text{original}} = \underbrace{\frac{b'_0 + b'_1 s_i + b'_2 s_i^2 + \dots + b'_{q-1} s_i^{q-1}}{1 + a'_1 s_i + \dots + a'_q s_i^q}}_{\hat{H}(s)}$$



• Can match 2q points $\hat{H}(s_i) = \frac{\hat{b}_0 + \hat{b}_1 s_i + \dots + \hat{b}_{q-1} s_i^{q-1}}{1 + \hat{a}_1 s_i + \dots + \hat{a}_q s_i^q}$

• cross multiplying generates a linear system

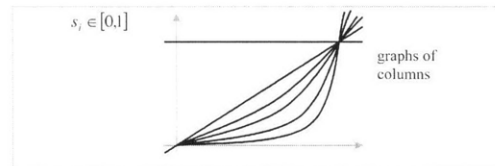
$$\text{For } i = 1 \text{ to } 2q \quad (1 + \hat{a}_1 s_i + \dots + \hat{a}_q s_i^q) \hat{H}(s_i) - (\hat{b}_0 + \hat{b}_1 s_i + \dots + \hat{b}_{q-1} s_i^{q-1}) = 0$$

- Point Matching can be ill-conditioned

$$\begin{bmatrix} s_1 H(s_1) & \dots & s_1^q H(s_1) & -1 & -s_1 & \dots & -s_1^{q-1} & \vdots \\ s_2 H(s_2) & \dots & s_2^q H(s_2) & -1 & -s_2 & \dots & -s_2^{q-1} & \vdots \\ \vdots & & \vdots & & \vdots & & \vdots & \vdots \\ s_{2q} H(s_{2q}) & \dots & s_{2q}^q H(s_{2q}) & -1 & -s_{2q} & \dots & -s_{2q}^{q-1} & \vdots \end{bmatrix} \begin{bmatrix} a_1 \\ \vdots \\ a_q \\ b_0 \\ b_1 \\ \vdots \\ b_{q-1} \end{bmatrix} = \begin{bmatrix} H(s_1) \\ H(s_2) \\ \vdots \\ H(s_{2q}) \end{bmatrix}$$

- Columns contain progressively higher powers of the test frequencies: problem is numerically ill-conditioned
- Also... missing data can cause severe accuracy problems

$$s_1 H(s_1) a_1 + \dots + s_1^q H(s_1) a_q - (b_0 + b_1 s_1 + \dots + b_{q-1} s_1^{q-1}) = 0$$

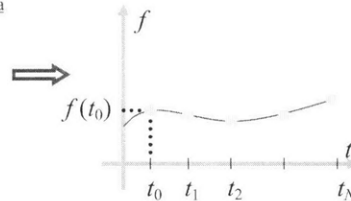


- Hard to Solve Systems – Fitting Example

Polynomial Interpolation

Table of Data

t_0	$f(t_0)$
t_1	$f(t_1)$
\vdots	\vdots
t_N	$f(t_N)$



Problem fit data with an N^{th} order polynomial

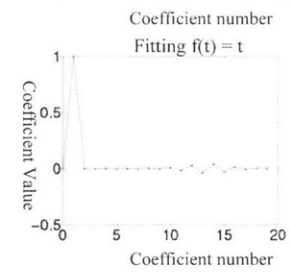
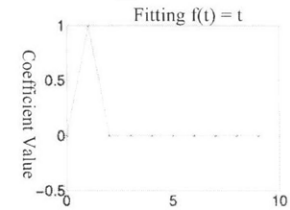
$$f(t) = \alpha_0 + \alpha_1 t + \alpha_2 t^2 + \dots + \alpha_N t^N$$

479

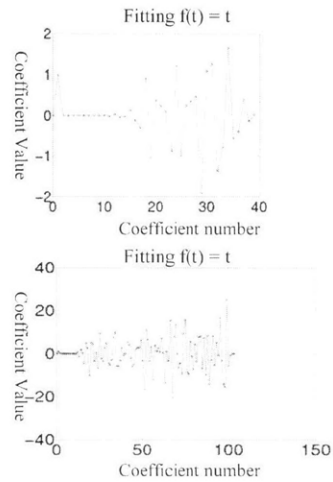
Matrix Form

$$\begin{bmatrix} 1 & t_0 & t_0^2 & \dots & t_0^N \\ 1 & t_1 & t_1^2 & \dots & t_1^N \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & t_N & t_N^2 & \dots & t_N^N \end{bmatrix} \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_N \end{bmatrix} = \begin{bmatrix} f(t_0) \\ f(t_1) \\ \vdots \\ f(t_N) \end{bmatrix}$$

M_{interp}



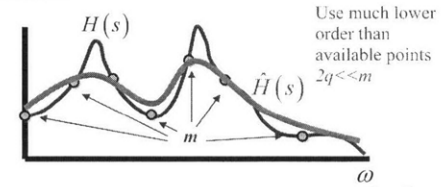
LECTURE 25.
Model-Order Reduction II



480

LECTURE 25.
Model-Order Reduction II

❖ Least Square



$$\begin{bmatrix} s_1 H(s_1) & \cdots & s_1^q H(s_1) & -1 & -s_1 & \cdots & -s_1^{q-1} \\ s_2 H(s_2) & \cdots & s_2^q H(s_2) & -1 & -s_2 & \cdots & -s_2^{q-1} \\ \vdots & & \vdots & & \vdots & & \vdots \\ s_m H(s_m) & \cdots & s_m^q H(s_m) & -1 & -s_m & \cdots & -s_m^{q-1} \end{bmatrix} \begin{bmatrix} a_1 \\ \vdots \\ a_q \\ b_0 \\ b_1 \\ \vdots \\ b_{q-1} \end{bmatrix} = \begin{bmatrix} H(s_1) \\ H(s_2) \\ \vdots \\ H(s_m) \end{bmatrix}$$

- Cross-multiplying generates a linear TALL SKINNY system
- Use, for instance, QR to solve it or Gauss-Newton

❖ Quasi-Convex Optimization Method
o Optimization Based Rational Fit Model Order Reduction Setup

$$\underset{p(s), q(s)}{\text{minimize}} \left\| H(s) - \frac{p(s)}{q(s)} \right\|$$

from field solver
OR measurements
small stable and passive
reduced order model

Least Square Method

- Cast as nonlinear least squares (solved by e.g. Gauss-Newton)
- Do not consider stability or passivity while finding poles (need post-processing)

Quasi Convex Method

- Cast as quasi-convex program (solved by convex optimization algorithm)
- Explicitly take care of stability and passivity while finding poles

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 25.

Model-Order Reduction II

- o Relaxation of the H-inf norm MOR setup [Sou, Megretski, Daniel]

$$\text{minimize}_{p,q,r} \left\| H(z) - \frac{p(z)}{q(z)} - \frac{r(z^{-1})}{q(z^{-1})} \right\|_{\infty}$$

Anti-stable term

$$\text{subject to} \quad \deg(q) = m, \quad \deg(p) \leq m,$$

$$\deg(r) < m$$

Stability: $q(z)$ Schur polynomial (roots inside unit circle)

Passivity, and possibly other constraints

Benefit: Relaxation *equivalent* to a quasi-convex program

Drawback: May obtain suboptimal solutions

- o How bad is our relaxation?

