

MIT Open Access Articles

Approximate policy iteration: A survey and some new methods

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

Citation: Bertsekas, Dimitri P. "Approximate Policy Iteration: a Survey and Some New Methods." Journal of Control Theory and Applications 9.3 (2011): 310–335. Web.

As Published: <http://dx.doi.org/10.1007/s11768-011-1005-3>

Publisher: Springer-Verlag

Persistent URL: <http://hdl.handle.net/1721.1/73485>

Version: Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

Terms of use: Creative Commons Attribution-Noncommercial-Share Alike 3.0



Approximate Policy Iteration: A Survey and Some New Methods

Dimitri P. Bertsekas †

Abstract

We consider the classical policy iteration method of dynamic programming (DP), where approximations and simulation are used to deal with the curse of dimensionality. We survey a number of issues: convergence and rate of convergence of approximate policy evaluation methods, singularity and susceptibility to simulation noise of policy evaluation, exploration issues, constrained and enhanced policy iteration, policy oscillation and chattering, and optimistic and distributed policy iteration.

Our discussion of policy evaluation is couched in general terms, and aims to unify the available methods in the light of recent research developments, and to compare the two main policy evaluation approaches: projected equations and temporal differences (TD), and aggregation. In the context of these approaches, we survey two different types of simulation-based algorithms: matrix inversion methods such as LSTD, and iterative methods such as LSPE and $TD(\lambda)$, and their scaled variants. We discuss a recent method, based on regression and regularization, which rectifies the unreliability of LSTD for nearly singular projected Bellman equations. An iterative version of this method belongs to the LSPE class of methods, and provides the connecting link between LSTD and LSPE.

Our discussion of policy improvement focuses on the role of policy oscillation and its effect on performance guarantees. We illustrate that policy evaluation when done by the projected equation/TD approach may lead to policy oscillation, but when done by aggregation it does not. This implies better error bounds and more regular performance for aggregation, at the expense of some loss of generality in cost function representation capability. Hard aggregation provides the connecting link between projected equation/TD-based and aggregation-based policy evaluation, and is characterized by favorable error bounds.

1. INTRODUCTION

In this paper we aim to survey and place in broad context a number of issues relating to approximate policy iteration methods for finite-state, discrete-time, stochastic dynamic programming (DP) problems.

† The author is with the Dept. of Electr. Engineering and Comp. Science, M.I.T., Cambridge, Mass., 02139. His research was supported by NSF Grant ECCS-0801549, by the LANL Information Science and Technology Institute, and by the Air Force Grant FA9550-10-1-0412. Many thanks are due to Huizhen (Janey) Yu for extensive helpful discussions and suggestions.

These methods are one of the major approaches for approximate DP, a field that has attracted substantial research interest, and has a wide range of applications, because of its potential to address large and complex problems that may not be treatable in other ways. Among recent works in the extensive related literature, we mention textbooks and research monographs: Bertsekas and Tsitsiklis [BeT96], Sutton and Barto [SuB98], Gosavi [Gos03], Cao [Cao07], Chang, Fu, Hu, and Marcus [CFH07], Meyn [Mey07], Powell [Pow07], Borkar [Bor08], Busoniu, Babuska, De Schutter, and Ernst [BBD10], and the author's text in preparation [Ber10a]; edited volumes and special issues: White and Sofge [WhS92], Si, Barto, Powell, and Wunsch [SBP04], Lewis, Lendaris, and Liu [LLL08], and the 2007-2009 Proceedings of the IEEE Symposium on Approximate Dynamic Programming and Reinforcement Learning; and surveys: Barto, Bradtke, and Singh [BBS95], Borkar [Bor09], Lewis and Vrabie [LeV09], and Szepesvari [Sze09].

For an overview of policy iteration methods, let us focus on the α -discounted n -state Markovian Decision Problem (MDP) with states $1, \dots, n$, controls $u \in U(i)$ at state i , transition probabilities $p_{ij}(u)$, and cost $g(i, u, j)$ for transition from i to j under control u . A (stationary) policy μ is a function from states i to admissible controls $u \in U(i)$, and $J_\mu(i)$ is the cost starting from state i and using policy μ . It is well-known (see e.g., Bertsekas [Ber07] or Puterman [Put94]) that the costs $J_\mu(i)$, $i = 1, \dots, n$, are the unique solution of Bellman's equation

$$J_\mu(i) = \sum_{j=1}^n p_{ij}(\mu(i)) (g(i, \mu(i), j) + \alpha J_\mu(j)), \quad i = 1, \dots, n.$$

Equivalently, the vector $J_\mu \in \mathfrak{R}^n$, which has components $J_\mu(i)$,[†] is the unique fixed point of the mapping $T_\mu : \mathfrak{R}^n \mapsto \mathfrak{R}^n$, which maps $J \in \mathfrak{R}^n$ to the vector $T_\mu J \in \mathfrak{R}^n$ that has components

$$(T_\mu J)(i) = \sum_{j=1}^n p_{ij}(\mu(i)) (g(i, \mu(i), j) + \alpha J(j)), \quad i = 1, \dots, n. \quad (1.1)$$

Similarly, the optimal costs starting from $i = 1, \dots, n$, are denoted $J^*(i)$, and are the unique solution of Bellman's equation

$$J^*(i) = \min_{u \in U(i)} \sum_{j=1}^n p_{ij}(u) (g(i, u, j) + \alpha J^*(j)), \quad i = 1, \dots, n.$$

Equivalently, the optimal cost vector $J^* \in \mathfrak{R}^n$, which has components $J^*(i)$, is the unique fixed point of the mapping $T : \mathfrak{R}^n \mapsto \mathfrak{R}^n$ defined by

$$(TJ)(i) = \min_{u \in U(i)} \sum_{j=1}^n p_{ij}(u) (g(i, u, j) + \alpha J(j)), \quad i = 1, \dots, n. \quad (1.2)$$

An important property is that T is a sup-norm contraction, so the iteration $J^{k+1} = TJ^k$ converges to J^* from any starting point J^0 - this is known as *value iteration*.

Policy iteration is a major alternative to value iteration. It produces a sequence of policies and associated cost functions through iterations that have two phases: *policy evaluation* (where the cost function of a policy is evaluated), and *policy improvement* (where a new policy is generated). In the exact form of the algorithm, the current policy μ is improved by finding $\bar{\mu}$ that satisfies $T_{\bar{\mu}} J_\mu = TJ_\mu$ (i.e., by minimizing in

[†] In our notation, \mathfrak{R}^n is the n -dimensional Euclidean space, all vectors in \mathfrak{R}^n are viewed as column vectors, and a prime denotes transposition. The identity matrix is denoted by I .

the right-hand side of Eq. (1.2) with J_μ in place of J). The improved policy $\bar{\mu}$ is evaluated by solving the linear system of equations $J_{\bar{\mu}} = T_{\bar{\mu}}J_{\bar{\mu}}$, and $(J_{\bar{\mu}}, \bar{\mu})$ becomes the new cost vector-policy pair, which is used to start a new iteration.

In a variant of the method, the improved policy $\bar{\mu}$ is evaluated by applying $T_{\bar{\mu}}$ a finite number of times to an approximate evaluation of the preceding policy μ , which we denote by \tilde{J}_μ . This is known as “optimistic” or “modified” policy iteration, and its motivation is that in problems with a large number of states, the linear system $J_{\bar{\mu}} = T_{\bar{\mu}}J_{\bar{\mu}}$ cannot be practically solved directly by matrix inversion, so it is best solved iteratively by value iteration, that is, by repeated application of $T_{\bar{\mu}}$ to some initial vector (most conveniently the vector \tilde{J}_μ). If the number m of applications of $T_{\bar{\mu}}$ is very large, the exact form of policy iteration is essentially obtained, but practice has shown that it is most efficient to use a moderate value of m . In this case, the algorithm looks like a hybrid of value and policy iteration, involving a sequence of alternate applications of T and T_μ , with μ changing over time.

It is well-known that the exact form of policy iteration converges to an optimal cost-policy pair for any initial conditions (in a finite number of iterations when the number of policies is finite, as in the discounted MDP). The convergence of optimistic policy iteration is more fragile, and requires that the initial pair $(\tilde{J}_{\mu^0}, \mu^0)$ satisfy $T_{\mu^0}\tilde{J}_{\mu^0} \leq \tilde{J}_{\mu^0}$, particularly when implemented in totally asynchronous form (irregularly, one state at a time); see Williams and Baird [WiB93], [BeT96], Section 2.2.3, or [Ber07], Section 1.3.3.† Since simulation-based policy iteration algorithms are typically asynchronous, this restriction has motivated algorithmic modifications with guaranteed asynchronous convergence from arbitrary initial conditions, which will be partly discussed in Section 3.8.

Policy Iteration Methods with Cost Function Approximation

In policy iteration methods with cost function approximation, we evaluate μ by approximating J_μ with a vector Φr_μ from the subspace S spanned by the columns of an $n \times s$ matrix Φ , which may be viewed as basis functions:‡

$$S = \{\Phi r \mid r \in \mathbb{R}^s\}.$$

We generate an “improved” policy $\bar{\mu}$ using the formula $T_{\bar{\mu}}(\Phi r_\mu) = T(\Phi r_\mu)$, i.e.,

$$\bar{\mu}(i) \in \arg \min_{u \in U(i)} \sum_{j=1}^n p_{ij}(u) (g(i, u, j) + \alpha \phi(j)' r_\mu), \quad i = 1, \dots, n, \quad (1.3)$$

† Convergence of optimistic policy iteration with some form of synchronous or regular order of state updates (including Gauss-Seidel variants) is obtained for a variety of DP models involving contraction mappings, without the condition $T_{\mu^0}\tilde{J}_{\mu^0} \leq \tilde{J}_{\mu^0}$, as first shown by Rothblum [Rot79]. However, this is not so under totally asynchronous implementations.

‡ We focus on linear/subspace cost function approximation in this paper. Naturally, the choice of the matrix Φ is critical for effective policy evaluation. The i th row of Φ consists of s numbers, which may be viewed as *features* that characterize state i , and are weighted with the corresponding s components of r_μ to provide an approximation of the cost starting from i under the given policy. Features, when well-crafted, can capture the dominant nonlinearities of the cost vector, and their linear combination can work well as an approximation architecture. In this paper, we do not discuss the selection of Φ , but we note the possibility of its optimal choice within some restricted class by using gradient and random search algorithms (see Menache, Mannor, and Shimkin [MMS06], Yu and Bertsekas [YuB09], Busoniu et. al. [BES10], and Di Castro and Mannor [DiM10] for recent work on this subject). Note that some of the methods discussed in this paper can be adapted so that they can use “sampled” features, that is, features whose values are corrupted by stochastic noise (see Bertsekas and Yu [BeY09], Section 6).

where $\phi(j)'$ is the row of Φ that corresponds to state j [the method terminates with μ if $T_\mu(\Phi r_\mu) = T(\Phi r_\mu)$]. We then repeat with μ replaced by $\bar{\mu}$.

In this paper, for generality and notational simplicity, we consider MDP where for each policy μ , the mapping $T_\mu : \mathfrak{R}^n \mapsto \mathfrak{R}^n$ has the form

$$T_\mu J = g_\mu + A_\mu J, \quad (1.4)$$

for given $n \times n$ matrix A_μ with nonnegative components $a_{ij}(\mu(i))$, and vector $g_\mu \in \mathfrak{R}^n$ with components $g_\mu(i)$. The α -discounted n -state MDP is the special case where $A_\mu = \alpha P_\mu$, with P_μ being the matrix of transition probabilities $p_{ij}(\mu(i))$ corresponding to μ , and g_μ the vector with i th component $g_\mu(i) = \sum_{j=1}^n p_{ij}(\mu(i))g(i, \mu(i), j)$, the expected one-stage cost at state i under control $\mu(i)$. In policy iteration with cost function approximation, we compute approximations Φr_μ to J_μ using some policy evaluation method, and we use them to compute new policies based on (approximate) policy improvement, in analogy with Eq. (1.3).

In Section 2, we focus just on the approximate evaluation of a single fixed policy μ , so for notational simplicity we drop the subscript μ from Eq. (1.4). We thus focus on approximations of the form Φr for the fixed point of the linear mapping

$$TJ = g + AJ, \quad (1.5)$$

where A is an $n \times n$ matrix with nonnegative components a_{ij} , and $g \in \mathfrak{R}^n$. Several algorithmic approaches have been proposed for this purpose. The most popular are:

- (a) *Projected equation approach*: Here the approximation Φr is obtained by solving the projected equation

$$\Phi r = \Pi T(\Phi r), \quad (1.6)$$

where Π denotes projection onto the subspace S . The projection is with respect to a weighted Euclidean norm $\|\cdot\|_\xi$, where $\xi = (\xi_1, \dots, \xi_n)$ is a probability distribution with positive components (i.e., $\|J\|_\xi^2 = \sum_{i=1}^n \xi_i x_i^2$, where $\xi_i > 0$ for all i). The distribution vector ξ is sometimes (but by no means always) the steady-state probability vector of a Markov chain associated with the MDP.

- (b) *Aggregation approach*: Here the approximation Φr is obtained by solving an equation of the form

$$r = DT(\Phi r), \quad (1.7)$$

where D is a matrix whose rows are probability distributions. Contrary to the projected equation approach where Φ is arbitrary, an important restriction is that Φ must also be a matrix whose rows are probability distributions. The components of D and Φ are known as the disaggregation and aggregation probabilities. The vector r can be viewed as a cost vector of an aggregate problem that has s states and is defined by D and Φ (see Section 4 for details and interpretation).

There are also optimistic versions of these two approaches, where the corresponding equation (1.6) or (1.7) is solved approximately, with a finite number of iterations of some iterative method. Generally, the projected equation approach is associated with *temporal difference* (TD) methods, which originated in reinforcement learning with the works of Samuel [Sam59], [Sam67] on a checkers-playing program. The papers by Barto, Sutton, and Anderson [BSA83], and Sutton [Sut88] proposed the TD(λ) method, which motivated a lot of research in simulation-based DP, particularly following an early success with the backgammon playing program of Tesauro [Tes92]. The original papers did not make the connection of TD methods with the projected equation. Indeed for quite a long time it was not clear which mathematical problem TD(λ) was

aiming to solve! The mathematical convergence and other properties of $TD(\lambda)$, and its connections with the projected equation were clarified in the mid 90s through the works of several authors, including Gurvits, Lin, and Hanson [GLH94], Jaakkola, Jordan, and Singh [JJS94], Pineda [Pin97], Tsitsiklis and Van Roy [TsV97], [TsV99a]. More recent works have focused on the use of least squares-based TD methods, such as the LSTD (Least Squares Temporal Differences) method (Bradtke and Barto [BrB96]), and the LSPE (Least Squares Policy Evaluation) method (Bertsekas and Ioffe [BeI96]), which will be discussed later.

Note that the projected equation framework has a long history in the context of Galerkin methods for the approximate solution of high-dimensional or infinite-dimensional linear equations (partial differential, integral, inverse problems, etc; see e.g., [Kra72], [Fle84]). In fact much of the policy evaluation theory discussed in this paper applies to general projected equations arising in contexts beyond DP (see [BeY07], [BeY09], [Ber09a], [YuB10], [Yu10a,b]). However, the use of the Monte Carlo simulation and Monte Carlo integration ideas that are central in approximate DP is an important characteristic that differentiates the methods of the present paper from the Galerkin methodology, as currently practiced in the numerical analysis field.

The aggregation approach also has a long history in scientific computation and operations research. It was introduced in the simulation-based approximate DP context, mostly in the form of value iteration; see Singh, Jaakkola, and Jordan [SJJ94], [SJJ95], Gordon [Gor95], Tsitsiklis and Van Roy [TsV96], and Van Roy [Van06]. Currently the aggregation approach seems to be less popular, but as we will argue in this paper, it has some interesting advantages over the projected equation approach, even though there are restrictions in its applicability (see also the discussion in [Ber10a]).

Let us discuss briefly some alternative policy iteration approaches. In one of them, evaluation is based on the minimization of the (squared) Bellman equation error, as expressed by

$$\min_{r \in \mathbb{R}^s} \|\Phi r - T(\Phi r)\|^2, \quad (1.8)$$

where T is the mapping of Eq. (1.5). The summation over all the states in the squared norm above is usually approximated by a summation over a smaller “representative” set of states, often obtained in part through simulation. The optimality conditions for this minimization can be shown to be equivalent to the projected equation $\Phi r = \Pi \hat{T}(\Phi r)$, where \hat{T} is the mapping

$$\hat{T}J = TJ + A'(J - TJ),$$

which has the same fixed points as T , provided $I - A'$ is an invertible matrix (see [Ber10a], Section 6.8.4). Thus the Bellman equation error approach is a special case of the projected equation approach, although it leads to different simulation-based implementations. Sometimes it gives a better approximation error and sometimes worse. It has some interesting characteristics; see e.g., the discussion in [BeT96], Sections 6.10 and 6.11. However, it is less efficient when implemented by simulation (it requires the simultaneous generation of two state trajectories instead of one), and it is more susceptible to simulation noise, and for this reason it is not as popular as the projected equation approach at present. We do not consider this approach in this paper, referring to the literature for further discussion (e.g., [BeT96], Sections 6.10, 6.11, [Ber10a], Section 6.8.4, and Scherrer [Sch10]).