THEOREM:

$$\text{Let } (q^*, p^*, r^*) = \arg \min_{q,p,r} \left\| H(z) - \frac{p(z)}{q(z)} - \frac{r(z^{-1})}{q(z^{-1})} \right\|_{\infty}$$

Such that $\deg(q) = m$, $q(z)$ is Schur polynomial

$$\text{Then } \left\| H(z) - \frac{p(z^*)}{q(z^*)} \right\|_{\infty} < m \sigma_{m+1}(H)$$

$m+1^{\text{th}}$ Hankel singular value

- o Equivalent Quasi-Convex Setup

$$\text{minimize}_{a,b,c} \left\| H(e^{j\omega}) - \frac{b(e^{j\omega}) + jc(e^{j\omega})}{a(e^{j\omega})} \right\|$$

quasi-convex function

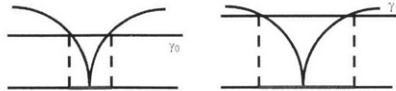
$$\deg(a) = m, \quad \deg(b) \leq m,$$

$$\deg(c) \leq m,$$

subject to Stability: $a(e^{j\omega}) > 0, \forall \omega \in [0, \pi]$

Passivity: $b(e^{j\omega}) > 0, \forall \omega \in [0, \pi]$

convex set



481

INTRODUCTION TO NUMERICAL SIMULATION

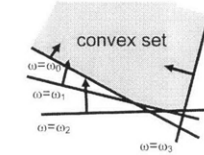
LECTURE 25.

Model-Order Reduction II

This is a *quasi-convex* program, because

$$a(e^{j\omega}) = 2 \cos(m\omega) + 2 \cos((m-1)\omega)a_{m-1} + \dots + a_0 > 0$$

defines an intersection of halfspaces



This is a *quasi-convex* program, because

$$2 \cos \left(\begin{bmatrix} m\omega_0 & (m-1)\omega_0 & \dots & 0 \end{bmatrix} \begin{bmatrix} 1 & a_{m-1} & \dots & a_0 \end{bmatrix} \right) > 0$$

$$2 \cos \left(\begin{bmatrix} m\omega_1 & (m-1)\omega_1 & \dots & 0 \end{bmatrix} \begin{bmatrix} 1 & a_{m-1} & \dots & a_0 \end{bmatrix} \right) > 0$$

$$2 \cos \left(\begin{bmatrix} m\omega_2 & (m-1)\omega_2 & \dots & 0 \end{bmatrix} \begin{bmatrix} 1 & a_{m-1} & \dots & a_0 \end{bmatrix} \right) > 0$$

$$2 \cos \left(\begin{bmatrix} m\omega_3 & (m-1)\omega_3 & \dots & 0 \end{bmatrix} \begin{bmatrix} 1 & a_{m-1} & \dots & a_0 \end{bmatrix} \right) > 0$$

defines an intersection of halfspaces

- o Solving the Quasi-Convex Program

$$\text{minimize}_{a,b,c} \left\| H(e^{j\omega}) - \frac{b(e^{j\omega}) + jc(e^{j\omega})}{a(e^{j\omega})} \right\|$$

quasi-convex set

$$\deg(a) = m, \quad \deg(b) \leq m,$$

$$\deg(c) \leq m,$$

subject to Stability: $a(e^{j\omega}) > 0, \forall \omega \in [0, \pi]$

Passivity: $b(e^{j\omega}) > 0, \forall \omega \in [0, \pi]$

convex set

Standard problem.

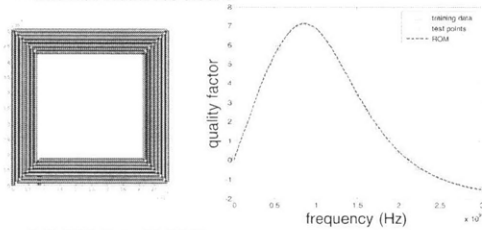
Use for example by the ellipsoid algorithm

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 25.

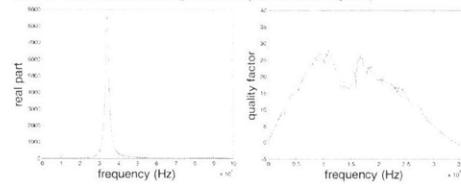
Model-Order Reduction II

- Example 2: RF Inductor with Substrate (from field solver)
 - RF inductor with substrate effect captured by layered Green's function [Hu Dac 05]
 - System matrices are frequency dependent
 - Full model has infinite order
 - Reduced model has order 6



- Example 3: RF Inductor Model (from measurement)

Fabricated 7 turn spiral inductor
 Blue: measurement
 Red: 10^6 order reduced model (positive real part constraint imposed)



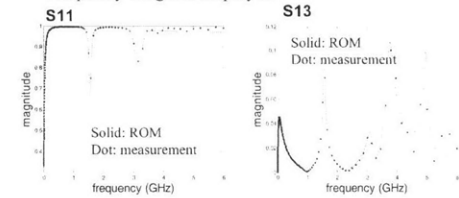
482

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 25.

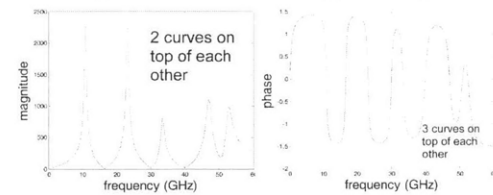
Model-Order Reduction II

- Example 4: Model of graphic card package (from measurement)
 - Industry example of a multi-port device (390 frequency samples)
 - 12^{th} order SISO reduced models are constructed
 - Bounded realness constraint is imposed
 - Frequency weight is employed



- Example 5: Large IC power distribution grid (from field solver)

- Power distribution grid (dimension size = 7mm, wire width = $2 \mu\text{m}$)
- Blue: full model (order 2046)
- Red: QCO 40^{th} order reduced model (positive real)



REDUCTION VIA PROJECTION FRAMEWORK

- Counting Degrees of Freedom

Dynamical System

$$\frac{dx}{dt} = Ax(t) + bu(t)$$

$$y(t) = c^T x(t)$$

Degrees of Freedom

$$n^2 + 2n$$

coefficients

Transfer Function

$$H(s) = \frac{r_1}{(s - \lambda_1)} + \dots + \frac{r_n}{(s - \lambda_n)}$$

$$= k \frac{(s - \zeta_1) \dots (s - \zeta_{n-1})}{(s - \lambda_1) \dots (s - \lambda_n)}$$

$$= \frac{b_0 + b_1 s + \dots + b_{n-1} s^{n-1}}{1 + a_1 s + \dots + a_n s^n}$$

Degrees of Freedom

$$2n$$

coefficients

- Fully Invertible Change of Coordinates

Reduced Model Transfer Function

$$\frac{dx}{dt} = Ax(t) + bu(t) \Rightarrow H(s) = c^T (sI - A)^{-1} b$$

$$y(t) = c^T x(t)$$

Apply any invertible change of coordinates to the state $x(t) = V\tilde{x}(t)$

$$V^{-1} \frac{d\tilde{x}}{dt} = V^{-1} A V \tilde{x}(t) + V^{-1} b u(t) \Rightarrow \tilde{H}(s) = c^T (sI - A)^{-1} b = H(s)$$

$$y(t) = c^T V \tilde{x}(t)$$

Many Dynamical System have the same transfer function!!!