Among other approximate policy iteration methods, let us note the ones proposed recently by Thiery and Scherrer [ThS10], which are simulation-based implementations of λ -policy iteration (see also [Ber10a], Section 6.3.9). This latter algorithm combines elements of value and policy iteration, and was proposed by Bertsekas and Ioffe [BeI96] in conjunction with the LSPE method (see also [BeT96], Section 2.3.1). The

implementations of [ThS10] and [Ber10a] (Section 6.3.9) combine elements of several of the ideas that we will discuss in Sections 2 and 3, including solution of linear systems by simulation (Section 2.3), a λ parameter to control the bias-variance tradeoff (Section 2.3), convenience of exploration (Section 3.2), and optimistic operation (Section 3.3).

Contrary to approximate policy evaluation, which is an extensively researched and reasonably well understood subject, policy improvement with cost function approximation may exhibit complex behavior that is hard to analyze and can have seriously adverse impact on practical performance. An important question relates to the issue of inadequate exploration, i.e., the potential bias of the simulation through underrepresentation of states that are unlikely to occur under the policy being simulated. Another important question is whether the policy iteration process is seriously hampered by oscillations between poor policies, roughly similar to the attraction of gradient methods to poor local minima. As we will discuss later, there has been little apparent concern in the approximate DP/reinforcement learning literature about this possibility, even though it has been documented with several simple examples, and may be responsible for the poor performance of approximate policy iteration recently demonstrated in case studies involving the game of tetris (see the discussion of Section 3.5).

Let us also note that approximate policy iteration methods are based on the idea of “approximation in value space” and hence also on the hypothesis that a more accurate cost-to-go approximation will yield a better one-step or multistep lookahead policy through the policy improvement equation (1.3). This is a reasonable but by no means self-evident hypothesis, and may in fact not even be true in a given problem. Competitor methods are based on the idea of “approximation in policy space” and use simulation in conjunction with a gradient or random search method to approximate directly an optimal policy with a policy of a given parametric form. These methods do not aim at good cost function approximation through which a well-performing policy may hopefully be obtained. Rather they aim directly at finding a policy with good performance. The circumstances under which one type of method may yield a policy with superior performance over the other are unclear at present.

In this survey paper, we first address various computational issues of policy evaluation, including convergence and rate of convergence, and sensitivity to simulation noise in the case where Φ is (nearly) rank-deficient (Section 2). We then discuss additional issues of policy improvement, arising when policy evaluation is embedded within policy iteration, such as exploration, limited sampling/optimistic policy iteration, policy oscillation, and error bounds (Section 3). For policy evaluation, we first focus on the projected equation approach, since this is the one that is better known, and has the most complex and multifaceted behavior. We return to policy evaluation using the aggregation approach in Section 4. However, we take special care to introduce sufficient generality and abstraction into our presentation of Sections 2-3, so that a coherent and unified view of the commonalities and differences of the two approaches emerges.

Since we cover a lot of ground, to keep the size of the paper reasonable, we adopt a somewhat abstract and research-oriented style of exposition. We give proofs of some results, we refer to the published literature for others, and we speculate on a few other results that are plausible but not yet rigorously established.

2. PROJECTED EQUATION METHODS FOR POLICY EVALUATION

In this section, we discuss general aspects of policy evaluation using the projected equation

$$\Phi r = \Pi T(\Phi r) = \Pi(g + A\Phi r),$$

[cf. Eqs. (1.5)-(1.6)] and we give the forms of the two basic methods: *matrix inversion* and *iterative*. While these methods apply to general Galerkin approximation, we focus primarily on TD-like, simulation-based algorithms that were developed in approximate DP. For frameworks that are based on a more general viewpoint, we refer to our work on the solution of linear systems of equations [BeY07], [YuB10], [BeY09], variational inequalities [Ber09a], and least squares and inverse problems [WPB09], [PWB09]. One of the advantages of using a general matrix A in the definition of T is that it allows simultaneous application of our results to all types of Bellman equations (e.g., discounted, stochastic shortest path, average cost, semi-Markov, single-step, multistep, etc.).

We recall the basic characterization of the projection: $\Pi T(J)$ is the unique vector $Y \in S$ such that $Y - TJ$ is orthogonal to all vectors of S in the scaled geometry of the inner product $\langle x_1, x_2 \rangle = x_1' \Xi x_2$, where Ξ is the diagonal matrix with the components of the probability distribution vector $\xi = (\xi_1, \dots, \xi_n)$ on the diagonal. Since the projected equation is to find $J \in S$ such that TJ projects onto J , its solutions are characterized by the orthogonality principle: *J is a solution if and only if $J \in S$ and $J - TJ$ is orthogonal to all the columns of Φ which span the subspace S , i.e., $\Phi' \Xi (J - TJ) = 0$.* Thus, a vector J solves the projected equation $J = \Pi T(J)$ if and only if it is of the form Φr and

$$0 = \Phi' \Xi (J - TJ) = \Phi' \Xi (\Phi r - T(\Phi r)) = \Phi' \Xi (\Phi r - (g + A\Phi r)).$$

It follows that solving the projected equation is equivalent to solving the linear system $Cr = d$, where[†]

$$C = \Phi' \Xi (I - A) \Phi, \quad d = \Phi' \Xi g. \quad (2.1)$$

In this paper when referring to the projected equation, we make no distinction between the two equivalent forms $Cr = d$ and $J = \Pi T(J)$. The set of solutions of $J = \Pi T(J)$ is obtained from the set of solutions of $Cr = d$ by multiplication with Φ . Note that if Φ does not have rank s , C will not be invertible, and $Cr = d$ will not have a unique solution even if $J = \Pi T(J)$ has a unique solution (equivalently, $I - \Pi A$ is invertible). Moreover, solutions of $J = \Pi T(J)$ depend only on the projection norm and the subspace S , and not on the matrix Φ , which provides just an algebraic representation of S . Thus if $J = \Pi T(J)$ has a unique solution \tilde{J} , we have $\Phi r^* = \tilde{J}$ for every matrix Φ whose range space is S and solution r^* of the corresponding equation $Cr = d$.

2.1 Contraction Properties of the Projected Equation

The projected equation $Cr = d$ may in principle be solved by matrix inversion, $r^* = C^{-1}d$. This does not require any contraction properties for the operators T or ΠT ; invertibility of C is the only requirement. The equation may also be solved by one of several possible iterative methods [an important example is the fixed point iteration $\Phi r_{k+1} = \Pi T(\Phi r_k)$]. For some of these methods ΠT must be a contraction in order to guarantee convergence. However, even if T is a contraction with respect to some norm, ΠT need not be a contraction, unless Π is nonexpansive with respect to that norm, and this is a major source of complications for iterative methods. The following proposition provides a condition, which guarantees that ΠT is a contraction.

[†] This line of development of the projected equation can be generalized to the case where Π is the operator that projects on a convex strict subset of S . In this case, the projected equation is equivalent to a variational inequality (in place of $Cr = d$), thereby leading to a broad extension of the overall approximation approach; see [Ber09a] and the discussion in Section 3.4.

Proposition 2.1: Let $\xi = (\xi_1, \dots, \xi_n)$ be the steady-state probability vector of an irreducible transition probability matrix Q . Then T and ΠT are contraction mappings with respect to $\|\cdot\|_\xi$ if $a_{ij} \leq q_{ij}$ for all (i, j) , and there exists an index \bar{i} such that $a_{\bar{i}j} < q_{\bar{i}j}$ for all $j = 1, \dots, n$.

The proof of the proposition is given in [BeY09], and is based on the inequality

$$\|QJ\|_\xi \leq \|J\|_\xi, \quad \forall J \in \Re^n,$$

which holds for any irreducible transition matrix Q with invariant distribution ξ (see e.g., [BeT96], Lemma 6.4, or [Ber07], Lemma 6.3.1). In a discounted MDP case where $A = \alpha P$ and P is irreducible, by taking $Q = P$, the inequality implies that T is a contraction with respect to $\|\cdot\|_\xi$ of modulus α , which together with the nonexpansiveness of Π with respect to $\|\cdot\|_\xi$, implies that the same is true for ΠT . The more general condition $a_{\bar{i}j} < q_{\bar{i}j}$, $j = 1, \dots, n$, is satisfied in discounted MDP where Q is a somewhat modified version of P to allow for exploration (see [Ber10a]), and in a stochastic shortest path MDP, where A is a substochastic matrix corresponding to a proper policy (one for which the termination state is reached with probability 1 from any starting state); see [Ber07], [Ber10a]. There are also other related conditions that guarantee that ΠT is a contraction, and apply among others to average cost MDP (see [YuB06], [BeY09], or [Ber10a]).

If ΠT is a contraction, then the projected equation $J = \Pi T(J)$ has a unique solution Φr , which implies that if Φ has rank s , then $Cr = d$ has a unique solution so C is invertible. This and other properties are collected in the following proposition.

Proposition 2.2: Let ΠT be a contraction with respect to $\|\cdot\|_\xi$ and let Φ have rank s . Then:

(a) C is invertible and positive definite in the sense that

$$r'Cr > 0, \quad \forall r \neq 0.$$

(b) C has eigenvalues with positive real parts.

(c) The unique solution r^* of the projected equation $Cr = d$ and any solution J^* of the original equation $J = TJ$ satisfy

$$\|J^* - \Phi r^*\|_\xi \leq \frac{1}{\sqrt{1 - \alpha^2}} \|J^* - \Pi J^*\|_\xi, \quad (2.2)$$

where α is the modulus of contraction of ΠT .

The results of the proposition are well known in the theory of TD methods, and have been given in various forms by Tsitsiklis and Van Roy [TsV97], [TsV99a]. Proofs can also be found in [Ber10a]. The error bound of Eq. (2.2) is quite conservative (usually in MDP α is very close to 1). There are sharper error bounds, due to Yu and Bertsekas [YuB10] (see also Scherrer [Sch10]), which depend on the finer structure of A , Ξ , and the subspace S (not just on the contraction modulus of ΠT), and also apply to the case where ΠT is not a contraction.

2.2 Deterministic Iterative Methods for Projected Equations

We will now discuss iterative methods as alternatives to matrix inversion for solving the projected equation $Cr = d$. We first discuss deterministic methods that do not use simulation. These methods are not practical when the number of states is large, since they require high-dimensional matrix-vector computations. However, we will subsequently use Monte-Carlo integration ideas to construct simulation-based versions that require low-dimensional calculations only.

We focus on iterative methods of the form

$$r_{k+1} = r_k - \gamma G(Cr_k - d), \quad (2.3)$$

where G is a scaling matrix and γ is a positive stepsize.[†] This iteration is convergent if $I - \gamma GC$ is a contraction with respect to some norm (equivalently has all its eigenvalues strictly within the unit circle), but (under certain conditions) may also converge even if $I - \gamma GC$ is just nonexpansive with respect to some norm, which may happen if C is singular. Some interesting special choices of γ and G are given in the following three propositions. The last two apply to cases where C is singular.

Proposition 2.3: If ΠT is a contraction with respect to $\|\cdot\|_\xi$ and Φ has rank s , the iteration

$$r_{k+1} = r_k - (\Phi' \Xi \Phi)^{-1} (Cr_k - d) \quad (2.4)$$

is equivalent to the fixed point iteration $J_{k+1} = \Pi T(J_k)$, and converges to the unique solution of the projected equation $Cr = d$.

Proof: We use the projection formula $\Pi = \Phi(\Phi' \Xi \Phi)^{-1} \Phi' \Xi$, which is valid when Φ has rank s , together with the definition (2.1) of C and d , to write the iteration (2.4) as $J_{k+1} = \Pi T(J_k)$, where $J_k = \Phi r_k$, so convergence follows from the contraction property of ΠT . **Q.E.D.**

[†] Iterative methods that involve incremental changes along directions of the form $Gf(r)$ are common and fundamental for solving a system of equations $f(r) = 0$. They arise prominently in cases where $f(r)$ is the gradient of a cost function, or the case where the equation $f(r) = 0$ can be viewed as a monotone variational inequality [cf. Prop. 2.2(a); see [BeT89], Section 3.5.3]. The iterative methods (2.3) with G being positive definite symmetric, coincides with the class of projection methods for monotone variational inequalities; see [Ber09a], which adopts a variational inequality view for the projected equation in a more general context that allows projection on a convex subset of the subspace S . The paper [Ber09a] also interprets the scaling matrix G in terms of “feature scaling” (a representation of the subspace S in a different coordinate system). See also Yao and Liu [YaL08] for some algorithms related to (2.3), but developed from a different viewpoint.

Proposition 2.4: If ΠT is a contraction with respect to $\|\cdot\|_\xi$ and G is positive definite symmetric, there exists $\bar{\gamma} > 0$ such that for all $\gamma \in (0, \bar{\gamma}]$, the iteration

$$r_{k+1} = r_k - \gamma G(Cr_k - d)$$

converges to some solution of the projected equation.

Proof: We prove the result for the case where Φ has rank s . The proof for the general case is more complicated, and is based on the monotonicity property of the mapping $I - T$ (a consequence of the contraction property of ΠT) and the theory of projection methods for monotone variational inequalities (see [BeG82], [Ber09a]).

If G is symmetric and positive definite, the matrix $G^{1/2}$ exists and is symmetric and positive definite. Let $M = G^{1/2}CG^{1/2}$, and note that since C is positive definite, M is also positive definite, so from Prop. 2.2(b) it follows that its eigenvalues have positive real parts. The eigenvalues of M and GC are equal (with eigenvectors that are multiples of $G^{1/2}$ or $G^{-1/2}$ of each other), so the eigenvalues of GC have positive real parts. It follows that the eigenvalues of $I - \gamma GC$ lie strictly within the unit circle for sufficiently small $\gamma > 0$.

Q.E.D.

Proposition 2.5: The iteration $r_{k+1} = r_k - G(Cr_k - d)$, where

$$G = (C'\Sigma^{-1}C + \beta I)^{-1}C'\Sigma^{-1},$$

Σ is a positive definite symmetric matrix, and β is a positive scalar, converges to some solution of the projected equation, assuming at least one solution exists.

Proof: This follows from standard convergence results about the proximal point algorithm.† We give a proof for the easier case where C is invertible, so the projected equation has the unique solution $r^* = C^{-1}d$. Assuming C is invertible, the matrix $C'\Sigma^{-1}C$ is positive definite, so its eigenvalues $\lambda_1, \dots, \lambda_s$ are positive.

† One may write the iteration as

$$r_{k+1} = \arg \min_{r \in \mathbb{R}^s} [(Cr - d)'\Sigma^{-1}(Cr - d) + \beta\|r - r_k\|^2]$$

which is the proximal point algorithm applied to the problem of minimizing the cost function $(Cr - d)'\Sigma^{-1}(Cr - d)$. Then Prop. 2.5 follows from the results of Martinet [Mar70] and Rockafellar [Roc76]; see also [Ber09b], Prop. 6.5.1. An additional result is that $\{r_k\}$ will converge to a minimizing point of $(Cr - d)'\Sigma^{-1}(Cr - d)$ even if the projected equation $Cr = d$ has no solution.

Let $U\Lambda U'$ be the singular value decomposition of $C'\Sigma^{-1}C$, where $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_s\}$ and U is a unitary matrix ($UU' = I$; see [TrB97], [Str09]). We also have $C'\Sigma^{-1}C + \beta I = U(\Lambda + \beta I)U'$, so

$$GC = (U(\Lambda + \beta I)U')^{-1}U\Lambda U' = U(\Lambda + \beta I)^{-1}\Lambda U'.$$

It follows that the eigenvalues of GC are $\lambda_i/(\lambda_i + \beta)$, $i = 1, \dots, s$, and lie in the interval $(0, 1)$, so the eigenvalues of $I - GC$ also lie in the interval $(0, 1)$. **Q.E.D.**

Generally, the algorithms of the preceding three propositions converge linearly, at a rate that depends strongly on the choices of γ and G . However, we will argue later that when implemented by simulation, they all converge at the same asymptotic rate, which is dictated by the slow convergence rate of the simulation. Note that the algorithms converge under different conditions and offer alternative advantages. In particular, the algorithm of Prop. 2.5 is the most generally applicable, as it does not require that ΠT is a contraction or that C is invertible, and moreover does not require a stepsize choice. At the other extreme, the simplest algorithm is

$$r_{k+1} = r_k - \gamma(Cr_k - d), \quad (2.5)$$

the special case of Prop. 2.4 that corresponds to $G = I$. However, this algorithm requires that ΠT be a contraction and a choice of stepsize from within a range that is usually unknown a priori. The iteration (2.4) of Prop. 2.3 may be viewed as intermediate between the preceding two in terms of applicability and overhead per iteration.

2.3 Iterative Methods Using Simulation-Based Approximations

Unfortunately, for problems of very high dimension n , it is difficult to calculate C and d explicitly, because of the high-dimensional inner products involved in their definition [cf. Eq. (2.1)]. An alternative that has a long history in approximate DP (see the textbook references in Section 1), is to calculate simulation-based estimates \hat{C} and \hat{d} of C and d , respectively, based on a number of samples, and obtain an approximate solution

$$\hat{r} = \hat{C}^{-1}\hat{d}$$

by matrix inversion, assuming the inverse exists (this is the LSTD method, first proposed by Bradtke and Barto [BrB96], and followed up by Boyan [Boy02], and Nedić and Bertsekas [NeB03]).[†]

An alternative to LSTD is to solve by iteration a simulation-based approximation to the projected equation $Cr = d$: we approximate the deterministic iteration

$$r_{k+1} = r_k - \gamma G(Cr_k - d)$$

[cf. Eq. (2.3)] with

$$r_{k+1} = r_k - \gamma \hat{G}(\hat{C}r_k - \hat{d}), \quad (2.6)$$

where \hat{C} and \hat{d} are simulation-based estimates of C and d , γ is a positive stepsize, and \hat{G} is an $s \times s$ matrix, which may also be obtained by simulation. Assuming that $I - \gamma \hat{G} \hat{C}$ is a contraction, this iteration will

[†] Estimating C and d by simulation is motivated by the advantage that Monte-Carlo summation (or integration) holds in similar circumstances: using sampling to estimate sums of a very large number of terms (or complicated integrals, respectively), with computational complexity that depends on the “variance” of the terms added, and not on the number of terms.

yield a solution to the system $\hat{C}r = \hat{d}$, which will serve as a simulation-based approximation to a solution of the projected equation $Cr = d$. Chronologically, the first method of this type is LSPE [BeI96], where $\gamma = 1$ and \hat{G} is an approximation to $(\Phi'\Xi\Phi)^{-1}$ [cf. Eq. (2.4)]. The general method (2.6) may be viewed as a scaled version of LSPE, obtained by multiplying the LSPE direction with a suitable scaling matrix and by introducing a stepsize γ (see [Ber09a]).

Like LSTD, the iteration (2.6) may be viewed as a batch simulation approach: we first simulate to obtain \hat{C} , \hat{d} , and \hat{G} , and then solve the system $\hat{C}r = \hat{d}$ by the iterative method (2.6) rather than direct matrix inversion. An alternative is to iteratively update r as simulation samples are collected and used to form ever improving approximations to C and d . In particular, one or more iterations of the form (2.6) may be performed between collections of additional batches of simulation samples, to improve the accuracy of \hat{C} and \hat{d} . In the most extreme type of such an algorithm, the iteration (2.6) is used after a single new sample is collected. In this case, we approximate the iteration (2.3) by

$$r_{k+1} = r_k - \gamma G_k(C_k r_k - d_k), \quad (2.7)$$

where G_k , C_k , and d_k are simulation-based approximations, which are updated after a single new simulation sample is collected. For the purposes of further discussion, we will focus on this type of algorithm, with the understanding that there are related versions that use forms of batch simulation between iterations of the form (2.7) and have similar properties.

The convergence properties of iteration (2.7) under the condition

$$C_k \rightarrow C, \quad d_k \rightarrow d, \quad G_k \rightarrow G,$$

can be for the most part inferred from the results of Props. 2.3-2.5. In particular, if ΠT is a contraction and Φ has rank s , the proof of Prop. 2.3 carries through when G_k converges to $(\Phi'\Xi\Phi)^{-1}$. Similarly, if ΠT is a contraction, C is invertible, and G_k converges to a positive definite symmetric matrix, the proof of Prop. 2.4 carries through. Also if C is invertible and

$$G_k = (C_k'\Sigma_k^{-1}C_k + \beta I)^{-1}C_k'\Sigma_k^{-1},$$

where $C_k \rightarrow C$, and $\Sigma_k \rightarrow \Sigma$ for a symmetric positive definite matrix Σ , the proof of Prop. 2.5 carries through (a formal convergence analysis for the case where C is singular is not available at present, even though the deterministic version of the iteration is convergent as per Prop. 2.5).

We now discuss a simulation mechanism for estimating C and d , proposed in [BeY09]. It is similar to methods used traditionally in TD methods for policy evaluation, which are based on sampling the Markov chain corresponding to the policy. However, our simulation mechanism may use two Markov chains, and is well suited for sampling with enhanced exploration (see the discussion in Section 3.2). We write C and d as expected values with respect to ξ :

$$C = \sum_{i=1}^n \xi_i \phi(i) \left(\phi(i) - \sum_{j=1}^n a_{ij} \phi(j) \right)', \quad d = \sum_{i=1}^n \xi_i \phi(i) g(i), \quad (2.8)$$

where ξ_i , $g(i)$, and a_{ij} are the components of ξ , g , and A , respectively, and $\phi(i)$ denotes the column of Φ' that corresponds to state i [$\phi(i)'$ is the i th row of Φ]. We generate a sequence of indices $\{i_0, i_1, \dots\}$, and a sequence of transitions $\{(i_0, j_0), (i_1, j_1), \dots\}$. We use any probabilistic mechanism for this, subject to the following two requirements (cf. Fig. 2.1):

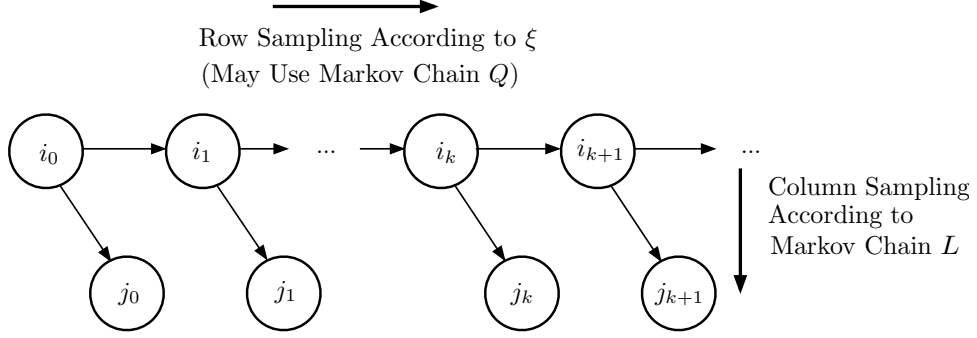


Figure 2.1. The Markov chain-based simulation methodology consists of (a) generating a sequence of indices $\{i_0, i_1, \dots\}$ according to the distribution ξ (a suitable Markov chain Q may be used for this, but this is not a requirement), and (b) generating a sequence of transitions $\{(i_0, j_0), (i_1, j_1), \dots\}$ using a Markov chain L . It is possible that $j_k = i_{k+1}$, but this is not necessary. Moreover, even when A is related to the transition probability matrix of a Markov chain of an underlying MDP, Q and L need not be equal to that matrix.

- (1) **Row sampling:** The sequence $\{i_0, i_1, \dots\}$ is generated according to the distribution ξ , which defines the projection norm $\|\cdot\|_\xi$, in the sense that with probability 1,

$$\lim_{k \rightarrow \infty} \frac{\sum_{t=0}^k \delta(i_t = i)}{k+1} = \xi_i, \quad i = 1, \dots, n, \quad (2.9)$$

where $\delta(\cdot)$ denotes the indicator function [$\delta(E) = 1$ if the event E has occurred and $\delta(E) = 0$ otherwise]. Row sampling requires that each state i is generated with a relative frequency ξ_i specified by the projection norm $\|\cdot\|_\xi$. We will discuss later (Section 3.2) how one may design row sampling schemes corresponding to values ξ_i that adequately address the issue of exploration.

- (2) **Column sampling:** The sequence $\{(i_0, j_0), (i_1, j_1), \dots\}$ is generated according to a certain stochastic matrix L with transition probabilities l_{ij} in the sense that with probability 1,

$$\lim_{k \rightarrow \infty} \frac{\sum_{t=0}^k \delta(i_t = i, j_t = j)}{\sum_{t=0}^k \delta(i_t = i)} = l_{ij}, \quad i, j = 1, \dots, n. \quad (2.10)$$

We require that

$$l_{ij} > 0 \quad \text{if} \quad a_{ij} > 0, \quad (2.11)$$

so transitions (i, j) with $a_{ij} > 0$ are generated with positive relative frequency. Other than this requirement, l_{ij} need not be related to a_{ij} . For example, in an α -discounted problem, the column sampling probabilities need not be equal to the transition probabilities of the policy under evaluation. This turns out to be important because it allows exploration enhancement through the choice of the column sampling probabilities l_{ij} (see the subsequent discussion of λ -methods).

At time k , we approximate C and d with

$$C_k = \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) \left(\phi(i_t) - \frac{a_{i_t j_t}}{l_{i_t j_t}} \phi(j_t) \right)', \quad d_k = \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) g(i_t). \quad (2.12)$$

To show that this is a valid approximation, we count the number of times an index occurs and after collecting terms, we write Eq. (2.12) as

$$C_k = \sum_{i=1}^n \hat{\xi}_{i,k} \phi(i) \left(\phi(i) - \sum_{j=1}^n \hat{l}_{ij,k} \frac{a_{ij}}{l_{ij}} \phi(j) \right)', \quad d_k = \sum_{i=1}^n \hat{\xi}_{i,k} \phi(i) g(i), \quad (2.13)$$

where

$$\hat{\xi}_{i,k} = \frac{\sum_{t=0}^k \delta(i_t = i)}{k+1}, \quad \hat{\ell}_{ij,k} = \frac{\sum_{t=0}^k \delta(i_t = i, j_t = j)}{\sum_{t=0}^k \delta(i_t = i)}.$$

In view of the assumption

$$\hat{\xi}_{i,k} \rightarrow \xi_i, \quad \hat{\ell}_{ij,k} \rightarrow \ell_{ij}, \quad i, j = 1, \dots, n,$$

[cf. Eqs. (2.9) and (2.10)], by comparing Eqs. (2.8) and (2.13), we see that $C_k \rightarrow C$ and $d_k \rightarrow d$.[†]

In the iterative method

$$r_{k+1} = r_k - \gamma G_k (C_k r_k - d_k),$$

[cf. Eq. (2.7)], in addition to C_k , d_k , one may simultaneously obtain G_k with the preceding simulation mechanism. In particular, one may choose $\gamma = 1$ and

$$G_k = \left(\frac{1}{k+1} \sum_{t=0}^k \phi(i_t) \phi(i_t)' \right)^{-1}, \quad (2.17)$$

which is a simulation-based approximation to $(\Phi' \Xi \Phi)^{-1}$ [cf. the iteration (2.4) of Prop. 2.3]. This is the LSPE method first proposed by Bertsekas and Ioffe [BeI96], and followed up by Nedić and Bertsekas [NeB03], Bertsekas, Borkar, and Nedić [BBN04], and Yu and Bertsekas [YuB06] (the method was also described in [BeI96] and [BeT96] as an implementation of the λ -policy iteration method, and was used to address a challenging tetris application; the name LSPE was introduced in [NeB03]).

The LSTD and LSPE methods were preceded and motivated by the TD(0) method, which for an α -discounted problem has the form

$$r_{k+1} = r_k - \gamma_k \phi(i_k) q_{k,k}, \quad (2.18)$$

where $q_{k,k}$ is the TD given by Eq. (2.15), and γ_k is a diminishing stepsize that diminishes at an appropriate rate, such as $\gamma_k = O(1/k)$. Thus, in view of Eq. (2.14), TD(0) has a similar form to the simple iteration (2.5), but differs in two important ways: first, the stepsize γ_k is diminishing, and second, C_k and d_k are single-sample estimates of C and d (only the last simulation sample is used rather than all of the k samples). TD(0) may be viewed as a stochastic approximation/Robbins-Monro scheme for solving the projected equation

[†] In the discounted case where $A = \alpha P$, the approximate projected equation $C_k r = d_k$ can be written as

$$C_k r - d_k = \sum_{t=0}^k \phi(i_t) q_{k,t} = 0, \quad (2.14)$$

where

$$q_{k,t} = \phi(i_t)' r_k - \alpha \phi(i_{t+1})' r_k - g(i_t). \quad (2.15)$$

The scalar $q_{k,t}$ is the so-called *temporal difference*, associated with r_k and transition (i_t, i_{t+1}) . It may be viewed as a sample of a residual term arising in the projected Bellman's equation. More specifically, we have

$$C r_k - d = \Phi' \Xi (\Phi r_k - \alpha P \Phi r_k - g). \quad (2.16)$$

The three terms in the definition (2.15) of the temporal difference $q_{k,t}$ can be viewed as samples [associated with the transition (i_t, i_{t+1})] of the corresponding three terms in the expression $\Xi (\Phi r_k - \alpha P \Phi r_k - g)$ in Eq. (2.16). For the purposes of this paper, temporal differences and the form of the approximate projected equation (2.14) play no substantive role, and they will not be used in any way.

$Cr = d$, and in this sense it is fundamentally different than the LSTD and LSTE methods, which are Monte-Carlo versions of matrix inversion and iterative deterministic methods, which use simulation-based estimates in place of the hard-to-compute quantities C and d . Thus TD(0) requires an entirely different type of mathematical convergence analysis than LSTD and LSPE (see e.g., [TsV97]). Let us also note a scaled form of TD(0), which uses a positive definite symmetric scaling matrix G in Eq. (2.18), and has been proposed and analyzed by Choi and Van Roy [ChV06]. Similar to TD(0), their method uses single-sample estimates of C and d in the algorithm of Prop. 2.4.

Note that there is a lot of flexibility for row and column sampling to satisfy Eqs. (2.9)-(2.10). For example, to satisfy Eq. (2.9), the indices i_t do not need to be sampled independently according to ξ . Instead, it may be convenient to introduce an irreducible Markov chain with transition matrix Q , states $1, \dots, n$, and ξ as its steady-state probability vector, and to start at some state i_0 and generate the sequence $\{i_0, i_1, \dots\}$ as a single infinitely long trajectory of the chain. For the transition sequence, we may optionally let $j_k = i_{k+1}$ for all k , in which case L would be identical to Q , but in general this is not essential. In the MDP literature the most common choice is $Q = L = P$, where P is the transition matrix of the associated Markov chain. However, in the context of policy iteration it soon became clear that it is important to do row sampling using a different Markov chain, one that generates by simulation all states with sufficient frequency (see the discussion of exploration in Section 3.2).

Note also that multiple simulated sequences can be used to form the simulation estimates (2.12) of C and d . For example, in the Markov chain-based sampling schemes, we can generate multiple infinitely long trajectories of the chain, starting at several different states. This will work even if the chain has multiple recurrent classes, as long as there are no transient states and at least one trajectory is started from within each recurrent class. Again ξ will be a steady-state probability vector of the chain, and need not be known explicitly. Note also that using multiple trajectories may be interesting even if there is a single recurrent class, for at least two reasons:

- (a) The generation of trajectories may be parallelized among multiple processors, resulting in significant speedup.
- (b) The empirical frequencies of occurrence of the states may approach the steady-state probabilities more quickly; this is particularly so for large and “stiff” Markov chains.

We finally note that there are several variants of the Markov chain-based simulation method outlined above (see [BeY07] or more details). For example, if instead of $g(i_t)$ we obtain the sample $g(i_t, j_t)$, the vector d_k in Eq. (2.12) should be

$$d_k = \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) g(i_t, j_t).$$

As another example, zero mean noise with appropriate independence properties may be added to a_{ij} and $g(i)$. In this case, the estimates of C and d are given by

$$C_k = \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) \left(\phi(i_t) - \sum_{j=1}^n (a_{ij} + \zeta_t(j)) \phi(j) \right)', \quad d_k = \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) (g(i_t) + \theta_t), \quad (2.19)$$

where for each j , $\zeta_t(j)$ is a sequence of random variables such that, with probability 1,

$$\lim_{k \rightarrow \infty} \frac{\sum_{t=0}^k \delta(i_t = i) \zeta_t(j)}{\sum_{t=0}^k \delta(i_t = i)} = 0, \quad \forall i, j = 1, \dots, n, \quad (2.20)$$

and θ_t is a sequence of random variables such that, with probability 1,

$$\lim_{k \rightarrow \infty} \frac{\sum_{t=0}^k \delta(i_t = i) \theta_t}{\sum_{t=0}^k \delta(i_t = i)} = 0, \quad \forall i = 1, \dots, n. \quad (2.21)$$

This variant can be used in situations where the components a_{ij} and $g(i)$ represent the expected values of random variables whose samples can be conveniently simulated with additive “noises” $\zeta_t(j)$ and θ_t , respectively, such that Eqs. (2.20) and (2.21) hold with probability 1.