$$(sI - V^{-1} A V) \tilde{X}(s) = V^{-1} b U(s)$$

$$\tilde{X}(s) = (sI - V^{-1} A V)^{-1} V^{-1} b U(s)$$

$$\tilde{X}(s) = (sV^{-1} V - V^{-1} A V)^{-1} V^{-1} b U(s)$$

$$\tilde{X}(s) = V^{-1} (sI - A)^{-1} V V^{-1} b U(s)$$

$$Y(s) = C^T V V^{-1} (sI - A)^{-1} V V^{-1} b U(s)$$

483

- Eigenanalysis

Consider an ODE: $\frac{dx(t)}{dt} = Ax(t) + bu(t), x(0) = 0$

Eigendecomposition:

$$A = \underbrace{\begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ E_1 & E_2 & \dots & E_N \\ \downarrow & \downarrow & \dots & \downarrow \end{bmatrix}}_E \begin{bmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \lambda_N & \\ & & & \ddots \end{bmatrix} \underbrace{\begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ E_1 & E_2 & \dots & E_N \\ \downarrow & \downarrow & \dots & \downarrow \end{bmatrix}}^{-1}$$

Change of Variables: $Ew(t) = x(t) \Leftrightarrow w(t) = E^{-1}x(t)$

Substituting: $\frac{dEw(t)}{dt} = AEw(t) + bu(t), Ew(0) = 0$

Multiply by E^{-1} : $\frac{dw(t)}{dt} = E^{-1}AEw(t) + E^{-1}bu(t)$

- Eigenvalue Method - Reduced Models via Mode Truncation

$$\begin{bmatrix} \tilde{w}_1 \\ \vdots \\ \tilde{w}_q \\ \vdots \\ \tilde{w}_N \end{bmatrix} = \begin{bmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \lambda_q & \\ & & & \ddots \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_q \\ \vdots \\ w_N \end{bmatrix} + \begin{bmatrix} \tilde{b}_1 \\ \vdots \\ \tilde{b}_q \\ \vdots \\ \tilde{b}_N \end{bmatrix} u(t)$$

$$\begin{bmatrix} \tilde{w}_1 \\ \vdots \\ \tilde{w}_q \\ \vdots \\ \tilde{w}_N \end{bmatrix} = \begin{bmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \lambda_q & \\ & & & \ddots \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_q \\ \vdots \\ w_N \end{bmatrix} + \begin{bmatrix} \tilde{b}_1 \\ \vdots \\ \tilde{b}_q \\ \vdots \\ \tilde{b}_N \end{bmatrix} u(t)$$

- Projection Framework: Noninvertible Change of Coordinates

Note: $q \ll N$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{N-1} \\ x_N \end{bmatrix} \approx \underbrace{\begin{bmatrix} \uparrow & \dots & \uparrow \\ \tilde{v}_1 & \dots & \tilde{v}_q \\ \downarrow & \dots & \downarrow \end{bmatrix}}_{V_q} \begin{bmatrix} \hat{x}_1 \\ \vdots \\ \hat{x}_q \end{bmatrix}$$

Note: $q \ll N$

reduced state

original state

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 25.
Model-Order Reduction II

$$\frac{d}{dt} \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} \mathbf{A} \end{bmatrix} \begin{bmatrix} x \end{bmatrix} + \begin{bmatrix} b \end{bmatrix} u(t)$$

$$\frac{d}{dt} \begin{bmatrix} \hat{x} \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{A}} \end{bmatrix} \begin{bmatrix} \hat{x} \end{bmatrix} + \begin{bmatrix} \hat{b} \end{bmatrix} u(t)$$

$$\begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} \bar{v}_1 \end{bmatrix} \begin{bmatrix} \bar{v}_2 \end{bmatrix} \cdots \begin{bmatrix} \bar{v}_q \end{bmatrix} \begin{bmatrix} \hat{x} \end{bmatrix}$$

- Original System

$$\frac{dx}{dt} = Ax(t) + bu(t)$$

$$y(t) = c^T x(t)$$

- Substitute $x \leftarrow V_q \hat{x}(t)$

$$V_q \frac{d\hat{x}}{dt} = AV_q \hat{x}(t) + bu(t)$$

$$\hat{y}(t) = c^T V_q \hat{x}(t)$$

- Notes: now few variables ($q \ll N$) in the state, but still thousands of equations (N)

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 25.
Model-Order Reduction II

$$\hat{\mathbf{A}} = \underbrace{\begin{bmatrix} \mathbf{U}^T & \mathbf{A} & \mathbf{V} \end{bmatrix}}_{q \times q} \begin{bmatrix} \mathbf{A} \\ \mathbf{V} \end{bmatrix}$$

$$\mathbf{U}^T \mathbf{V} = \mathbf{I}$$

$$V_q \frac{d\hat{x}(t)}{dt} = AV_q \hat{x}(t) + bu(t)$$

$$y(t) = c^T V_q \hat{x}(t)$$

Reduction of number of equations:
test by multiplying by V_q^T

$$U_q^T V_q \frac{d\hat{x}(t)}{dt} = U_q^T AV_q \hat{x}(t) + U_q^T b u(t)$$

$$U_q^T V_q = \mathbf{I}$$

If U_q^T and V_q^T are
chosen biorthogonal

$$y(t) = c^T V_q \hat{x}(t)$$

$$\frac{d\hat{x}(t)}{dt} = \hat{\mathbf{A}} \hat{x}(t) + \hat{b} u(t)$$

$$y(t) = \hat{c}^T \hat{x}(t)$$

$$V_q^T = \begin{bmatrix} 1 & 0 & 0 & \vdots & \vdots & \vdots \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ \vdots & \vdots & \vdots & 0 & 0 & 1 \end{bmatrix}$$

eliminate all the equations except for the first three and last three

$$V_q^T = \begin{bmatrix} 1 & 1 & 1 & & & \\ & & & 1 & 1 & 1 & \\ & & & & & & 1 & 1 & 1 \end{bmatrix}$$

Add groups of three

$$\frac{dx}{dt} = \begin{bmatrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{bmatrix} x + bu$$

nxn

$$\begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \frac{d\hat{x}}{dt} = \begin{bmatrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{bmatrix} \hat{x} + bu$$

nxq nxn nxq

485

$$\frac{d\hat{x}}{dt} = \begin{bmatrix} U_q^T & & \\ \text{qxn} & & \end{bmatrix} \begin{bmatrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{bmatrix} \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \hat{x} + U_q^T bu$$

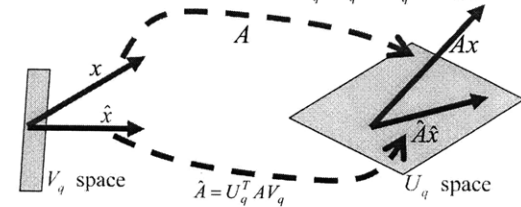
nxn nxq

$$\frac{d\hat{x}}{dt} = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \hat{x}(t) + \hat{b}u(t)$$

qxq

$\frac{dx}{dt} = Ax(t) + bu(t)$
Non-invertible change of coordinates (Projection)
 $x = U_q \hat{x}$

$\frac{d\hat{x}}{dt} = \hat{A} \hat{x}(t) + \hat{b}u(t)$
Equation Testing (Projection)
 $U_q^T A V_q \hat{x} = U_q^T A x = \hat{A} \hat{x}$



INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 25.