λ -Versions of Simulation Methods

We will now briefly discuss “ λ -versions” of the LSTD, LSPE, and TD algorithms, which use a parameter $\lambda \in [0, 1)^\dagger$ and aim to solve the geometrically weighted multistep version of Bellman’s equation $J = T^{(\lambda)}J$, where T is defined by

$$T^{(\lambda)} = (1 - \lambda) \sum_{t=0}^{\infty} \lambda^t T^{t+1}. \quad (2.22)$$

The best known example is TD(λ) for α -discounted problems [Sut88]. The motivation for using $\lambda > 0$ is that the error $\|J^* - \Phi r^*\|_\infty$ corresponding to the solution r_λ^* of the projected equation

$$\Phi r = \Pi T^{(\lambda)}(\Phi r)$$

tends to decrease as λ increases, and indeed an improved error bound similar to the one of Prop. 2.2(c) can be shown for $\lambda > 0$ (see e.g., [Ber07], Prop. 6.3.5). However, the number of samples needed to approximate well the projected equation tends to increase; this is the well-known bias-variance tradeoff, and we refer to textbook sources, such as [BeT96], [SuB98], and [Ber07] for a discussion. In this survey, we mostly focus on $\lambda = 0$, but our algorithms and qualitative conclusions generally apply to $\lambda \in [0, 1)$ [for further discussion, see [BeY09], where extensions of LSTD(λ), LSPE(λ), and TD(λ) beyond the MDP context are also provided]. Let us note that the projected equation $\Phi r = \Pi T^{(\lambda)}(\Phi r)$ can be written as a linear equation of the form $C^{(\lambda)}r = d^{(\lambda)}$, where similar to Eq. (2.1), $C^{(\lambda)}$ and $d^{(\lambda)}$ are defined by

$$C^{(\lambda)} = \Phi' \Xi (I - A^{(\lambda)}) \Phi, \quad d^{(\lambda)} = \Phi' \Xi g^{(\lambda)},$$

with

$$A^{(\lambda)} = (1 - \lambda) \sum_{t=0}^{\infty} \lambda^t A^{t+1}, \quad g^{(\lambda)} = \sum_{t=0}^{\infty} \lambda^t A^t g.$$

Similar to the case $\lambda = 0$, there are simulation procedures that produce approximations $C_k^{(\lambda)}$ and $d_k^{(\lambda)}$ to $C^{(\lambda)}$ and $d^{(\lambda)}$, respectively, thereby leading to corresponding matrix inversion and iterative methods. However, for these procedures to be algorithmically convenient, some restrictions must be observed. In particular, *it is essential that row and column sampling are generated by a common process*, so that

$$j_k = i_{k+1}, \quad k = 0, 1, \dots$$

where each transition (i_k, i_{k+1}) is generated using the matrix L [which we recall may be unrelated to the matrix A of the system $J = g + AJ$, other than $\ell_{ij} > 0$ if $a_{ij} > 0$, cf. Eq. (2.11)]. In this case, approximations

[†] It is also possible to let $\lambda = 1$ with appropriate changes in various formulas that follow, but for simplicity we will not consider this possibility here.

$C_k^{(\lambda)}$ and $d_k^{(\lambda)}$ to $C^{(\lambda)}$ and $d^{(\lambda)}$, respectively, are obtained from the following recursive algorithm, first given by Bertsekas and Yu [BeY09]:[‡]

$$C_k^{(\lambda)} = \left(1 - \frac{1}{k+1}\right) C_{k-1}^{(\lambda)} + \frac{1}{k+1} z_k \left(\phi(i_k) - \frac{a_{i_k i_{k+1}}}{\ell_{i_k i_{k+1}}} \phi(i_{k+1}) \right)', \quad (2.23)$$

$$d_k^{(\lambda)} = \left(1 - \frac{1}{k+1}\right) d_{k-1}^{(\lambda)} + \frac{1}{k+1} z_k g(i_k), \quad (2.24)$$

where z_k is generated by

$$z_k = \lambda \frac{a_{i_{k-1} i_k}}{\ell_{i_{k-1} i_k}} z_{k-1} + \phi(i_k). \quad (2.25)$$

An analysis of the convergence $C_k^{(\lambda)} \rightarrow C^{(\lambda)}$ and $d_k^{(\lambda)} \rightarrow d^{(\lambda)}$ has been given in several sources under different assumptions:

- (a) In Nedić and Bertsekas [NeB03] for the case of policy evaluation in an α -discounted problem, where A/α is the transition probability matrix P and row and column sampling are done according to that matrix ($Q = L = P$). This is the case where there is no exploration (cf. the subsequent discussion in Section 3.2).
- (b) In Bertsekas and Yu [BeY09], assuming that $\lambda \max_{(i,j)}(a_{ij}/\ell_{ij}) < 1$ (where we adopt the convention $0/0 = 0$), in which case the vectors z_k of Eq. (2.25) are generated by a contractive process and remain bounded. This covers the common case in MDP where $L = (1 - \epsilon)P + \epsilon\bar{P}$ where $\epsilon > 0$ is a small constant, P is a transition matrix corresponding to the policy under evaluation, and \bar{P} is a transition probability matrix that induces exploration. The essential restriction here is that λ should be no more than $(1 - \epsilon)$.
- (c) In Yu [Yu10a,b], for all $\lambda \in [0, 1]$ and $\lambda \sum_{j=1}^n a_{ij} < 1$ for all i , in which case the vectors z_k typically become unbounded as k increases when $\lambda \max_{(i,j)}(a_{ij}/\ell_{ij}) > 1$. Mathematically, this is the most challenging analysis, and involves interesting stochastic phenomena.

Much of what will be said later for projected equation methods where $\lambda = 0$ also applies to methods with $\lambda > 0$.

We finally mention a potential drawback of the row and column sampling methods: they generally require explicit knowledge of the ratios a_{ij}/ℓ_{ij} . Sometimes, in “model-free” applications of approximate DP, the matrix A is not explicitly known, and can only be computed through simulation, so in this case row and column sampling may not apply, at least in the general form given here. We will discuss this issue further in Section 3.8, in the context of Q-learning.

2.4 Simulation by Importance Sampling

We have seen that a core issue in approximate policy evaluation using the projected equation is to estimate the matrix C by sampling (also to estimate d , but this is similar and simpler). We may thus consider general

[‡] If $g(i_k)$ is replaced by a simulation sample $g(i_k, i_{k+1})$ the formula for updating $d_k^{(\lambda)}$ should be

$$d_k^{(\lambda)} = \left(1 - \frac{1}{k+1}\right) d_{k-1}^{(\lambda)} + \frac{1}{k+1} z_k \frac{a_{i_k i_{k+1}}}{\ell_{i_k i_{k+1}}} g(i_k, i_{k+1}).$$

sampling methods for estimating matrices of the form

$$Z = DH\Phi$$

for an $s \times n$ matrix D , an $n \times n$ matrix H , and an $n \times s$ matrix Φ (in our case $D = \Phi'\Xi$ and $H = I - A$, or $H = I$ and $H = A$). When n is very large and s is relatively small, we may estimate each of the components of Z separately. Denoting by $d_{\ell i}$, h_{ij} , and ϕ_{jq} the components of D , H , and Φ , respectively, the (ℓ, q) -th component of Z is

$$z_{\ell q} = \sum_{i=1}^n \sum_{j=1}^n d_{\ell i} h_{ij} \phi_{jq}, \quad (2.26)$$

so it is a sum of a very large number of terms that can be estimated by simulation using Monte Carlo integration techniques aimed at variance reduction, such as importance sampling.

In particular, we generate samples $(i_0, j_0), \dots, (i_k, j_k)$ according to some sampling distribution $\{\zeta_{ij} \mid i, j = 1, \dots, n\}$, and we estimate $z_{\ell q}$ by

$$\hat{z}_{\ell q} = \frac{1}{k+1} \sum_{t=0}^k \frac{d_{\ell i_t} h_{i_t j_t} \phi_{j_t q}}{\zeta_{i_t j_t}}.$$

According to the principles of importance sampling (see e.g., Liu [Liu01], Rubinstein and Kroese [RuK08]), for fixed ℓ, q , the distribution ζ_{ij} should approximate in form the generic term $d_{\ell i} h_{ij} \phi_{jq}$ of the sum, i.e.,

$$\zeta_{ij} \approx \frac{d_{\ell i} h_{ij} \phi_{jq}}{N},$$

where N is a normalization constant [we assume here that $d_{\ell i} h_{ij} \phi_{jq}$ is nonnegative for all (i, j) – otherwise the sum (2.26) should be split into the two sums corresponding to the positive and the negative terms].

A key issue here is the design of the distribution ζ_{ij} , which by necessity should exploit the special/sparsity structure of the matrices D , B , and Φ , to localize the distribution in a relatively small portion of the set of pairs (i, j) . When a suitable importance sampling distribution can be obtained, the variance of the estimates $\hat{z}_{\ell q}$ and the attendant number of samples for a given accuracy can be dramatically reduced (by several orders of magnitude). We do not pursue this topic further, but refer to the papers [WPB09] and [PWB09] for sampling distribution design techniques that have been used successfully in some very large problems within a related least squares estimation context.

2.5 Convergence Rate Issues

Let us now discuss the choice of γ and G_k in the iterative method

$$r_{k+1} = r_k - \gamma G_k (C_k r_k - d_k), \quad (2.27)$$

from the convergence rate point of view. It can be easily verified with simple examples that in the deterministic case where $G_k \equiv G$, $C_k \equiv C$, and $d_k \equiv d$, the values of γ and G affect significantly the convergence rate of the deterministic iteration (2.3). This is evident when $I - \gamma GC$ is a contraction, since then the convergence rate is linear and is determined by the spectral radius of $I - \gamma GC$. Surprisingly, however, when G_k , C_k , and d_k are generated by the preceding simulation schemes, the asymptotic convergence rate of the simulation-based iteration (2.27) does not depend on the choices of γ and G . Indeed it can be proved that *the iteration (2.27) converges at the same asymptotic rate, for all choices of γ and G , as long as $I - \gamma GC$*

is a contraction (although the short-term convergence rate may be significantly affected by the choices of γ and G).

The reason is that the deterministic iteration has a linear convergence rate (since it involves the contraction $I - \gamma GC$), which is fast relative to the slow convergence rate of the simulation-generated G_k , C_k , and d_k . Thus the simulation-based iteration (2.27) operates on two time scales (see, e.g., Borkar [Bor08], Ch. 6): the slow time scale at which G_k , C_k , and d_k change, and the fast time scale at which r_k adapts to changes in G_k , C_k , and d_k . As a result, essentially, there is convergence in the fast time scale before there is appreciable change in the slow time scale. Roughly speaking, r_k “sees G_k , C_k , and d_k as effectively constant,” so that for large k , r_k is essentially equal to the corresponding limit of iteration (2.27) with G_k , C_k , and d_k held fixed. When C_k is invertible, this limit is $C_k^{-1}d_k$. It follows that the sequence r_k generated by the scaled LSPE iteration (2.27) “tracks” the sequence $C_k^{-1}d_k$ generated by the LSTD iteration in the sense that

$$\|r_k - C_k^{-1}d_k\| \ll \|r_k - r^*\|, \quad \text{for large } k,$$

independent of the choice of γ and the scaling matrix G that is approximated by G_k (see also [Ber09a] for further discussion). Results of this type were conjectured for the LSPE method in [BBN04] and were formally proved in [YuB06] for discounted and average cost problems, and for the special choice (2.17) of G_k . The proof for more general cases where C is invertible should be similar, but has not yet been documented.

2.6 Regression-Based Methods for Projected Equations

An important concern in LSTD is to ensure that the simulation-induced error

$$e_k = \Phi(C_k^{-1}d_k - C^{-1}d)$$

is not excessively large. Then the low-dimensional error $C_k^{-1}d_k - C^{-1}d$ is typically also large (the reverse is not true: $C_k^{-1}d_k - C^{-1}d$ may be large without e_k being large). In the lookup table representation case ($\Phi = I$) a large error e_k may be traced directly to the simulation error in evaluating C and d , combined with near singularity of $\Xi(I - A)$. In the compact representation case ($\Phi \neq I$), the effect of near singularity of C on the high-dimensional error e_k is more complex, but is also primarily due to the same causes.† In what follows we consider approaches to reduce the low-dimensional error $C_k^{-1}d_k - C^{-1}d$ with the understanding that they will also be effective in reducing the high-dimensional error e_k , when the latter is very large.

Example 2.1

To get a rough sense of the potential effect of the simulation error in LSTD, consider the approximate inversion of a small nonzero number c , which is estimated with simulation error ϵ . The absolute and relative errors are

$$E = \frac{1}{c + \epsilon} - \frac{1}{c}, \quad E_r = \frac{E}{1/c}.$$

† Near-singularity of C , causing large low-dimensional errors $C_k^{-1}d_k - C^{-1}d$, may be due either to the columns of Φ being nearly linearly dependent or to the matrix $\Xi(I - A)$ being nearly singular [cf. the formula $C = \Phi'\Xi(I - A)\Phi$]. However, near-linear dependence of the columns of Φ will not affect the high-dimensional error e_k . The reason is that e_k depends only on the subspace S and not its representation in terms of the matrix Φ . In particular, if we replace Φ with a matrix ΦB where B is an $s \times s$ invertible scaling matrix, the subspace S will be unaffected and it can be verified that the errors e_k will also be unaffected. On the other hand, near singularity of the matrix $I - A$ may affect significantly e_k . Large variations in the size of the diagonal components of Ξ may also affect significantly e_k , although this dependence is complicated by the fact that Ξ appears not only in the formula $C = \Phi'\Xi(I - A)\Phi$ but also in the formula $d = \Phi'\Xi g$.

By a first order Taylor series expansion around $\epsilon = 0$, we obtain for small ϵ

$$E \approx \frac{\partial(1/(c+\epsilon))}{\partial\epsilon} \Big|_{\epsilon=0} \epsilon = -\frac{\epsilon}{c^2}, \quad E_r \approx -\frac{\epsilon}{c}.$$

Thus for the estimate $\frac{1}{c+\epsilon}$ to be reliable, we must have $|\epsilon| \ll |c|$. If N independent samples are used to estimate c , the variance of ϵ is proportional to $1/N$, so for a small relative error, N must be much larger than $1/c^2$. Thus as c approaches 0, the amount of sampling required for reliable simulation-based inversion increases very fast. This is similar to the roundoff error-related difficulties of nearly singular matrix inversion.

To counter the potential difficulties associated with a near-singular matrix C , we may estimate r^* by a form of regularized regression, which works even if the projected equation is singular. In this approach, suggested by Wang, Polydorides, and Bertsekas [WPB09], instead of solving the system $C_k r = d_k$, we use regression based on an associated linear model that properly encodes the effect of the simulation noise. We will see, however, that this method comes at a price: it introduces a bias of the estimate of r towards some fixed vector/prior guess \bar{r} . To reduce or eliminate this bias it is necessary to adopt an iterative regularization approach: start with some \bar{r} , obtain an estimate \hat{r}_k , replace \bar{r} by \hat{r}_k , and repeat for any number of times. This gives rise to an iterative method, which turns out to be a special case of the method (2.27).

Given simulation-based estimates C_k and d_k to C and d , respectively, let us write the projected form of Bellman's equation $d = Cr$ as

$$d_k = C_k r + e_k, \tag{2.28}$$

where e_k is the vector

$$e_k = (C - C_k)r + d_k - d,$$

which we view as ‘‘simulation noise.’’ We then estimate the solution $r^* = C^{-1}d$ based on Eq. (2.28) by using regression. In particular, we choose r by solving the least squares problem:

$$\min_r [(d_k - C_k r)' \Sigma^{-1} (d_k - C_k r) + \beta \|r - \bar{r}\|^2], \tag{2.29}$$

where \bar{r} is an *a priori* estimate of r^* , Σ is some positive definite symmetric matrix, and β is a positive scalar. By setting to 0 the gradient of the least squares objective in Eq. (2.29), we can find the solution in closed form:

$$\hat{r}_k = (C_k' \Sigma^{-1} C_k + \beta I)^{-1} (C_k' \Sigma^{-1} d_k + \beta \bar{r}). \tag{2.30}$$

A suitable choice of \bar{r} may be some heuristic guess based on intuition about the problem, or it may be the parameter vector corresponding to the estimated cost vector $\Phi \bar{r}$ of a similar policy (for example a preceding policy in an approximate policy iteration context). One may try to choose Σ in special ways to enhance the quality of the estimate of r^* , but we will not consider this issue here, and the subsequent analysis in this section does not depend on the choice of Σ , as long as it is positive definite and symmetric.

The quadratic $\beta \|r - \bar{r}\|^2$ in Eq. (2.29) is a *regularization term*, and has the effect of ‘‘biasing’’ the estimate \hat{r}_k towards the *a priori* guess \bar{r} . The proper size of β is not clear (a large size reduces the effect of near singularity of C_k , and the effect of the simulation errors $C_k - C$ and $d_k - d$, but may cause a large ‘‘bias’’). However, this is typically not a major difficulty in practice, because trial-and-error experimentation with different values of β involves low-dimensional linear algebra calculations once C_k and d_k become available.

To eliminate the bias of \hat{r}_k towards \bar{r} , one possibility is to adopt an iterative regularization approach: start with some \bar{r} , obtain \hat{r}_k , replace \bar{r} by \hat{r}_k , and repeat for any number of times. This turns LSTD to the iterative method,

$$r_{k+1} = (C_k' \Sigma_k^{-1} C_k + \beta I)^{-1} (C_k' \Sigma_k^{-1} d_k + \beta r_k),$$

where Σ_k are positive definite symmetric matrices that converge to a positive definite matrix Σ , and β is a positive scalar [cf. Eq. (2.30)]. Equivalently, this method can be written as

$$r_{k+1} = r_k - G_k(C_k r_k - d_k), \quad (2.31)$$

where

$$G_k = (C'_k \Sigma_k^{-1} C_k + \beta I)^{-1} C'_k \Sigma_k^{-1}, \quad (2.32)$$

and can be recognized as a simulation-based version of the special case of the class of LSPE-type methods (2.27), whose convergence properties have been discussed in Prop. 2.5.

The salient feature of the method (2.31)-(2.32) is that, based on Prop. 2.5 and experimental evidence, it can be expected to converge to a solution of the projected equation $Cr = d$ even if the mapping ΠT is not a contraction, as long as C is invertible. We will see in the next section that this generality is particularly important in the context of policy iteration where the simulation is done by using exploration-enhanced Markov chains. Furthermore, the iteration is well-suited for problems where C is nearly singular.

We also note an alternative and somewhat simpler regularization approach, whereby we approximate the equation $C_k r = d_k$ by

$$(C_k + \beta I)\hat{r} = d_k + \beta \bar{r}, \quad (2.33)$$

where β is a positive scalar and \bar{r} is some guess of the solution $r^* = C^{-1}d$. This suggests an iteration of the form $r_{k+1} = r_k - G_k(C_k r_k - d_k)$ where the scaling matrix G_k is $G_k = (C_k + \beta I)^{-1}$, i.e.,

$$r_{k+1} = r_k - (C_k + \beta I)^{-1}(C_k r_k - d_k). \quad (2.34)$$

The convergence of this iteration to r^* can be proved, assuming that C is positive definite, based on the fact $C_k \rightarrow C$ and standard convergence results for the proximal point algorithm ([Mar70], [Roc76]). Note that by contrast with the iteration (2.31)-(2.32), the positive definiteness of C is essential for invertibility of $C_k + \beta I$ and for convergence of Eq. (2.34). We refer to [Ber09a] for more details on this method.

3. POLICY ITERATION ISSUES

We will now discuss various issues that arise within a policy iteration context: the effects of inadequate exploration, of inexact policy evaluation, and of policy oscillation. For concreteness, we focus on the α -discounted MDP given in the introduction, where the cost vector J_μ of a policy μ is the unique fixed point of the mapping T_μ of Eq. (1.1), and the optimal cost vector J^* is the unique fixed point of the mapping T of Eq. (1.2).

We consider a form of policy iteration where we compute simulation-based approximations Φr_μ to the cost vectors J_μ of stationary policies μ with some policy evaluation method. We use Φr_μ to compute new policies based on (approximate) policy improvement using the formula $T_{\bar{\mu}}(\Phi r_\mu) = T(\Phi r_\mu)$, i.e.,

$$\bar{\mu}(i) \in \arg \min_{u \in U(i)} \sum_{j=1}^n p_{ij}(u) (g(i, u, j) + \alpha \phi(j)' r_\mu), \quad i = 1, \dots, n, \quad (3.1)$$

where $\phi(j)'$ is the row of Φ that corresponds to state j . The method terminates with μ if $T_\mu(\Phi r_\mu) = T(\Phi r_\mu)$. We then repeat with μ replaced by $\bar{\mu}$.

The theoretical basis for the preceding method was given in Prop. 6.2 of Bertsekas and Tsitsiklis [BeT96], where it was shown that if the policy evaluation is accurate to within δ (in the sup-norm sense $\|\Phi r_\mu - J_\mu\|_\infty \leq \delta$), then the method will yield a sequence of policies $\{\mu^k\}$ such that

$$\limsup_{k \rightarrow \infty} \|J_{\mu^k} - J^*\|_\infty \leq \frac{2\alpha\delta}{(1-\alpha)^2}. \quad (3.2)$$

An alternative bound based on a Euclidean norm is given by Munos [Mun03]. Experimental evidence indicates that the bound (3.2), while tight in theory ([BeT96], Example 6.4), is usually conservative in practice. Furthermore, often just a few policy evaluations are needed before the bound is attained. Some insight about the looseness of the bound (3.2) may be obtained by noting that we may add any multiple of the unit vector to Φr_μ , thereby making δ arbitrarily large, without affecting the policy improvement process and hence also the generated sequence $\{\mu^k\}$.

When the policy sequence $\{\mu^k\}$ terminates with some $\bar{\mu}$, the much sharper bound

$$\|J_{\bar{\mu}} - J^*\|_\infty \leq \frac{2\alpha\delta}{1-\alpha} \quad (3.3)$$

holds. To show this, let \bar{J} be the cost vector $\Phi r_{\bar{\mu}}$ obtained by policy evaluation of $\bar{\mu}$, and note that it satisfies $\|\bar{J} - J_{\bar{\mu}}\|_\infty \leq \delta$ (by our assumption on the accuracy of policy evaluation) and $T\bar{J} = T_{\bar{\mu}}\bar{J}$ (since μ^k terminates at $\bar{\mu}$). We write

$$TJ_{\bar{\mu}} \geq T(\bar{J} - \delta e) = T\bar{J} - \alpha\delta e = T_{\bar{\mu}}\bar{J} - \alpha\delta e \geq T_{\bar{\mu}}(J_{\bar{\mu}} - \delta e) - \alpha\delta e = T_{\bar{\mu}}J_{\bar{\mu}} - 2\alpha\delta e = J_{\bar{\mu}} - 2\alpha\delta e,$$

where e is the unit vector, from which by applying T to both sides, we obtain

$$T^2J_{\bar{\mu}} \geq TJ_{\bar{\mu}} - 2\alpha^2\delta e \geq J_{\bar{\mu}} - 2\alpha\delta(1+\alpha)e,$$

and by similar continued application of T to both sides,

$$J^* = \lim_{m \rightarrow \infty} T^m J_{\bar{\mu}} \geq J_{\bar{\mu}} - \frac{2\alpha\delta}{1-\alpha}e,$$

thereby showing the error bound (3.3).

An important fact, to be shown later, is that when the aggregation approach is used for policy evaluation, the policy sequence $\{\mu^k\}$ terminates with some $\bar{\mu}$, while this is generally not true for the projected equation approach.

3.1 The Issue of Policy Computations

We first note briefly an important practical issue: the computation of the improved policy via the minimization of Eq. (3.1) may present a serious difficulty when the expected value in Eq. (3.1) is hard to compute, and/or when the number of controls is large (or infinite in an extended version of the method where the control constraint set is infinite). Moreover, this issue arises generically when an approximate cost function is translated to an on-line implementation of a corresponding policy.

A number of devices have been proposed to address this difficulty, such as various forms of Q-learning, model simplifications (such as post-decision states), and various forms of approximate minimization in Eq. (3.1) (such as adaptive sampling). We refer to the literature, and particularly the books [CFH07], [Pow07], and [Ber10a] for an explanation of these terms and further discussion. Note that in some cases, the choice of the approximation architecture may be specially designed to exploit special problem structure and facilitate the minimization.

3.2 The Issue of Exploration

We next discuss an important generic difficulty with simulation-based policy iteration: to evaluate a policy μ , we need to generate cost samples using that policy, but this biases the simulation by underrepresenting states that are unlikely to occur under μ . As a result, the cost-to-go estimates of these underrepresented states may be highly inaccurate, causing potentially serious errors in the calculation of the improved control policy $\bar{\mu}$ via the policy improvement Eq. (3.1). A related issue is the presence of transient states in the Markov chain used for simulation (in which case the components of ξ corresponding to transient states are 0), or multiple recurrent classes (in which case states from some recurrent classes may never be generated during the simulation).

The difficulty just described is known as *inadequate exploration* of the system’s dynamics because of the use of a fixed policy. It is a particularly acute difficulty when the system is deterministic, or more generally when the randomness embodied in the transition probabilities is “relatively small.” One possibility for guaranteeing adequate exploration of the state space is to frequently restart the simulation and to ensure that the initial states employed form a rich and representative subset. There are also other related possibilities. For example, an approach called *iterative resampling*, is to enrich the sampled set of states in evaluating the current policy μ as follows: derive an initial cost evaluation of μ , simulate the next policy $\bar{\mu}$ obtained on the basis of this initial evaluation to obtain a set of representative states \bar{S} visited by $\bar{\mu}$, and repeat the evaluation of μ using additional trajectories initiated from \bar{S} .

One of the most frequently used approaches is to artificially introduce some extra randomization in the simulation, by occasionally modifying the state sequence generated by the policy μ . In a common scheme of this type we do row sampling according to an irreducible transition probability matrix

$$Q = (I - B)P + B\bar{P}, \quad (3.4)$$

where P is the transition matrix corresponding to μ , B is a diagonal matrix with diagonal components $\beta_i \in [0, 1]$, and \bar{P} is another transition matrix. Thus, at state i , the next state is generated with probability $1 - \beta_i$ according to transition probabilities p_{ij} corresponding to μ , and with probability β_i according to transition probabilities \bar{p}_{ij} whose purpose is to enhance exploration.†

Column sampling can be done using $L = P$, in which case the policy evaluation phase finds a solution to the projected equation $\Phi r = \Pi T(\Phi r)$, but with an important difference: the projection is with respect to the norm $\|\cdot\|_\xi$ where ξ is the steady-state distribution of Q rather than P . As a result of this norm mismatch, ΠT may not be a contraction, and a policy evaluation method that does not rely on contractions is needed.† LSTD is one such method, and the regression-based scaled LSPE method (2.31)-(2.32) is another.

† In the literature, e.g., [SuB98], the policy being evaluated is sometimes called the *target policy* to distinguish from the matrix \bar{P} that may correspond to a policy used for exploration, which is called *behavior policy*. Also, methods that use a behavior policy are called *off-policy* methods, while methods that do not are called *on-policy* methods. Note, however, that \bar{P} need not correspond to an admissible policy, and indeed there may not exist a suitable admissible policy that can induce sufficient exploration. For example in many queueing applications the transition matrices corresponding to all policies (including randomized ones) are very sparse and cannot induce adequate exploration.

† Proposition 2.2(a) applies only for small values of the probabilities β_i along the diagonal of B . It also applies only for λ sufficiently close to 1, in the case of the λ -methods of Section 2.3. The reason is that the operator $T^{(\lambda)}$ of Eq. (2.22) [and hence also $\Pi T^{(\lambda)}$] can be shown to be a contraction of modulus arbitrarily close to 0 if λ is sufficiently close to 1 (see e.g., [Ber07], Prop. 6.3.3, or [Ber10]).

There is a further restriction when a λ -method is used in conjunction with simulation. As discussed in Section 2.3, one must then choose the same Markov chain for row and column sampling ($Q = L$) in order to use the algorithm (2.23)-(2.25) for computing simulation-based $C_k^{(\lambda)} \approx C^{(\lambda)}$ and $d_k^{(\lambda)} \approx d^{(\lambda)}$. Still in this case both LSTD and the regression-based scaled LSPE method (2.31)-(2.32) may be used, even if $\Pi T^{(\lambda)}$ is not known to be a contraction (we refer to [BeY09] and [Yu10ab] for an analysis and discussion). It is then necessary, however, to know explicitly the ratios $a_{ij}/\ell_{ij} = \alpha p_{ij}/q_{ij}$ that appear in the algorithm (2.23)-(2.25). In practical applications reported in the literature, this difficulty has been bypassed by disregarding the ratios p_{ij}/q_{ij} from the formulas (2.23)-(2.25) (i.e., replacing a_{ij}/ℓ_{ij} with just α).

We mention that there are also some other possibilities for exploration enhancement for which we refer to [BeY09] and [Ber10a]. Moreover, the TD(0) method has been modified to address the exploration-related convergence difficulties (Sutton, Szepesvari, and Maei [SSM08], Maei et. al. [MSB08], Sutton et. al. [SMP09]). We finally note that exploration becomes an even more acute difficulty in the context of policy iteration for computing Q-factors, as we will discuss in Section 3.8.

3.3 Limited Sampling/Optimistic Versions

In the LSTD and LSPE methods discussed so far, the underlying assumption is that each policy is evaluated with a very large number of samples, so that accurate approximations of C and d are obtained. There are also optimistic versions, where the policy μ is replaced by an “improved” policy $\bar{\mu}$ after only a limited number of simulation samples have been processed. The motivation for such versions is potentially faster progress towards better policies, and less susceptibility to exploration difficulties (one may argue that mixing samples from several past policies may have a beneficial exploration effect). While experimental evidence lends some support to this speculation, solid evidence for the validity of optimistic methods is currently lacking, at least for the case of projected equation-based policy evaluation, when the chattering phenomenon may occur, as will be explained shortly. The validity of optimistic policy iteration is better established in the case of a lookup table representation (where $\Phi = I$; see [BeT96], Section 5.4) and the related case of aggregation-based policy evaluation (see Section 4). However, even in that case, the convergence properties of optimistic policy iteration are fragile and depend on the implementation details, as discussed in [BeT96], Section 5.4; see also Williams and Baird [WiB93] who provide remarkable counterexamples to the convergence of a version of asynchronous optimistic policy iteration where costs and policies are updated one at a time (these examples do not involve cost function approximation or simulation).

A natural form of optimistic LSTD is $\hat{r}_{k+1} = C_k^{-1}d_k$, where C_k and d_k are obtained by averaging samples collected using the controls corresponding to the (approximately) improved policy. By this we mean that C_k and d_k are time averages of the matrices and vectors

$$\phi(i_t)(\phi(i_t) - \alpha\phi(i_{t+1}))', \quad \phi(i_t)g(i_t, i_{t+1}),$$

corresponding to simulated transitions (i_t, i_{t+1}) that are generated using the policy μ^{k+1} whose controls are given by

$$\mu^{k+1}(i) \in \arg \min_{u \in U(i)} \sum_{j=1}^n p_{ij}(u)(g(i, u, j) + \alpha\phi(j)'\hat{r}_k).$$

Unfortunately, this method requires the collection of many samples between policy updates, as it is susceptible to simulation noise in C_k and d_k , particularly when C_k is nearly singular. The optimistic version of the (scaled) LSPE method

$$r_{k+1} = r_k - \gamma G_k(C_k r_k - d_k)$$

is based on similar ideas, but seems less susceptible to simulation noise because of its iterative nature. Note that because the theoretical convergence guarantees of LSPE apply only to the nonoptimistic version, it may be essential to experiment with various values of the stepsize γ .

To improve the reliability of the optimistic LSTD method it seems necessary to turn it into an iterative method, which then brings it very close to LSPE. In particular, an iterative version of the regression-based LSTD method (2.30) is given by Eq. (2.31), and is the special case of LSPE, corresponding to the special choice of the scaling matrix G_k of Eq. (2.32). This method, as well as other related versions of LSPE are less susceptible to the high level of simulation noise associated with limited number of samples between policy updates in optimistic policy iteration.

3.4 Policy Oscillation and Chattering

Projected equation-based variants of policy iteration methods are popular in practice, and have been tested extensively, dating to the early nineties (see e.g., the books [BeT96], [SuB08], and the references quoted there; for a sample of more recent experimental studies, see Lagoudakis and Parr [LaP03], Jung and Polani [JuP07], Busoniu et al. [BED09], and Thiery and Scherrer [ThS10]). However, the associated convergence behavior is complex, and involves potentially damaging oscillations, which despite their initial description long ago [Ber96], are still not well understood at present.

To get a sense of this behavior, we introduce the so called *greedy partition*, discussed in Section 6.4 of [BeT96]. This is a partition of the space \mathfrak{R}^s of parameter vectors r into subsets R_μ , each subset corresponding to a stationary policy μ , and defined by

$$R_\mu = \{r \mid T_\mu(\Phi r) = T(\Phi r)\} \quad (3.5)$$

or equivalently

$$R_\mu = \left\{ r \mid \mu(i) \in \arg \min_{u \in U(i)} \sum_{j=1}^n p_{ij}(u) (g(i, u, j) + \alpha \phi(j)'r), i = 1, \dots, n \right\}.$$

Thus, R_μ is the set of parameter vectors r for which μ is “greedy” with respect to Φr . Note that the greedy partition depends only on Φ , and is independent of the particular method used for policy evaluation (i.e., whether a projection/TD or different method is used).

For simplicity, let us assume that we use a policy evaluation method that for each given μ produces a unique parameter vector r_μ . Nonoptimistic policy iteration starts with a parameter vector r_0 , which specifies μ^0 as a greedy policy with respect to Φr_0 , and generates r_{μ^0} by using the given policy evaluation method. It then finds a policy μ^1 that is greedy with respect to Φr_{μ^0} , i.e., a μ^1 such that $r_{\mu^0} \in R_{\mu^1}$. It then repeats the process with μ^1 replacing μ^0 . If some policy μ^k satisfying $r_{\mu^k} \in R_{\mu^k}$ is encountered, the method keeps generating that policy. This is the necessary and sufficient condition for policy convergence in the nonoptimistic policy iteration method. Of course, the mere fact that a policy iteration method is guaranteed to converge is not in itself a guarantee of good performance when cost function approximation is used, beyond the fact that a better error bound holds in this case [Eq. (3.3) versus Eq. (3.2)].

In the case of a lookup table representation where the parameter vectors r_μ are equal to the cost-to-go vector J_μ , the condition $r_{\mu^k} \in R_{\mu^k}$ is equivalent to $r_{\mu^k} = T r_{\mu^k}$, and is satisfied if and only if μ^k is optimal. When there is function approximation, however, this condition need not be satisfied for any policy.

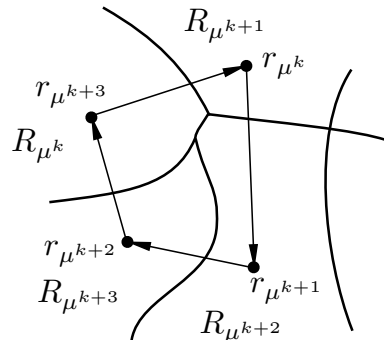


Figure 3.1. Greedy partition and cycle of policies generated by nonoptimistic policy iteration. In particular, μ yields $\bar{\mu}$ by policy improvement if and only if $r_\mu \in R_{\bar{\mu}}$. In this figure, the method cycles between four policies and the corresponding four parameters $r_{\mu^k}, r_{\mu^{k+1}}, r_{\mu^{k+2}}, r_{\mu^{k+3}}$.