Model-Order Reduction II

- ❖ Approaches for picking V and U
 - Use Eigenvectors of the system matrix (modal analysis)

Eigenmode Analysis

$$\mathbf{A} = \mathbf{A}^T$$

$$\mathbf{A} = \mathbf{E}\mathbf{\Lambda}\mathbf{E}^{-1} = \begin{bmatrix} \uparrow & \uparrow & & \uparrow & \uparrow & & \uparrow \\ E_1 & E_2 & \dots & E_N & & & \\ \downarrow & \downarrow & & \downarrow & \downarrow & & \downarrow \end{bmatrix} \begin{bmatrix} \lambda_1 & & & & & & \\ & \lambda_2 & & & & & \\ & & \dots & & & & \\ & & & \lambda_N & & & \end{bmatrix} \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ E_1 & E_2 & \dots & E_N \\ \downarrow & \downarrow & & \downarrow \end{bmatrix}^{-1}$$

$$\mathbf{x}(t) = \mathbf{E}\tilde{\mathbf{x}} = \begin{bmatrix} \uparrow & \uparrow & & \uparrow & \uparrow & & \uparrow \\ E_1 & E_2 & \dots & E_q & E_{q+1} & \dots & E_N \\ \downarrow & \downarrow & & \downarrow & \downarrow & & \downarrow \end{bmatrix} \tilde{\mathbf{x}}$$

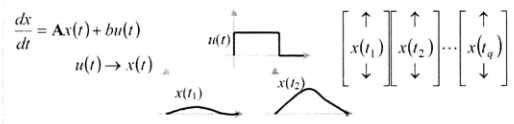
pick some of these

Invertible change of coordinates

$$\mathbf{x} = \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ E_1 & E_2 & \dots & E_q \\ \downarrow & \downarrow & & \downarrow \end{bmatrix} \hat{\mathbf{x}} \quad \text{Non-invertible change of coordinates}$$

V_q q eigenvectors

$$\mathbf{V}^T = ? \quad \mathbf{V}^T \mathbf{V} = \mathbf{I} \quad \mathbf{V}_q = \begin{bmatrix} \uparrow & \uparrow & \uparrow & \dots & \uparrow \\ E_1 & E_2 & E_3 & \dots & E_q \\ \downarrow & \downarrow & \downarrow & & \downarrow \end{bmatrix}$$



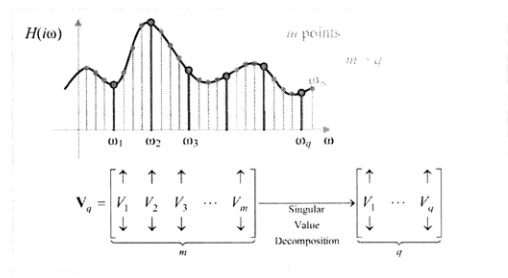
$$\mathbf{V}_q = \begin{bmatrix} \uparrow & \uparrow & \uparrow & \dots & \uparrow \\ x(t_1) & x(t_2) & x(t_3) & \dots & x(t_m) \\ \downarrow & \downarrow & \downarrow & & \downarrow \end{bmatrix} \xrightarrow{\text{Singular Value Decomposition}} \begin{bmatrix} \uparrow & \dots & \uparrow \\ x(t_1) & \dots & x(t_q) \\ \downarrow & & \downarrow \end{bmatrix}$$

486

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 25.

Model-Order Reduction II



- Use Frequency Domain Data
 - Compute $x(s_1), x(s_2), \dots, x(s_k)$
 - Use the SVD to pick $q < k$ important vectors
- Use Time Series Data
 - Compute $x(t_1), x(t_2), \dots, x(t_k)$
 - Use the SVD to pick $q < k$ important vectors

SVD Singular Value Decomposition Method
other names:

- POD Proper Orthogonal Decomposition
- KLD Karhunen-Loeve Decomposition
- PCA Principal Component Analysis
- "Poor Man" TBR (Truncated Balance Realizations)

- Use Singular Vectors of System Grammians Product (Truncated Balance Realizations)
- Use Krylov Subspace Vectors (Moment Matching)

INTRODUCTION TO NUMERICAL SIMULATION

LECTURE 26.

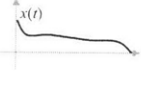
Model-Order Reduction III

TODAY'S OUTLINE:

- ❖ Truncated Balance Realizations
- ❖ Krylov Subspace Moment Matching
 - Preserving passivity
 - Need for orthogonalization (Arnoldi process)
 - Overall computational complexity
- ❖ Comparing Truncated Balance Realizations and Krylov Subspace Moment Matching

TRUNCATED BALANCE REALIZATIONS

- ❖ Approaches for picking V and U
 - Use Eigenvectors of the system matrix (modal analysis)
 - Use Frequency Domain Data
 - Compute $x(s_1), x(s_2), \dots, x(s_k)$
 - Use the SVD to pick $q < k$ important vectors
 - Use Time Series Data
 - Compute $x(t_1), x(t_2), \dots, x(t_k)$
 - Use the SVD to pick $q < k$ important vectors
 - Use Singular Vectors of System Grammians Product (Truncated Balance Realizations)
 - Use Krylov Subspace Vectors (Moment Matching)
- Observability Gramian

$$u(t) = 0 \Rightarrow \begin{cases} \frac{dx}{dt} = Ax(t) & x(0) = x \\ y(t) = C^T x(t) & y(t) = e^{At} Cx \end{cases}$$


$$x^T \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_{N \times N} x$$

Energy of the output $y(t)$ starting from state x with no input:

$$\|y(t)\|_2^2 = \int_0^\infty y(t)^T y(t) dt = \int_0^\infty (C e^{At} x)^T C e^{At} x dt = x^T \underbrace{\int_0^\infty e^{At} C^T C e^{At} dt}_{W_O} x$$

Note: it is also the solution of $A^T W_O + W_O A = -C^T C$

Note: If $x = x_i$ the i^{th} eigenvector of W_O :

$$\|y(t)\|_2^2 = x_i^T W_O x_i = \lambda_{O,i}$$

Hence: eigenvectors of W_O corresponding to small eigenvalues do NOT produce much energy at the output (i.e. they are not very observable):

Idea: let's get rid of them!

$$W_O = \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ x_1 & \dots & x_N & \\ \downarrow & & \downarrow & \end{bmatrix} \begin{bmatrix} \lambda_{O,1} & & & \\ & \ddots & & \\ & & \lambda_{O,N} & \\ & & & \ddots \end{bmatrix} \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ x_1 & \dots & x_N & \\ \downarrow & & \downarrow & \end{bmatrix}^{-1}$$

eigenvectors

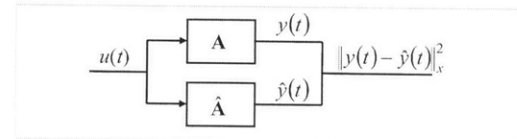
$$\text{Note } W_O^{-1} = W_O^T \Rightarrow \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ x_1 & \dots & x_N & \\ \downarrow & & \downarrow & \end{bmatrix}^{-1} = \begin{bmatrix} \leftarrow & x_1 & \rightarrow \\ & \vdots & \\ \leftarrow & x_N & \rightarrow \end{bmatrix}$$