Since there is a finite number of possible vectors r_μ , one generated from another in a deterministic way, the algorithm ends up repeating some cycle of policies $\mu^k, \mu^{k+1}, \dots, \mu^{k+m}$ with

$$r_{\mu^k} \in R_{\mu^{k+1}}, r_{\mu^{k+1}} \in R_{\mu^{k+2}}, \dots, r_{\mu^{k+m-1}} \in R_{\mu^{k+m}}, r_{\mu^{k+m}} \in R_{\mu^k}; \quad (3.6)$$

(see Fig. 3.1). Furthermore, there may be several different cycles, and the method may end up converging to any one of them. The actual cycle obtained depends on the initial policy μ^0 . This is similar to gradient methods applied to minimization of functions with multiple local minima, where the limit of convergence depends on the starting point.

In the case of optimistic policy iteration, the trajectory of the method is less predictable and depends on the fine details of the iterative policy evaluation method, such as the frequency of the policy updates and the stepsize used. Generally, given the current policy μ , optimistic policy iteration will move towards the corresponding “target” parameter r_μ , for as long as μ continues to be greedy with respect to the current cost-to-go approximation Φr , that is, for as long as the current parameter vector r belongs to the set R_μ . Once, however, the parameter r crosses into another set, say $R_{\bar{\mu}}$, the policy $\bar{\mu}$ becomes greedy, and r changes course and starts moving towards the new “target” $r_{\bar{\mu}}$. Thus, the “targets” r_μ of the method, and the corresponding policies μ and sets R_μ may keep changing, similar to nonoptimistic policy iteration. Simultaneously, the parameter vector r will move near the boundaries that separate the regions R_μ that the method visits, following reduced versions of the cycles that nonoptimistic policy iteration may follow. Furthermore, as Fig. 3.1 suggests, if diminishing parameter changes are made between policy updates (such as for example when a diminishing stepsize is used by the policy evaluation method) and the method eventually cycles between several policies, the parameter vectors will tend to converge to the common boundary of the regions R_μ corresponding to these policies. This is the so-called *chattering* phenomenon for optimistic policy iteration, whereby there is simultaneously oscillation in policy space and convergence in parameter space. The following is a simple example of policy oscillations and chattering. Other examples are given in Section 6.4.2 of [BeT96] (Examples 6.9 and 6.10).

Example 3.1: (Policy Oscillation and Chattering)

Consider a discounted problem with two states, 1 and 2, illustrated in Fig. 3.2(a). There is a choice of control only at state 1, and there are two policies, denoted μ^* and μ . The optimal policy μ^* , when at state 1, stays at 1 with probability $p > 0$ and incurs a cost $c < 0$, while when at state 2, returns to state 1 with 0 cost. The

other policy is μ and cycles between the two states with 0 cost. We consider linear approximation with a single feature $\phi(i)' = i$ for each of the states $i = 1, 2$, i.e.,

$$\Phi = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \quad \tilde{J} = \Phi r = \begin{pmatrix} r \\ 2r \end{pmatrix}.$$

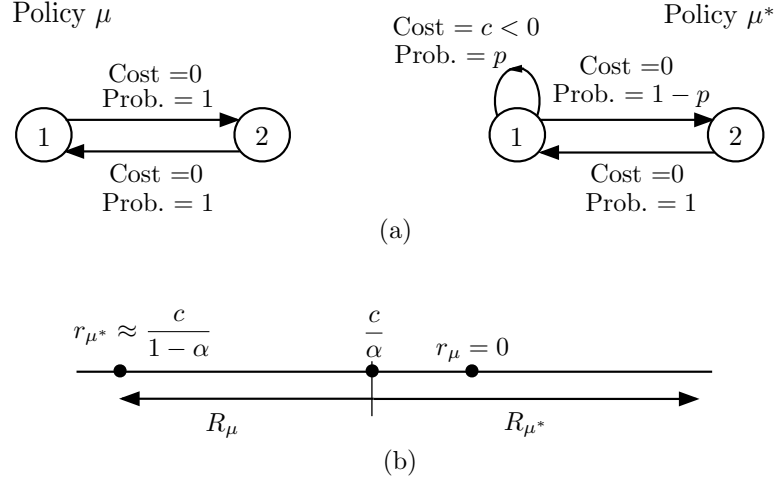


Figure 3.2. The problem of Example 3.1. (a) Costs and transition probabilities for the policies μ and μ^* . (b) The greedy partition and the solutions of the projected equations corresponding to μ and μ^* . Nonoptimistic policy iteration oscillates between r_μ and r_{μ^*} .

We first derive the greedy partition using Eq. (3.5):

$$R_{\mu^*} = \left\{ r \mid p(c + \alpha(1 \cdot r)) + (1-p)\alpha(2 \cdot r) \leq \alpha(2 \cdot r) \right\} = \{r \mid c \leq \alpha r\}, \quad R_\mu = \{r \mid c \geq \alpha r\}.$$

We next calculate the points r_μ and r_{μ^*} that solve the projected equations

$$C_\mu r_\mu = d_\mu, \quad C_{\mu^*} r_{\mu^*} = d_{\mu^*},$$

which correspond to μ and μ^* , respectively. We have

$$C_\mu = \Phi' \Xi_\mu (1 - \alpha P_\mu) \Phi = (1 \ 2) \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & -\alpha \\ -\alpha & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = 5 - 9\alpha,$$

$$d_\mu = \Phi' \Xi_\mu g_\mu = (1 \ 2) \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \end{pmatrix} = 0,$$

so

$$r_\mu = 0.$$

Similarly, with some calculation,

$$C_{\mu^*} = \Phi' \Xi_{\mu^*} (1 - \alpha P_{\mu^*}) \Phi = (1 \ 2) \begin{pmatrix} \frac{1}{2-p} & 0 \\ 0 & \frac{1-p}{2-p} \end{pmatrix} \begin{pmatrix} 1-\alpha p & -\alpha(1-p) \\ -\alpha & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \frac{5 - 4p - \alpha(4 - 3p)}{2-p},$$

$$d_{\mu^*} = \Phi' \Xi_{\mu^*} g_{\mu^*} = (1 \ 2) \begin{pmatrix} \frac{1}{2-p} & 0 \\ 0 & \frac{1-p}{2-p} \end{pmatrix} \begin{pmatrix} c \\ 0 \end{pmatrix} = \frac{c}{2-p},$$

so

$$r_{\mu^*} = \frac{c}{5 - 4p - \alpha(4 - 3p)}.$$

We now note that since $c < 0$, $r_\mu = 0 \in R_{\mu^*}$, while for $p \approx 1$ and $\alpha > 1 - \alpha$, we have $r_{\mu^*} \approx \frac{c}{1-\alpha} \in R_\mu$; cf. Fig. 3.2(b). In this case, policy iteration cycles between μ and μ^* . Its optimistic version uses some algorithm that moves the current value r towards r_{μ^*} if $r \in R_{\mu^*}$, and towards r_μ if $r \in R_\mu$. Thus optimistic policy iteration starting from a point in R_μ moves towards r_{μ^*} and once it crosses the boundary point c/α of the greedy partition, it reverses course and moves towards r_μ . If the method makes small incremental changes in r before checking whether to change the current policy, it will incur a small oscillation around c/α . If the incremental changes in r are diminishing, the method will converge to c/α . Yet c/α does not correspond to any one of the two policies and has no meaning as a desirable parameter value.

Notice that it is hard to predict when an oscillation will occur and what kind of oscillation it will be. For example if $c > 0$, we have $r_\mu = 0 \in R_\mu$, while for $p \approx 1$ and $\alpha > 1 - \alpha$, we have $r_{\mu^*} \approx \frac{c}{1-\alpha} \in R_{\mu^*}$. Then, approximate and optimistic policy iteration converge to μ (or μ^*) if started with r in R_μ (or R_{μ^*} , respectively).

When chattering occurs, the limit of optimistic policy iteration tends to be on a common boundary of several subsets of the greedy partition and *may not meaningfully represent a cost approximation of any of the corresponding policies*, as illustrated by the preceding example. Thus, the limit to which the method converges cannot always be used to construct an approximation of the cost-to-go of any policy or the optimal cost-to-go. As a result, at the end of optimistic policy iteration and in contrast with the nonoptimistic version, one must go back and perform a screening process; that is, evaluate by simulation the many policies generated by the method starting from the initial conditions of interest and select the most promising one. This is a disadvantage of optimistic policy iteration that may nullify whatever practical rate of convergence advantages the method may have over its nonoptimistic counterpart.

An additional insight is that the choice of the iterative policy evaluation method (e.g., LSTD, LSPE, or TD for various values of λ) makes a difference in rate of convergence, but does not seem crucial for the quality of the final policy obtained (as long as the methods converge). The choice of method affects the targets r_μ somewhat, but leaves the greedy partition unchanged. As a result, different methods “fish in the same waters” and tend to yield similar ultimate cycles of policies. Generally, there is no guarantee that any of the policies involved in oscillation is particularly good; the policy iteration process may be just cycling in a “bad” part of the greedy partition. Moreover, it is hard to know how far from the optimum the oscillation is occurring, similar to the case of local minima in gradient-based optimization.

The full ramifications of policy oscillation in practice are not fully understood at present, but it is clear that they give serious reason for concern. An interesting case in point is the game of tetris, which has been used as a testbed for approximate DP methods [Van95], [TsV96], [BeI96], [Kak02], [FaV06], [SzL06], [DFM09].[†] Using policy iteration with policy evaluation based on the projected equation and a set of 22 features, an average score of a few thousands was achieved (originally the LSPE method was used [BeI96]; similar results were obtained with the LSTD method [LaP03]). Using the same features and a random search method (the cross-entropy method [BKM05], [RuK08]) in the space of weight vectors r , an average score

[†] Local attraction-type phenomena may be causing similar difficulties in other related approximate DP methodologies: policy iteration with function approximation using the Bellman error method, policy gradient methods, and approximate linear programming (the tetris problem, using the same 22 features, has been addressed by approximate linear programming [FaV06], [DFM09], and with a policy gradient method [Kak02], also with an achieved average score of a few thousands, roughly comparable to the ones obtained using policy iteration). These connections are interesting subjects for investigation.

of 600,000-900,000 was achieved [SzL06], [ThS09]. The causes of this notable failure of approximate policy iteration have not been clarified yet. It is an open question whether it is due to the inherent difficulty of the tetris problem, or whether it can be attributed to inherent pathologies of approximate policy iteration, such as oscillations/chattering between relatively poor policies, or weaknesses of the cost function approximation philosophy.‡

3.5 Constrained Policy Iteration

It is natural in approximate DP to try to exploit whatever prior information is available about J^* , such as for example bounds on $J^*(i)$ for some or all of the states i . If it is known that J^* belongs to a subset of \mathfrak{R}^n , we may try to find an approximation Φr that belongs to that subset. This leads to projected equations of the form $\Phi r = \Pi T(\Phi r)$, where Π is the operator that projects a vector onto a closed convex subset \hat{S} of the subspace $S = \{\Phi r \mid r \in \mathfrak{R}^s\}$, a form of nonlinear approximation. Corresponding analogs of LSTD and LSPE-type methods for such projected equations involve the solution of linear variational inequalities rather than linear systems of equations, and are described in [Ber09a]. In particular, it is shown in [Ber09a] that \hat{r} solves the equation $\Phi \hat{r} = \Pi T(\Phi \hat{r})$ if and only if it solves the variational inequality

$$(C\hat{r} - d)'(r - \hat{r}) \geq 0, \quad \forall r \in \hat{R}, \quad (3.7)$$

where $\hat{R} = \{r \mid \Phi r \in \hat{S}\}$ [cf. Eq. (2.1)]. This is a generalization of the equation $Cr = d$, which corresponds to the case $\hat{R} = \mathfrak{R}^s$. Thus one may carry out simulation-based policy evaluation by replacing C and d with simulation-based estimates in the preceding variational inequality. A potential difficulty here is that even if the set \hat{S} is simple, the set \hat{R} may be complicated (e.g., it may be specified by a large number of inequalities), in which case one may consider approximating \hat{R} with a more convenient set.

There is a simpler alternative to solving the variational inequality (3.7) in the practically common case where upper and/or lower bounds of J^* are available. In this alternative, which stems from the work of Bertsekas and Yu [BeY10a], we introduce the bounds directly into the approximation architecture, and use a corresponding nonlinear iterative policy evaluation algorithm. As an example, suppose that we have access to the components of a vector \bar{J} with $\bar{J}(i) \geq J^*(i)$ for all i (for example, the cost vector of a known policy).

‡ For a given policy μ , the policy evaluation methods based on the projected equation estimate the cost function of μ by Φr_μ , yet the next policy $\bar{\mu}$ obtained by policy improvement is determined by the Q-factor differences $\tilde{Q}_\mu(i, u) - \tilde{Q}_\mu(i, \mu(i))$, where

$$\tilde{Q}_\mu(i, u) = \sum_{j=1}^n p_{ij}(u) (g(i, u, j) + \alpha \phi(j)' r_\mu), \quad i = 1, \dots, n, \quad u \in U(i),$$

is the corresponding approximation of the Q-factor $Q_\mu(i, u)$ of μ [cf. Eq. (3.1)]. Thus adding a constant to all components of Φr_μ does not affect $\bar{\mu}$. This has led to proposals of policy evaluation methods that aim to estimate directly Q-factor differences or cost function differences (see Werbos [Wer92a], [Wer92b], Baird, Harmon, and Klopf [Bai93], [HBK94], and Bertsekas [Ber97]). There is no solid theory regarding the merits of these methods over methods based on cost or Q-factor approximation, but it is clear that there are circumstances where they may offer an advantage (see the discussion of [BeT96], Sections 6.6.2 and 6.11). The problem is that policy evaluation based on least squares criteria tends to capture the large scale behavior of the function being approximated, and may thus be less sensitive to the fine scale variations in $Q_\mu(i, u)$. As a result, the dependence of Q_μ on u may be lost, when using instead the approximation \tilde{Q}_μ .

Then policy iteration with cost function approximation can be modified to incorporate this knowledge as follows. Given a policy μ , we evaluate it by finding an approximation Φr_μ to the solution \tilde{J}_μ of the equation

$$\tilde{J}_\mu(i) = \sum_{j=1}^n p_{ij}(\mu(i)) \left(g(i, \mu(i), j) + \alpha \min \{ \bar{J}(j), \tilde{J}_\mu(j) \} \right), \quad i = 1, \dots, n, \quad (3.8)$$

followed by the policy improvement operation [cf. Eq. (3.1)]

$$\bar{\mu}(i) \in \arg \min_{u \in U(i)} \sum_{j=1}^n p_{ij}(u) \left(g(i, u, j) + \alpha \min \{ \bar{J}(j), \phi(j)' r_\mu \} \right), \quad i = 1, \dots, n, \quad (3.9)$$

where $\phi(j)'$ is the row of Φ that corresponds to state j .

Note that Eq. (3.8) has a unique solution since its right-hand side represents a sup-norm contraction mapping. Indeed this is Bellman's equation for the Q-factor of an *optimal stopping problem* involving a stopping cost $\bar{J}(i)$ at state i , a problem that has received a lot of attention in the approximate DP literature because of its favorable structure (see [BeT96], Section 6.8, or [Ber07], Section 6.4.2). In particular, Eq. (3.8) can be solved approximately, with approximation of \tilde{J}_μ by Φr_μ , using variants of the TD(0) and LSPE(0) methods (see Tsitsiklis and Van Roy [TSV99b], Choi and Van Roy [ChV06], Yu and Bertsekas [YuB07]). One form of the LSPE(0) method of [YuB07] generates a long trajectory (i_0, i_1, \dots, i_m) using the policy $\bar{\mu}$ of Eq. (3.9), and determines the weight vector $r_{\bar{\mu}}$ as the limit of the iteration

$$r_{k+1} = \arg \min_{r \in \mathbb{R}^s} \sum_{t=0}^{m-1} \left(\phi(i_t)' r - g(i_t, \mu(i_t), i_{t+1}) - \alpha \min \{ \bar{J}(i_{t+1}), \phi(i_{t+1})' r_k \} \right)^2, \quad r_0 = r_\mu, \quad (3.10)$$

i.e., $r_{\bar{\mu}} = \lim_{k \rightarrow \infty} r_k$ (for this it is necessary that the length of the simulation trajectory is sufficiently large). This is the analog of the batch simulation approach of Eq. (2.6): first a batch of simulation data is generated, and then the corresponding Bellman equation (3.8) is solved approximately via the iteration (3.10) [an LSTD/matrix inversion approach is not possible because Eq. (3.8) is nonlinear]. The method (3.10) involves substantial overhead because it involves the entire simulation trajectory at each iteration. There are also alternative methods given by [YuB07], which are more iterative in character (iterate on r as new simulation data is generated), and also require less overhead per iteration; see also [Ber07], [Ber10a] for details.†

Under the assumption $\bar{J}(i) \geq J^*(i)$ for all i , and a lookup table representation ($\Phi = I$), it can be shown that the method (3.8)-(3.9) yields J^* in a finite number of iterations, just like the standard (exact) policy iteration method. The proof is based on the monotonicity and sup-norm contraction properties of the mappings \bar{T}_μ and \bar{T} given by

$$(\bar{T}_\mu J)(i) = \sum_{j=1}^n p_{ij}(\mu(i)) \left(g(i, \mu(i), j) + \alpha \min \{ \bar{J}(j), J(j) \} \right), \quad i = 1, \dots, n, \quad (3.11)$$

† An algorithm similar to (3.8)-(3.10) is also possible if we know instead a lower bounding vector to J^* , or both upper and lower bounding vectors: the terms $\min \{ \bar{J}(j), \tilde{J}_\mu(j) \}$ and $\min \{ \bar{J}(j), \phi(j)' r_\mu \}$ in Eqs. (3.8)-(3.9) should be replaced by the respective projections of $\tilde{J}_\mu(j)$ or $\phi(j)' r_\mu$ on the corresponding interval of upper and/or lower bounds. In this case a generalization of the optimal stopping algorithm of [YuB07], given in [BeY09], Section 7.3, may be used for policy evaluation. Note that this algorithm provides a means for approximating the cost $J_\mu(j)$ of a policy μ with a simple nonlinear architecture $f_j(\phi(j)' r)$, where f_j denotes projection onto the interval of upper and lower bounds to $J^*(j)$. This approach may make sense regardless of whether true upper and/or lower bounds of $J^*(j)$ are known.

$$(\bar{T}J)(i) = \min_{u \in U(i)} \sum_{j=1}^n p_{ij}(u) \left(g(i, u, j) + \alpha \min \{ \bar{J}(j), J(j) \} \right), \quad i = 1, \dots, n, \quad (3.12)$$

[cf. Eqs. (3.8) and (3.9)], and the fact that J^* is the unique fixed point of F (which rests on the hypothesis $\bar{J} \geq J^*$); see the subsequent Prop. 3.1.