Unitary Orthogonal Matrix

$$\|y(t)\|_2^2 = x_i^T W_O x_i$$

$$= \begin{bmatrix} \leftarrow & x_i^T & \rightarrow \\ & \vdots & \\ \leftarrow & x_N & \rightarrow \end{bmatrix} \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ x_1 & \dots & x_N & \\ \downarrow & & \downarrow & \end{bmatrix} \begin{bmatrix} \lambda_{O,1} & & & \\ & \ddots & & \\ & & \lambda_{O,N} & \\ & & & \ddots \end{bmatrix} \begin{bmatrix} \leftarrow & x_1 & \rightarrow \\ & \vdots & \\ \leftarrow & x_N & \rightarrow \end{bmatrix} \begin{bmatrix} \uparrow \\ \vdots \\ x_i \\ \downarrow \end{bmatrix}$$

$$= \begin{bmatrix} 0 & \dots & 0 & 1 & 0 & \dots & 0 \\ & & & \lambda_{O,1} & & & \\ & & & & \ddots & & \\ & & & & & \lambda_{O,N} & \\ & & & & & & \ddots \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow i = \lambda_{O,i}$$



➤ Controllability Gramian

Minimum amount input energy required to drive the system to a specific state x :

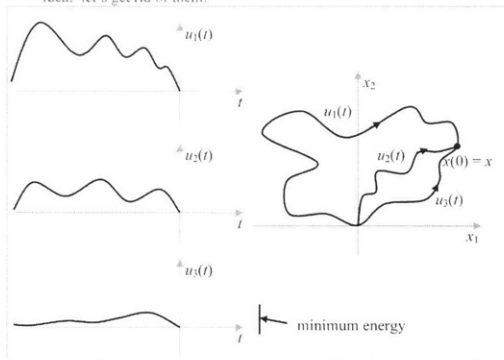
$$\min_{u(t)} \int_{-\infty}^0 u(t)^T u(t) dt = x^T \underbrace{\int_0^{\infty} e^{A^t} B B^T e^{A^t} dt}_{W_C^{-1}} x$$

Note: If $x = x_i$, the i^{th} eigenvector of W_C :
 It is also the solution of $A W_C + W_C A^T = -B B^T$

$$\min_{u(t)} u(t)^2_{x_i} = x_i^T W_C^{-1} x_i = \frac{1}{\lambda_{C,i}}$$

Hence: eigenvectors of W_C corresponding to small eigenvalues do NOT produce much energy at the output (i.e. they are not very controllable);

Idea: let's get rid of them!



eigenvectors

$$\begin{aligned} x_i^T W_C^{-1} x_i &= x_i^T \left[X_C \Lambda_C^{-1} X_C^T \right]^{-1} x_i \\ &= x_i^T X_C \Lambda_C X_C^T x_i \\ &= x_i^T \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ x_1 & \dots & x_N & \\ \downarrow & \downarrow & & \downarrow \end{bmatrix} \begin{bmatrix} \frac{1}{\lambda_{C,1}} & & & \\ & \ddots & & \\ & & \frac{1}{\lambda_{C,N}} & \\ & & & \ddots \end{bmatrix} \begin{bmatrix} \leftarrow x_1 \rightarrow \\ \vdots \\ \leftarrow x_N \rightarrow \end{bmatrix} \\ &= \begin{bmatrix} 0 & \dots & 0 & 1 & 0 & \dots & 0 \\ & & & \frac{1}{\lambda_{C,i}} & & & \\ & & & & \ddots & & \\ & & & & & \frac{1}{\lambda_{C,N}} & \\ & & & & & & 0 \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow i = \frac{1}{\lambda_{C,i}} \end{aligned}$$

- Native Controllability/Observability MOR
 - Suppose I could compute a basis for the strongly observable and/or strongly controllable spaces. Projection-based MOR can give a reduced model that deletes weakly observable and/or weakly controllable modes.
 - Problem:
 - What if the same mode is strongly controllable, but weakly observable?
 - Are the eigenvalues of the respective Gramians even unique?
- Changing Coordinate System
 - Consider an invertible change of coordinates: $x(t) = U \tilde{x}(t)$
 - We know that the input/output relationship will be unchanged.
 - But what about the Gramians, and their eigenvalues?

$$\tilde{W}_O \leftarrow U^T W_O U \quad \tilde{W}_C \leftarrow U^{-1} W_C U^{-T}$$
 - Gramians and their eigenvalues change! Hence the relative degrees of observability and controllability are properties of the coordinate system
 - A bad choice of coordinates will lead to bad reduced models if we look at controllability and observability separately.
 - What coordinate system should we use then?

$$W_o = \int_0^{\infty} e^{A^T t} C^T C e^{A t} dt \quad x = U\bar{x}$$

$$\begin{cases} \frac{d\bar{x}}{dt} = U^{-1} A U \bar{x}(t) + U^{-1} b u(t) \\ y(t) = C U \bar{x}(t) \end{cases}$$

$$\begin{aligned} \bar{W}_o &= \int_0^{\infty} e^{(U^{-1} A U)^T t} (C U)^T (C U) e^{U^{-1} A U t} dt \\ &= \int_0^{\infty} U^T e^{A^T t} U^{-T} U^T C^T C U U^{-1} e^{A t} U dt \\ &= U^T W_o U \end{aligned}$$

$$W_o = X \Lambda X^T \quad U \text{ invertible}$$

$$\bar{W}_o = \underbrace{U^T X \Lambda X^T U}_{\text{these are \textit{not} the new eigenvectors in general}} \quad X^T U \neq (U^T X)^{-1}$$

$$U^T X X^T U = U^T U = I \quad = I \text{ only if } U \text{ is orthonormal}$$

➤ **Balancing**
Fortunately the eigenvalues of the product of the Gramians (Hankel singular values) do not change when changing coordinates:

Diagonal matrix with eigenvalues of the product

$$W_c W_o = S \tilde{\Sigma}^2 S^{-1}$$

The eigenvectors change
But not the eigenvalues

$$U^{-1} W_c U^{-T} \quad U^T W_o U = U^{-1} W_c W_o U = (U^{-1} S) \tilde{\Sigma}^2 (U^{-1} S)^{-1}$$

And since W_c and W_o are symmetric a change of coordinate matrix U can be found that diagonalize both

$$\Sigma \Sigma = \tilde{\Sigma}^2 \quad \text{In Balanced Coordinates the Gramians are equal and diagonal}$$

➤ **Selection of Vectors for the Columns of the Reduced Order Projection Matrix.**
In balanced coordinates it is easy to select the best vectors for the reduced model: we want the subspace of vectors that are at the same time most controllable and observable:

$$U^{-1} W_c U^{-T} \quad U^T W_o U = \Sigma^2$$

Simply pick the eigenvectors corresponding to the largest entries on the diagonal (Hankel singular values)

In other words the ones corresponding to the largest eigenvalues of the controllability and observability Gramians product.

➤ **Truncated Balance Realization Summary**

- o The good news: $\|H(j\omega) - H_q(j\omega)\|_{\infty} \leq \sum_{q+1, q+1}^{\infty} (\Sigma_{q+1, q+1} + \dots + \Sigma_{N, N})$

- We even have bounds for the error
- Can do even a bit better with the optimal Hankel Reduction
- o The bad news: $A^T W_o + W_o A = -C^T C$
 - It is expensive:
 - Need to compute the Gramians (solve Lyapunov equation)
 - Need to compute eigenvalues of the product: $O(N^3)$
- o The bottom line:
 - If the size of your system allows $O(N^3)$ computation, Truncation Balance Realization or Hankel Reduction are a much better choice than the any other reduction method.
 - But if you cannot afford $O(N^3)$ computation (e.g. dense matrix with $N > 5000$) then PRIMA or PVL or Quasi-Convex-Optimization are better choices.

- ❖ **Approaches for picking V and U**
 - o Use Eigenvectors of the system matrix (modal analysis)
 - o Use Frequency Domain Data
 - Compute $x(s_1), x(s_2), \dots, x(s_k)$
 - Use the SVD to pick $q < k$ important vectors
 - o Use Time Series Data
 - Compute $x(t_1), x(t_2), \dots, x(t_k)$
 - Use the SVD to pick $q < k$ important vectors
 - o Use Singular Vectors of System Grammians Product (Truncated Balance Realizations)

o Use Krylov Subspace Vectors (Moment Matching)

Original System Transfer Function Moments

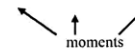
$$H(s) = c^T (sI - A)^{-1} b$$

$$= -c^T \underbrace{(I - sA^{-1})^{-1}}_{\text{Taylor expand with respect to } s} A^{-1} b$$

$$\frac{1}{1-x} = 1 + x + x^2 + x^3 + \dots$$

$$\left(I - \underbrace{sA^{-1}}_{\text{small}} \right)^{-1} = I + sA^{-1} + s^2 A^{-2} + s^3 A^{-3} + \dots$$

$$= \sum_{k=0}^{\infty} c^T A^{-(k+1)} b s^k$$



$$H_r(s) = \frac{b'_0 + b'_1 s + \dots + b'_{q-1} s^{q-1}}{1 + a'_1 s + \dots + a'_q s^q} = m_0 + m_1 s + \dots + m_{2q-1} s + \dots$$

$$c^T \underline{\Lambda}^{-1} b \quad c^T \underline{\Lambda}^{-2} b \quad c^T \underline{\Lambda}^{-3} b \quad \dots \quad c^T \underline{\Lambda}^{-k} b$$

$$\underline{\Lambda}^{-1} b = \underline{\Lambda}^{-1} [\alpha_1 \underline{e}_1 + \dots + \alpha_N \underline{e}_N]$$

$$= \alpha_1 \lambda_1^{-1} \underline{e}_1 + \dots + \alpha_N \lambda_N^{-1} \underline{e}_N$$

$$\underline{\Lambda}^{-k} b = \alpha_1 \lambda_1^{-k} \underline{e}_1 + \dots + \alpha_N \lambda_N^{-k} \underline{e}_N$$

dominates

Cross-Multiplying and Matching Terms

$$\begin{bmatrix} m_0 & m_1 & \dots & m_{k-1} \\ m_1 & \ddots & & \vdots \\ \vdots & & \ddots & m_{2q-3} \\ m_{k-1} & \dots & m_{2q-3} & m_{2q-2} \end{bmatrix} \begin{bmatrix} a_q \\ a_{q-1} \\ \vdots \\ a_1 \end{bmatrix} = \begin{bmatrix} m_q \\ m_{q+1} \\ \vdots \\ m_{2q-1} \end{bmatrix}$$

A Canonical Form for Model Order Reduction

$$s x = A x + \bar{b} u$$

$$y = c^T x$$



$$E = A^{-1}$$

$$b = A^{-1} \bar{b}$$

Assuming A is non-singular we can cast the dynamical linear system into a canonical form for moment matching model order reduction

$$s E x = x + b u$$

Note: this step is not necessary, it just makes the notation simple for educational purposes

$$y = c^T x$$

The Moment Matching Idea [Grimme PhD97]

$$s E x = x + b u \implies x = -(I - s E)^{-1} b u$$

Taylor series expansion:

$$x = -\sum_{k=0}^{\infty} s^k E^k b u \implies x \in \text{span} \{b, E b, E^2 b, \dots\}$$

change basis: use only first q vectors of the Taylor series expansion matching first q derivatives around expansion point

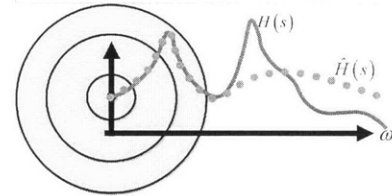
$$x = b + s \underline{E} b + s^2 \underline{E}^2 b + s^3 \underline{E}^3 b + \dots$$

Moment

$$\frac{\partial H}{\partial s} \quad \frac{1}{2} \frac{\partial^2 H}{\partial s^2} \quad \frac{1}{3!} \frac{\partial^3 H}{\partial s^3}$$

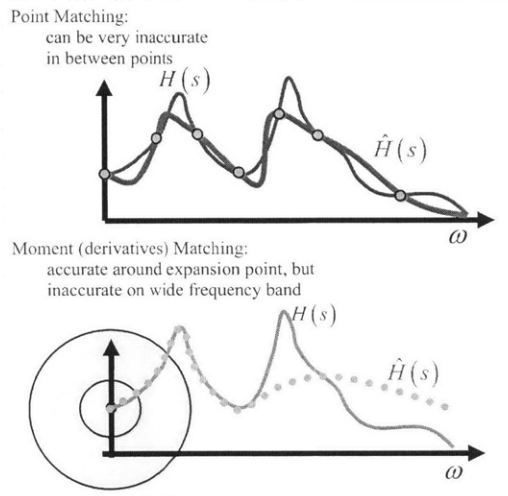
$$x \in \underbrace{\{b, \underline{E} b, \underline{E}^2 b, \dots\}}_N$$

$$x = \begin{bmatrix} | & | & | & | \\ b & \underline{E} b & \underline{E}^2 b & \dots \end{bmatrix} \hat{x}$$

$$x = \begin{bmatrix} | & | & | & | \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} \hat{x}$$


$$x = \begin{bmatrix} | & | & | & | \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} U_q \hat{x}$$

$$x = \begin{bmatrix} | & | & | & | \\ b & \underline{E} b & \underline{E}^2 b & \dots \end{bmatrix} \hat{x}$$



- Krylov Subspaces – Definition
The order k Krylov subspace generated from matrix E and vector b is defined as $\kappa_k(E, b) \equiv \text{span}\{b, Eb, E^2b, \dots, E^{k-1}b\}$
- Moment Matching around Non-zero Frequencies
 - Instead of expanding around only $s = 0$ we can expand around other points $s = 0, s = s_1, s = s_2, \dots, s = s_j$

$$\begin{array}{l}
 sx = Ax + bu \\
 y = c^T x
 \end{array}
 \quad
 \xrightarrow{\quad}
 \quad
 \begin{array}{l}
 (\tilde{s} + s_h)x = Ax + bu \\
 y = c^T x
 \end{array}$$

$$\begin{aligned}
 \tilde{s}x &= (Ax - s_h x) + bu \\
 \tilde{s} \underbrace{(A - s_h I)^{-1}}_{E_h} x &= x + \underbrace{(A - s_h I)^{-1} b}_b u
 \end{aligned}$$

$$x = \begin{bmatrix} b \\ Eb \\ E^2b \\ \vdots \\ b_h \\ E_h b_h \\ \vdots \\ E_k b_k \\ E_k^2 b_k \\ \vdots \end{bmatrix} \tilde{s}$$

- For each expansion point the problem can then be put again in the canonical form

$$\begin{array}{l}
 E_h = (A - s_h I)^{-1} \\
 b_h = (A - s_h I)^{-1} b
 \end{array}$$

$$\begin{array}{l}
 \tilde{s} E_h x = x + b_h u \\
 y = c^T x
 \end{array}$$

- Projection Framework: Moment Matching Theorem [E. Grimme 97]

$$\text{If } \text{Range}(U_q) \supseteq \bigcup_{h=1}^J \kappa_{k_h^u}(E_h, b_h)$$

$$\text{and } \text{Range}(V_q) \supseteq \bigcup_{h=1}^J \kappa_{k_h^v}(E_h^T, c_h)$$