When a compact feature-based representation is used ($\Phi \neq I$), the constrained policy iteration method (3.8)-(3.9) is still vulnerable to oscillation and chattering, regardless of which variant of Eq. (3.10) is used. The magnitude of the oscillations can be estimated with an error bound similar to Eq. (3.2). However, the method may produce substantially different results than its unconstrained counterpart, because it operates with a different greedy partition, so the policies involved in an oscillation may be different in the two methods.

3.6 Conditions for Policy Convergence

The preceding analysis has indicated that it is desirable to avoid policy oscillation in policy iteration methods. Moreover, as mentioned earlier, when policies converge there is a more favorable error bound associated with the method [cf. Eq. (3.3) versus Eq. (3.2)]. It is therefore interesting to investigate conditions under which we have convergence of policies. From the mathematical point of view, it turns out that policy oscillation is caused by the lack of monotonicity and the dependence (through Ξ) on μ of the projection operator. With this in mind, we will replace Π with a constant operator W that has a monotonicity property. We develop convergence conditions for an abstract approximate policy iteration algorithm that applied to general (beyond discounted) DP problems. This leads to conceptual and notational simplification as well as a more broadly applicable algorithm and analysis.

To this end, consider any (possibly nonlinear) mapping $H_\mu : \mathfrak{R}^n \mapsto \mathfrak{R}^n$, parametrized by the policy μ , and the mapping $H : \mathfrak{R}^n \mapsto \mathfrak{R}^n$, defined by

$$HJ = \min_{\mu \in \mathcal{M}} H_\mu J, \quad (3.13)$$

where \mathcal{M} is a finite subset of policies, and the minimization above is done separately for each component of $H_\mu J$; i.e.,

$$(HJ)(i) = \min_{\mu \in \mathcal{M}} (H_\mu J)(i), \quad \forall i = 1, \dots, n.$$

Abstract mappings of this type and their relation to DP have been studied in Denardo [Den67], Bertsekas [Ber77], and Bertsekas and Shreve [BeS78]. A special case is the discounted DP case, where $H_\mu = T_\mu$ and $H = T$. Another special case is a mapping H_μ that is similar to T_μ but arises in discounted semi-Markov problems. A special case of a nonlinear mapping H_μ arises in constrained policy iteration, cf. the mapping \bar{T}_μ of Eq. (3.11). Nonlinear mappings H_μ also arise in the context of minimax DP problems and sequential games; see Shapley [Sha53], and [Den67], [Ber77], [BeS78].

We will construct a policy iteration method that aims to find an approximation to a fixed point of H , and evaluates a policy $\mu \in \mathcal{M}$ with a solution \tilde{J}_μ of the following fixed point equation in the vector J :

$$(WH_\mu)(J) = J, \quad (3.14)$$

where $W : \mathfrak{R}^n \mapsto \mathfrak{R}^n$ is a mapping (possibly nonlinear, but independent of μ). Policy evaluation by solving the projected equation corresponds to $W = \Pi$. Rather than specify properties of H_μ under which H has a unique fixed point (as in [Den67], [Ber77], and [BeS78]), it is simpler for our purposes to introduce corresponding assumptions on the mappings W and WH_μ . In particular, we assume the following:

(a) For each J , the minimum in Eq. (3.13) is attained, in the sense that there exists $\bar{\mu} \in \mathcal{M}$ such that $HJ = H_{\bar{\mu}}J$.

(b) For each $\mu \in \mathcal{M}$, the mappings W and WH_{μ} are monotone in the sense that

$$WJ \leq W\bar{J}, \quad (WH_{\mu})(J) \leq (WH_{\mu})(\bar{J}), \quad \forall J, \bar{J} \in \mathfrak{R}^n \text{ with } J \leq \bar{J}. \quad (3.15)$$

(c) For each μ , the solution \tilde{J}_{μ} of Eq. (3.14) is unique, and for all J such that $(WH_{\mu})(J) \leq J$, we have

$$\tilde{J}_{\mu} = \lim_{k \rightarrow \infty} (WH_{\mu})^k(J).$$

Based on condition (a), we introduce a policy improvement operation that is similar to the case where $H_{\mu} = T_{\mu}$: the “improved” policy $\bar{\mu}$ satisfies $H_{\bar{\mu}}\tilde{J}_{\mu} = H\tilde{J}_{\mu}$. Note that condition (c) is satisfied if WH_{μ} is a contraction on S , while condition (b) is satisfied if W is a matrix with nonnegative components and H_{μ} is monotone for all μ . These conditions also hold for the constrained policy iteration method (3.8)-(3.9), where H_{μ} is the nonlinear mapping \tilde{T}_{μ} of Eq. (3.11).

Proposition 3.1: Let conditions (a)-(c) hold. Consider the approximate policy iteration method that uses the fixed point \tilde{J}_{μ} of the mapping WH_{μ} for evaluation of the policy μ [cf. Eq. (3.14)], and the equation $H_{\bar{\mu}}\tilde{J}_{\mu} = H\tilde{J}_{\mu}$ for policy improvement. Assume that the method is initiated with some policy in \mathcal{M} , and it is operated so that it terminates when a policy $\bar{\mu}$ is obtained such that $H_{\bar{\mu}}\tilde{J}_{\bar{\mu}} = H\tilde{J}_{\bar{\mu}}$. Then the method terminates in a finite number of iterations, and the vector $\tilde{J}_{\bar{\mu}}$ obtained upon termination is a fixed point of WH .

Proof: Similar to the standard proof of convergence of (exact) policy iteration, we use the policy improvement equation $H_{\bar{\mu}}\tilde{J}_{\mu} = H\tilde{J}_{\mu}$, the monotonicity of W , and the policy evaluation Eq. (3.14) to write

$$(WH_{\bar{\mu}})(\tilde{J}_{\mu}) = (WH)(\tilde{J}_{\mu}) \leq (WH_{\mu})(\tilde{J}_{\mu}) = \tilde{J}_{\mu}.$$

By iterating with the monotone mapping $WH_{\bar{\mu}}$ and by using condition (c), we obtain

$$\tilde{J}_{\bar{\mu}} = \lim_{k \rightarrow \infty} (WH_{\bar{\mu}})^k(\tilde{J}_{\mu}) \leq \tilde{J}_{\mu}.$$

There are finitely many policies, so we must have $\tilde{J}_{\bar{\mu}} = \tilde{J}_{\mu}$ after a finite number of iterations, which using the policy improvement equation $H_{\bar{\mu}}\tilde{J}_{\mu} = H\tilde{J}_{\mu}$, implies that $H_{\bar{\mu}}\tilde{J}_{\bar{\mu}} = H\tilde{J}_{\bar{\mu}}$. Thus the algorithm terminates with $\bar{\mu}$, and since $\tilde{J}_{\bar{\mu}} = (WH_{\bar{\mu}})(\tilde{J}_{\bar{\mu}})$, it follows that $\tilde{J}_{\bar{\mu}}$ is a fixed point of WH . **Q.E.D.**

Note the mechanism by which the preceding policy iteration method avoids the oscillations described in Section 3.5. Again the policy improvement step yields $\bar{\mu}$ from μ if and only if $\tilde{J}_{\mu} \in R_{\bar{\mu}}$, where $R_{\bar{\mu}} = \{J \mid H_{\bar{\mu}}J = HJ\}$ is the corresponding set of the greedy partition (which is independent of the method used for policy evaluation). However, the monotonicity inherent in the policy evaluation step guarantees that

$$\tilde{J}_{\bar{\mu}} \leq \tilde{J}_{\mu}, \quad \forall \mu, \bar{\mu} \in \mathcal{M} \text{ such that } \tilde{J}_{\mu} \in R_{\bar{\mu}},$$

so a policy oscillation is impossible. Note also that the preceding proof applies to the more general case where for policy improvement we use any policy $\bar{\mu}$ such that $H_{\bar{\mu}}\tilde{J}_{\mu} \leq \tilde{J}_{\mu}$, with strict inequality $H_{\bar{\mu}}\tilde{J}_{\mu} \neq \tilde{J}_{\mu}$ when $H_{\bar{\mu}}\tilde{J}_{\mu} \neq \tilde{J}_{\mu}$. This allows for approximate minimization over candidate policies in the policy improvement step.

An important special case where Prop. 3.1 applies and policies converge is when $H_{\mu} = T_{\mu}$, $H = T$, W is linear of the form $W = \Phi D$, where Φ and D are $n \times s$ and $s \times n$ matrices, respectively, whose rows are probability distributions, and the policy evaluation uses the linear feature-based approximation $\tilde{J}_{\mu} = \Phi r_{\mu}$. This is the case of policy evaluation by aggregation [cf. Eq. (1.7)], which will be discussed in the next section. Then it can be seen that W and WT_{μ} are monotone, and that WT_{μ} is a sup-norm contraction (since W is nonexpansive with respect to the sup norm), so that conditions (a)-(c) are satisfied. The same is true in the more general case $W = \Phi D$, where the matrix W has nonnegative components, and its row sums are less than or equal to 1, i.e.,

$$\sum_{m=1}^s \Phi_{im} D_{mj} \geq 0, \quad \forall i, j = 1, \dots, n, \quad \sum_{m=1}^s \Phi_{im} \sum_{j=1}^n D_{mj} \leq 1, \quad \forall i = 1, \dots, n.$$

If Φ and D have nonnegative components, the first of these relations is automatically satisfied, while the second is equivalent to the set of n linear inequalities

$$\phi(i)' \zeta \leq 1, \quad \forall i = 1, \dots, n, \quad (3.16)$$

where $\phi(i)'$ is the i th row of Φ , and $\zeta \in \mathfrak{R}^s$ is the column vector of row sums of D , i.e., the one that has components[†]

$$\zeta(m) = \sum_{j=1}^n D_{mj}, \quad \forall m = 1, \dots, s.$$

Even in this more general case, the policy evaluation Eq. (3.14) can be solved by using simulation and low order calculations (see [Ber10a], Section 6.8).

Given a “good” choice of $\Phi \geq 0$, an interesting question is how to construct effective algorithms for parametric optimization of a nonnegative matrix D , subject to the constraints (3.16). This is similar to the basis function optimization issue mentioned in the introduction. Finally, let us note that a special case when Eq. (3.16) may be easily satisfied arises when through a reordering of indexes, the matrix D can be partitioned in the form $D = (\Delta \ 0)$, where Δ is a positive definite diagonal matrix with diagonal elements δ_m , $m = 1, \dots, s$, satisfying

$$\sum_{m=1}^s \Phi_{im} \delta_m \leq 1, \quad \forall i = 1, \dots, n.$$

An example of such a structure arises in coarse grid discretization/aggregation schemes (see Section 4).

[†] A column of Φ that has both positive and negative components may be replaced with the two columns that contain its positive and the opposite of its negative components. This will create a new nonnegative matrix Φ with as many as twice the number of columns, and will also enlarge the approximation subspace S (leading to no worse approximation). Then the matrix D may be optimized subject to $D \geq 0$ and the constraints (3.16), with respect to some performance criterion.

3.7 The Case of Hard Aggregation

If the mapping H_μ is monotone in the sense that $H_\mu J \leq H_\mu \bar{J}$ for all $J, \bar{J} \in \mathfrak{R}^n$ with $J \leq \bar{J}$ (as it typically is in DP models), then the monotonicity assumption (3.15) is satisfied if W is a projection matrix Π (the same for all μ), which satisfies $\Pi J \geq 0$ for all J with $J \geq 0$ (or equivalently $\Pi \mathfrak{R}_+^n = S \cap \mathfrak{R}_+^n$, where \mathfrak{R}_+^n is the nonnegative orthant). This is so if and only if all the components of Π are nonnegative. A special case where this is true is when Φ has nonnegative components, and has linearly independent columns that are orthogonal with respect to the inner product $\langle x_1, x_2 \rangle = x_1' \Xi x_2$. This follows from the projection formula $\Pi = \Phi(\Phi' \Xi \Phi)^{-1} \Phi' \Xi$ and the fact that $\Phi' \Xi \Phi$ is positive definite and diagonal. An interesting special case is when policy evaluation is done by *hard aggregation*, one of the methods from the class to be discussed in the next section. Here the state space $\{1, \dots, n\}$ is partitioned in s nonempty subsets I_1, \dots, I_s and:

- (1) The ℓ th column of Φ has components that are 1 or 0 depending on whether they correspond to an index in I_ℓ or not.
- (2) The ℓ th row of D is a probability distribution $(d_{\ell 1}, \dots, d_{\ell n})$ whose components are positive depending on whether they correspond to an index in I_ℓ or not, i.e., $\sum_{j=1}^n d_{\ell j} = 1$, $d_{\ell j} > 0$ if $j \in I_\ell$, and $d_{\ell j} = 0$ if $j \notin I_\ell$.

With these definitions of Φ and D , it can be verified that ΦD is given by the projection formula

$$\Phi D = \Phi(\Phi' \Xi \Phi)^{-1} \Phi' \Xi,$$

where Ξ is the diagonal matrix with the nonzero components of D along the diagonal. In fact Π , the projection on S with respect to $\|\cdot\|_\xi$, can be written in the explicit form

$$(\Pi J)(i) = \sum_{j \in I_\ell} d_{\ell j} J(j), \quad \forall i \in I_\ell, \ell = 1, \dots, s.$$

Thus, assuming that D (and hence Π) is held constant, policies converge and the more favorable error bound (3.3) holds.†

3.8 Q-Factor Approximations, Exploration, and Optimistic Policy Iteration

The policy iteration methods with cost function approximation discussed so far rely on the calculation of the approximation Φr to the cost function of the current policy μ , which is then used for policy improvement using the minimization (3.1). Carrying out this minimization requires knowledge of the transition probabilities $p_{ij}(u)$ and calculation of the associated expected values for all controls $u \in U(i)$ (otherwise a time-consuming simulation of these expected values is needed). An interesting “model-free” alternative is to perform approximate policy evaluation by computing *approximate Q-factors*

$$\tilde{Q}(i, u, r) \approx (F_\mu Q)(i, u), \quad \forall (i, u), \quad (3.17)$$

† Van Roy [Van06] has established especially favorable error bounds for a hard aggregation-based approximate DP method. However, the policy iteration method considered in [Van06] is not the same as the one considered here, because the corresponding matrix Π is not constant and policy convergence is not guaranteed. It is not clear whether it is practically advantageous to select Π in the manner suggested by Van Roy [Van06], compared to keeping it constant as in the standard hard aggregation scheme of the next section.

where $F_\mu Q$ is the Q-learning mapping given by

$$(F_\mu Q)(i, u) = \sum_{j=1}^n p_{ij}(u) \left(g(i, u, j) + \alpha Q(j, \mu(j)) \right), \quad \forall (i, u), \quad (3.18)$$

and to use the minimization

$$\bar{\mu}(i) \in \arg \min_{u \in U(i)} \tilde{Q}(i, u, r), \quad \forall i, \quad (3.19)$$

for policy improvement. Here, r is an adjustable parameter vector and $\tilde{Q}(i, u, r)$ is a parametric architecture, possibly of the linear form

$$\tilde{Q}(i, u, r) = \sum_{k=1}^s r_k \phi_k(i, u),$$

where $\phi_k(i, u)$ are basis functions that depend on both state and control.