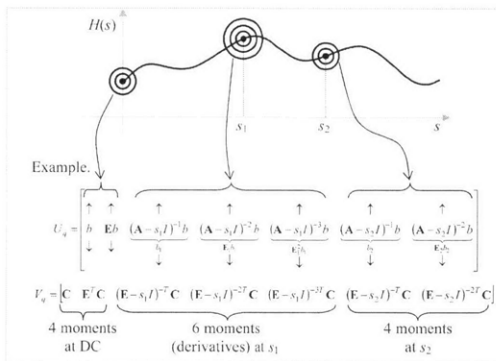
$$\text{Then } \left. \frac{\partial^l H}{\partial s^l} \right|_{s_h} = \left. \frac{\partial^l \hat{H}}{\partial s^l} \right|_{s_h} \text{ for } l = 0, \dots, k_h^u + k_h^v - 1, h = 1, \dots, J$$

Total of $2q$ moment of the transfer function will match

If $\text{span}\{\bar{u}_1, \dots, \bar{u}_q\} \supseteq \bigcup_{j=1}^J \kappa_{k_j^a} \left\{ (\mathbf{A} - s_j I)^{-1}, (\mathbf{A} - s_j I)^{-1} b \right\}$
 and $\text{span}\{\bar{v}_1, \dots, \bar{v}_q\} \supseteq \bigcup_{j=1}^J \kappa_{k_j^c} \left\{ (\mathbf{A} - s_j I)^{-T}, (\mathbf{A} - s_j I)^{-T} c \right\}$
 Then $\frac{\partial^l H(s_j)}{\partial s^l} = \frac{\partial^l H_r(s_j)}{\partial s^l}$ for $l = 0, \dots, k_j^b + k_j^c - 1$

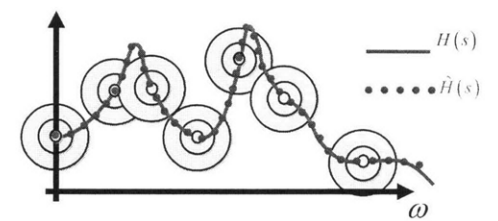
$$\text{span}\{\bar{v}_1, \dots, \bar{v}_q\} = \text{span}\{\bar{u}_1, \dots, \bar{u}_q\} \supseteq \mathbf{A}^{-1} b, \mathbf{A}^{-2} b, \mathbf{A}^{-3} b, \dots$$

$$H_r(s) = \underbrace{m_0 + m_1 s + m_2 s^2 + \dots}_{\text{original system moments}}$$



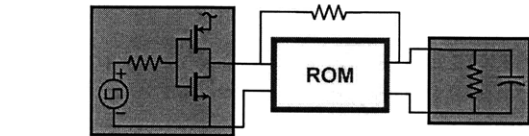
If \mathbf{U} and \mathbf{V} are such that
 $\mathbf{U} = \mathbf{V} = \{\bar{u}_1, \dots, \bar{u}_q\}$ and $\mathbf{U}^T \mathbf{U} = \mathbf{I}$
 $\text{span}\{\bar{u}_1, \dots, \bar{u}_q\} = \text{span}\{\mathbf{A}^{-1} b, \mathbf{A}^{-2} b, \dots, \mathbf{A}^{-q} b\}$
 Then the first q moments of reduced system match
 $H(s) = -c^T (I - s\mathbf{A}^{-1})^{-1} \mathbf{A}^{-1} b = \sum_{k=0}^{\infty} c^T \mathbf{A}^{-(k+1)} b s^k$
 $H_r(s) = -c_r^T (I - s\mathbf{A}_r^{-1})^{-1} \mathbf{A}_r^{-1} b_r = \sum_{k=0}^{\infty} c_r^T \mathbf{A}_r^{-(k+1)} b_r s^k$
 $c^T \mathbf{A}^{-(k+1)} b = c^T \mathbf{U}_q (\mathbf{U}_q^T \mathbf{A} \mathbf{U}_q)^{-(k+1)} \mathbf{U}_q^T b \quad k = \{0, \dots, q-1\}$
 $\mathbf{A}_r = \mathbf{U}_q^T \mathbf{A} \mathbf{U}_q$
 $b_r = \mathbf{U}_q^T b$
 $c_r = \mathbf{U}_q^T c$

- Combine Point and Moment Matching: Multipoint Moment Matching
- Multipole expansion points give larger band
- Moment (derivates) matching gives more accurate behavior in between expansion points

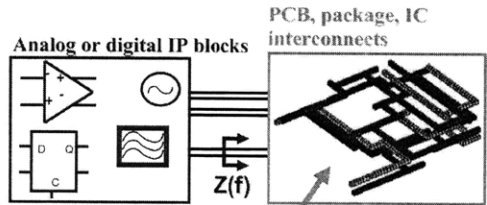


KRYLOV SUBSPACE MOMENT MATCHING

- Preserving Passivity
 - Interconnected Systems
 - In reality, reduced models are only useful when connected together with models of other components in a composite simulation
 - Consider a state-space model connected to external circuitry (possibly with feedback)

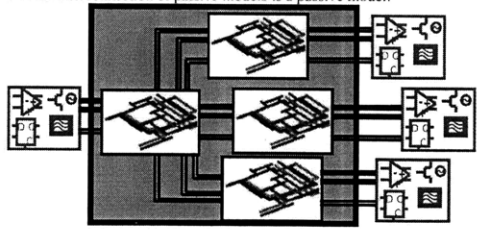


- Can we assure that the simulation of the composite system will be well-behaved? At least preclude non-physical behavior of the reduced model?
- o Need to Preserve Passivity for Models of Passive Interconnect



Note: passive!
Designers will connect models use them in ODE time domain simulators. If the models are not passive they can generate energy and the simulation may explode!!

- o Interconnecting Passive Systems
 - The interconnection of stable models is not necessarily stable
 - But the interconnection of passive models is a passive model:



- o Sufficient Conditions for Passivity

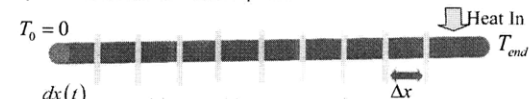
$$sEx = x + Bu$$

$$y = Cx$$

- Sufficient conditions for passivity
1. $C = B^T$
 2. $x^T E x \leq 0$, for all x
i.e. E is negative semidefinite

Note that these are NOT necessary conditions (common misconceptions)

- o Example Finite Difference on Poisson Equation



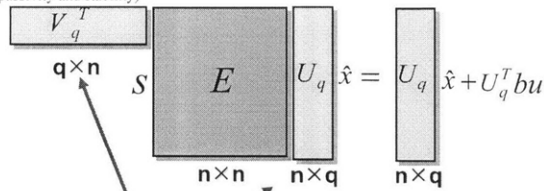
$$\frac{dx(t)}{dt} = \underbrace{A}_{N \times N} x(t) + \underbrace{b}_{N \times 1} u(t)$$

$$y(t) = \underbrace{c}_{1 \times N} x(t)$$

$$A = \begin{bmatrix} 2 & -1 & & & & & & & \\ -1 & 2 & -1 & & & & & & \\ & -1 & 2 & \ddots & & & & & \\ & & \ddots & \ddots & \ddots & & & & \\ & & & & & -1 & 2 & & \\ & & & & & -1 & 2 & & \end{bmatrix} \quad b = \begin{bmatrix} 0 \\ \vdots \\ \vdots \\ 0 \\ 1 \end{bmatrix} \quad c = \begin{bmatrix} 0 \\ \vdots \\ \vdots \\ 0 \\ 1 \end{bmatrix}$$

We already know the Finite Difference matrices is positive semidefinite. Hence $E = -A^{-1}$ is negative semidefinite.

- Congruence Transformation Preserves Negative Definiteness of E (hence passivity and stability)



If we use $V_q^T = U_q^T$

- Then we loose half of the degrees of freedom i.e. we match only q moments instead of 2q
- But if the original matrix E is negative semidefinite so is the reduced hence the system is passive and stable

$$x^T E x \leq 0 \quad \forall x$$

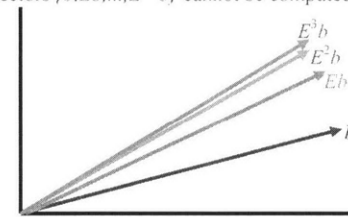
$$x^T U_q^T E U_q x \leq 0 \quad \forall x$$

Pick x_0

$$\underbrace{x_0^T U_q^T E U_q}_{\tilde{x}} x_0 \leq 0 \quad \text{use} \quad x^T E x \leq 0$$

- ❖ Need for Orthogonalization (Arnoldi process)
- Need for Orthonormalization of U

Vectors $\{b, Eb, \dots, E^{k-1}b\}$ cannot be computed directly



Vectors will quickly line up with dominant eigenspace!

$$E = V \Lambda U^T \quad U = \{b, Eb, E^2b, E^3b, \dots\}$$

$$b = \alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_N v_N \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N$$

$$Eb = \alpha_1 E v_1 + \alpha_2 E v_2 + \dots + \alpha_N E v_N = \alpha_1 \lambda_1 v_1 + \alpha_2 \lambda_2 v_2 + \dots + \alpha_N \lambda_N v_N$$

$$E^2 b = \alpha_1 \lambda_1^2 v_1 + \alpha_2 \lambda_2^2 v_2 + \dots + \alpha_N \lambda_N^2 v_N$$

$$E^{k-1} b = \alpha_1 \lambda_1^{k-1} v_1 + \alpha_2 \lambda_2^{k-1} v_2 + \dots + \alpha_N \lambda_N^{k-1} v_N \cong \alpha_1 \lambda_1^{k-1} v_1$$

very large

They all become linearly dependent and parallel to the eigenvector with the largest eigenvalue

- ❖ Overall Computational Complexity
 - Orthonormalization of U: The Arnoldi Algorithm Computational Complexity
 - $\tilde{u}_1 = b / \|b\|$ Normalize first vector $O(n)$
 - For $i = 1$ to k
 - $\tilde{u}_{i+1} = E\tilde{u}_i$ Generates new Krylov subspace vector $O(k^2n)$
 - For $j = 1$ to i
 - $\tilde{u}_{i+1} \leftarrow \tilde{u}_{i+1} - \left(\frac{\tilde{u}_{i+1}^T \tilde{u}_j}{\Gamma_{ji}} \right) \tilde{u}_j$ Orthogonalize new vector $O(k^2n)$
 - $\tilde{u}_{i+1} \leftarrow \frac{1}{\|\tilde{u}_{i+1}\|} \tilde{u}_{i+1}$ Normalize new vector $O(n)$
 - Generating Vectors for the Krylov Subspace
 - Most of the computation cost is spent in calculating:
 - $\tilde{u}_{i+1} = E_h \tilde{u}_i = (A - s_h I)^{-1} \tilde{u}_i$
 - $\tilde{u}_{i+1} = (A - s_h I)^{-1} \tilde{u}_i \Rightarrow (A - s_h I) \tilde{u}_{i+1} = \tilde{u}_i \quad O(n)$
 - Set up and solve a linear system using GCR
 - $(A - s_h I)^{-1} \tilde{u}_{i+1} = \tilde{u}_i$
 - If we have a good preconditioner and a fast matrix vector product each new vector is calculated $O(n)$
 - The total complexity for calculating the projection matrix U_q is $O(qn)$

COMPARING TRUNCATED BALANCE REALIZATIONS AND KRYLOV SUBSPACE MOMENT MATCHING

$$U_q^T E U_q \frac{d\hat{x}}{dt} = \hat{x}(t) + U^T b \hat{x}(t) \quad \text{Reduced Model}$$

$$\hat{y}(t) = C^T U_q \hat{x}(t)$$

- ❖ What about computing the reduced matrix $\hat{E} = U_q^T E U_q$?

$$E \begin{bmatrix} \tilde{u}_i \end{bmatrix} = \begin{bmatrix} \tilde{u}_1 & \dots & \tilde{u}_q \end{bmatrix} \begin{bmatrix} \Gamma_{1,i} \\ \vdots \\ \Gamma_{q,i} \end{bmatrix} \quad \text{Orthonormalization of the } i^{\text{th}} \text{ column of } U_q$$

$$E \begin{bmatrix} \tilde{u}_1 & \dots & \tilde{u}_q \end{bmatrix} = \begin{bmatrix} \tilde{u}_1 & \dots & \tilde{u}_q \end{bmatrix} \Gamma \quad \begin{matrix} \text{Orthonormalization of} \\ \text{all columns of } U_q \end{matrix}$$

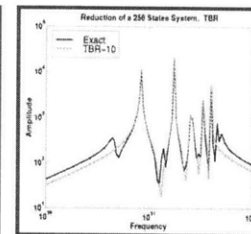
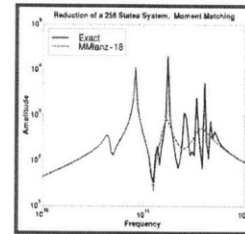
$$EU_q = U_q \Gamma$$

$$U_q^T E U_q = \Gamma$$

So we don't need to compute the reduced matrix. We have it already: $\hat{E} = \Gamma$

- ❖ Two Complementary Approaches

<p>Moment Matching Approaches</p> <ul style="list-style-type: none"> - Accurate over a narrow band. <ul style="list-style-type: none"> • Matching function values and derivatives. - Cheap: $O(qn)$ - Use it as a FIRST STAGE REDUCTION 	<p>Truncated Balanced Realization and Hankel Reduction</p> <ul style="list-style-type: none"> - Optimal (best accuracy for given size q, and a priori error bound. - Expensive: $O(n^3)$ - USE IT AS A SECOND STAGE REDUCTION
--	---



❖ Model Order Reduction Computational Complexity (time and memory)

$$\frac{dx}{dt} = \mathbf{A} x(t) + \mathbf{B} u(t)$$

A is DENSE!

MOR technique	Computational Complexity
TBR (*81) Hankel (*84) (have error bounds)	Bottleneck: SVD $O(N^3)$ e.g. 10months, 80GB, for $N=100,000$
Moment matching (*97) (no error bounds)	Bottleneck: matrix-vector product $O(qN^2)$ e.g. 7days, 80GB, for $N=100,000$ $q=10$
Moment Matching + pFFT matrix-vector	$O(qN \log N)$ e.g. 8hours, 0.3GB for $N=100,000$ $q=10$
QuasiConvex Optimization + pFFT field solver	$O(mN \log N)$ Same as above, but have error bound!

❖ And Remember, No Matter What...

NEVER INVERT A MATRIX!