An important point is that given the current policy μ , we can construct Q-factor approximations for that policy using any method for constructing cost vector approximations. The way to do this is to apply the latter method to the Markov chain whose states are the pairs (i, u) , and the probability of transition from (i, u) to (j, v) is

$$p_{(i,u)(j,v)} = p_{ij}(u) \delta(v = \mu(j)), \quad (3.20)$$

where as earlier $\delta(\cdot)$ denotes the indicator function: $\delta(v = \mu(j)) = 1$ if $v = \mu(j)$ and is equal to 0 otherwise. This is the probabilistic mechanism by which state-control pairs evolve under the stationary policy μ . However, a major concern here is that the state-control pairs (i, u) with $u \neq \mu(i)$ are never generated in this Markov chain, so they are not represented in the cost samples used to construct the approximation $\tilde{Q}(i, u, r)$. This creates an acute difficulty due to diminished exploration, which must be addressed in any simulation-based implementation.

To address the issue of exploration in evaluating (exactly or approximately) the Q-factors of a given policy, we may use the LSTD(λ) approach in conjunction with the algorithm (2.23)-(2.25) of Section 2.3. We recall that in this algorithm, row and column sampling is done by generating a sequence of state-control pairs $\{(i_0, u_0), (i_1, u_1) \dots\}$. The Q-factor approximation is $\Phi \hat{r}$, where \hat{r} is the solution of the system $C_k^{(\lambda)} r = d_k^{(\lambda)}$, with $C_k^{(\lambda)}$ and $d_k^{(\lambda)}$ generated recursively by

$$C_k^{(\lambda)} = \left(1 - \frac{1}{k+1} \right) C_{k-1}^{(\lambda)} + \frac{1}{k+1} z_k \left(\phi(i_k, u_k) - \alpha \frac{p_{(i_k, u_k)(i_{k+1}, u_{k+1})}}{\ell_{(i_k, u_k)(i_{k+1}, u_{k+1})}} \phi(i_{k+1}, u_{k+1}) \right)', \quad (3.21)$$

$$d_k^{(\lambda)} = \left(1 - \frac{1}{k+1} \right) d_{k-1}^{(\lambda)} + \frac{1}{k+1} z_k \alpha \frac{p_{(i_k, u_k)(i_{k+1}, u_{k+1})}}{\ell_{(i_k, u_k)(i_{k+1}, u_{k+1})}} g(i_k, u_k, i_{k+1}), \quad (3.22)$$

where z_k is generated by

$$z_k = \lambda \alpha \frac{p_{(i_{k-1}, u_{k-1})(i_k, u_k)}}{\ell_{(i_{k-1}, u_{k-1})(i_k, u_k)}} z_{k-1} + \phi(i_k, u_k). \quad (3.23)$$

Here $\phi(i, u)'$ is the row of Φ corresponding to (i, u) , $p_{(i,u)(j,v)}$ are the transition probabilities between state-control pairs given by Eq. (3.20), and $\ell_{(i,u)(j,v)}$ are the corresponding transition probabilities governing row and column sampling. A particularly interesting choice is

$$\ell_{(i,u)(j,v)} = p_{ij}(u) \nu(v | j), \quad (3.24)$$

where ν is a randomized policy that effects exploration by mapping a state j to control v with probability $\nu(v | j)$. Then from Eqs. (3.20) and (3.24), we see that the ratios in algorithm (3.21)-(3.23) simplify:

$$\frac{p_{(i,u)(j,v)}}{\ell_{(i,u)(j,v)}} = \frac{\delta(v = \mu(j))}{\nu(v | j)},$$

(again we adopt the convention $0/0 = 0$), so knowledge of $p_{ij}(u)$ is not needed, and the algorithm maintains its model-free character.

Recently, an alternative approach to Q-learning with exploration enhancement has been proposed, which is new even in the context of exact DP (Bertsekas and Yu [BeY10a]). It is based on replacing the Q-learning mapping F_μ of Eq. (3.18) with the mapping

$$(F_{J,\nu}Q)(i, u) = \sum_{j=1}^n p_{ij}(u) \left(g(i, u, j) + \alpha \sum_{v \in U(j)} \nu(v | j) \min\{J(j), Q(j, v)\} \right), \quad \forall (i, u), \quad (3.25)$$

which depends on a vector $J \in \Re^n$, with components denoted $J(i)$, and on a randomized policy ν , which for each state j defines a probability distribution

$$\{\nu(v | j) \mid v \in U(j)\}$$

over the feasible controls at j , and may depend on the “current policy” μ . Thus *solving a linear equation (the traditional policy evaluation method) is replaced by solving an optimal stopping problem*: finding a fixed point of the mapping (3.25) [a similarity with the constrained policy iteration (3.8)-(3.9) described earlier].

The exploration-enhanced policy iteration algorithm is similar to the ordinary version (3.17)-(3.19), except that it is based on the mapping $F_{J,\nu}$ in place of F_μ . Its policy evaluation phase is

$$\tilde{Q}(i, u, r) \approx (F_{J,\nu}Q)(i, u), \quad (3.26)$$

and its policy improvement phase is Eq. (3.19). At the end of an iteration, the vector J is set to

$$J(i) = \min_{u \in U(i)} (F_{J,\nu}\tilde{Q})(i, u), \quad \forall i. \quad (3.27)$$

The policy ν may be chosen arbitrarily at each iteration. It may encode aspects of the “current policy” μ , but may also allow for arbitrary and easily controllable amount of exploration. For extreme choices of ν and a lookup table representation, the algorithm yields as special cases the classical Q-learning/value iteration and policy iteration methods. Together with linear Q-factor approximation, the algorithm may be combined with the TD(0)-like method of Tsitsiklis and Van Roy [TsV99b], which can be used to solve the associated stopping problems with low overhead per iteration, thereby resolving the issue of exploration.

The enhanced policy iteration algorithm also has some interesting properties in the context of exact/lookup table policy iteration: it admits asynchronous and stochastic iterative implementations, which can be attractive alternatives to standard methods of asynchronous/optimistic policy iteration and Q-learning. Its advantage is more robust convergence properties (it is not susceptible to the fragility of asynchronous policy iteration demonstrated by Williams and Baird [WiB93], and of optimistic policy iteration with lookup table representation discussed by Bertsekas and Tsitsiklis [BeT96], Section 5.4); see Bertsekas and Yu [BeY10a], [BeY10b] for further elaboration. An important fact regarding these properties is that for every randomized policy ν , the mapping underlying the algorithm (3.26)-(3.27),

$$(Q, J) \mapsto (F_{J,\nu}Q, M F_{J,\nu}Q), \quad (3.28)$$

where

$$(M F_{J,\nu} Q)(i) = \min_{u \in U(i)} (F_{J,\nu} Q)(i, u), \quad \forall i,$$

is a sup-norm contraction of modulus α , and its unique fixed point is the optimal Q-factor and cost vector pair (Q^*, J^*) .

The contraction property just described, and other similar properties, form the basis for a variety of asynchronous stochastic Q-learning algorithms, which update Q-factors and relate to policy iteration with an optimistic character (a limited number of samples used to update Q for fixed J and ν , between updates of J). These algorithms may naturally involve enhanced exploration through the use of policies ν , which may relate but need not be identical to the policies under evaluation. We refer to Bertsekas and Yu [BeY10a] for a detailed discussion, convergence analysis, and also error bounds (in the case where Q-factors are compactly represented through a set of basis functions).

The contraction property of the mapping (3.28) may also be used to construct deterministic asynchronous policy iteration algorithms with an optimistic character, which update costs (rather than Q-factors); see [BeY10a], [BeY10b]. In these algorithms we use deterministic policies μ^k in place of ν , and for each state i there is a set of times $\mathcal{T}(i)$ at which we do a “policy evaluation” iteration at i , given by

$$V^{k+1}(i) = \sum_{j=1}^n p_{ij}(\mu^k(i)) \left(g(i, \mu^k(i), j) + \alpha \min\{J^k(j), V^k(j)\} \right),$$

and we leave $J^k(i)$, $\mu^k(i)$ unchanged. There is also an infinite set of times $\bar{\mathcal{T}}(i)$ [disjoint from $\mathcal{T}(i)$] at which we do a “policy improvement” at i , given by

$$J^{k+1}(i) = V^{k+1}(i) = \min_{u \in U(i)} \sum_{j=1}^n p_{ij}(u) \left(g(i, u, j) + \alpha \min\{J^k(j), V^k(j)\} \right),$$

while setting $\mu^{k+1}(i)$ to a u that attains the minimum above. Based on the sup-norm contraction property of mapping (3.28), it can be shown that this algorithm satisfies $\lim_{k \rightarrow \infty} J^k(i) = \lim_{k \rightarrow \infty} V^k(i) = J^*(i)$ for all i , with arbitrary initial conditions (J^0, V^0, μ^0) , thus overcoming the convergence difficulties demonstrated by Williams and Baird [WiB93] for the natural version of asynchronous policy iteration.

4. POLICY EVALUATION BY AGGREGATION

In this section we consider policy iteration with aggregation-based approximate policy evaluation. The aggregation approach resembles the projected equation approach because it also constructs cost approximations of the form Φr . However, there are some differences, the most important of which is that in aggregation methods there are restrictions in the form of Φ , so there is less freedom in choosing the approximation subspace. Moreover from a mathematical point of view, in the aggregation approach there are no projections with respect to Euclidean norms, the simulations can be done more flexibly, and the underlying contractions are with respect to the sup-norm rather than a Euclidean norm.

From a conceptual point of view, the most important difference from the projected equation approach is that aggregation is based on problem approximation (rather than Bellman equation approximation): the original problem is approximated with a related “aggregate” problem, which is then solved exactly to yield a cost-to-go approximation for the original problem. As a result, policy iteration with the aggregation approach is not afflicted by some of the pathologies that emanate from compact feature-based representations of cost

functions, such as policy oscillations and chattering. The sequence of generated policies converges (to an optimal policy of the aggregate problem) with monotonic cost improvement (for the aggregate problem), and a more favorable error bound can be derived than for the projected equation case.

To construct an aggregation framework for the discounted MDP of the preceding section, we introduce a finite set \mathcal{A} of aggregate states, and we introduce two sets of probabilities, which relate the original system states with the aggregate states:

- (1) For each aggregate state x and original system state i , we specify the *disaggregation probability* d_{xi} [we have $\sum_{i=1}^n d_{xi} = 1$ for each $x \in \mathcal{A}$]. Roughly, d_{xi} may be interpreted as the “degree to which x is represented by i .”
- (2) For each aggregate state y and original system state j , we specify the *aggregation probability* ϕ_{jy} (we have $\sum_{y \in \mathcal{A}} \phi_{jy} = 1$ for each $j = 1, \dots, n$). Roughly, ϕ_{jy} may be interpreted as the “degree of membership of j in the aggregate state y .” The vectors $\{\phi_{jy} \mid j = 1, \dots, n\}$ may also be viewed as basis functions that will be used to provide a cost approximation for the original problem.

We mention a few examples, and we refer to [Ber07], Vol. I, Section 6.3, and to [Ber10a] for further discussion and examples.

- (a) In *hard and soft aggregation*, we group the original system states into subsets, and we view each subset as an aggregate state. In hard aggregation each state belongs to one and only one subset, and the aggregation probabilities are

$$\phi_{jy} = 1 \quad \text{if system state } j \text{ belongs to aggregate state/subset } y.$$

The disaggregation probabilities could be

$$d_{xi} = 1/n_x \quad \text{if system state } i \text{ belongs to aggregate state/subset } x,$$

where n_x is the number of states of x (this implicitly assumes that all states that belong to aggregate state/subset y are “equally representative”). In soft aggregation, we allow the aggregate states/subsets to overlap, with the aggregation probabilities ϕ_{jy} quantifying the “degree of membership” of j in the aggregate state/subset y .

- (b) In various *discretization schemes*, each original system state j is associated with a convex combination of aggregate states:

$$j \sim \sum_{y \in \mathcal{A}} \phi_{jy} y,$$

for some nonnegative weights ϕ_{jx} , whose sum is 1, and which are viewed as aggregation probabilities (this makes geometrical sense if both the original and the aggregate states are associated with points in a Euclidean space).

- (c) In *coarse grid schemes*, a subset of representative states is chosen, each being an aggregate state. Thus, each aggregate state x is associated with a unique original state i_x , and we may use the disaggregation probabilities $d_{xi} = 1$ for $i = i_x$ and $d_{xi} = 0$ for $i \neq i_x$. The aggregation probabilities are chosen as in the preceding case (b).

The aggregation approach approximates the cost vector of a policy with Φr , where $r \in \mathbb{R}^s$ is a weight vector to be determined, and Φ is the matrix whose j th row consists of the aggregation probabilities $\phi_{j1}, \dots, \phi_{js}$. Thus aggregation involves an approximation architecture similar to the one of projected

equation methods: it uses as features the aggregation probabilities. There is a reciprocal possibility: starting from a set of s features for each state, we may construct a *feature-based hard aggregation* scheme by grouping together states with “similar features.” In particular, we may use a more or less regular partition of the feature space, which induces a possibly irregular partition of the original state space into aggregate states (all states whose features fall in the same set of the feature partition form an aggregate state). This is a general approach for passing from a feature-based approximation of the cost vector to an aggregation-based approximation (an interesting special case is when one of the features is a vector that approximates the cost vector of some known suboptimal policy). Unfortunately, in the resulting aggregation scheme the number of aggregate states may become very large.

The aggregation and disaggregation probabilities specify a dynamical system involving both aggregate and original system states. In this system:

- (i) From aggregate state x , we generate original system state i according to d_{xi} .
- (ii) We generate transitions from original system state i to original system state j according to $p_{ij}(u)$, with cost $g(i, u, j)$.
- (iii) From original system state j , we generate aggregate state y according to ϕ_{jy} .

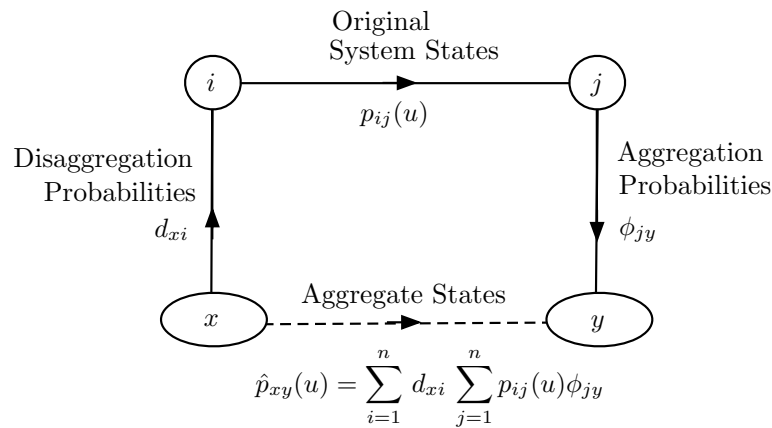


Figure 4.1. Illustration of the transition mechanism of a dynamical system involving both aggregate and original system states.

Let us now introduce the vectors \tilde{J}_0 , \tilde{J}_1 , and R^* where:

$R^*(x)$ is the optimal cost-to-go from aggregate state x .

$\tilde{J}_0(i)$ is the optimal cost-to-go from original system state i that has just been generated from an aggregate state (left side of Fig. 4.1).

$\tilde{J}_1(j)$ is the optimal cost-to-go from original system state j that has just been generated from an original system state (right side of Fig. 4.1).

Note that because of the intermediate transitions to aggregate states, \tilde{J}_0 and \tilde{J}_1 are different.

These three vectors satisfy the following three Bellman's equations:

$$\begin{aligned} R^*(x) &= \sum_{i=1}^n d_{xi} \tilde{J}_0(i), \quad x \in \mathcal{A}, \\ \tilde{J}_0(i) &= \min_{u \in U(i)} \sum_{j=1}^n p_{ij}(u) (g(i, u, j) + \alpha \tilde{J}_1(j)), \quad i = 1, \dots, n, \\ \tilde{J}_1(j) &= \sum_{y \in \mathcal{A}} \phi_{jy} R^*(y), \quad j = 1, \dots, n. \end{aligned}$$

By combining these equations, we obtain an equation for R^* :

$$R^*(x) = (FR^*)(x), \quad x \in \mathcal{A},$$

where F is the mapping defined by

$$(FR)(x) = \sum_{i=1}^n d_{xi} \min_{u \in U(i)} \sum_{j=1}^n p_{ij}(u) \left(g(i, u, j) + \alpha \sum_{y \in \mathcal{A}} \phi_{jy} R(y) \right), \quad x \in \mathcal{A}. \quad (4.1)$$

It can be seen that F is a sup-norm contraction mapping and has R^* as its unique fixed point. This follows from standard contraction arguments and the fact that d_{xi} , $p_{ij}(u)$, and ϕ_{jy} are all transition probabilities.†

Once R^* is found, the optimal-cost-to-go of the original problem may be approximated by $\tilde{J}_1 = \Phi R^*$, and a suboptimal policy may be found through the minimization defining \tilde{J}_0 . Again, the optimal cost function approximation \tilde{J}_1 is a linear combination of the columns of Φ , which may be viewed as basis functions.

4.1 Aggregation-Based Policy Iteration

A policy iteration algorithm to find R^* may be constructed, using aggregation for policy evaluation. It starts with a stationary policy μ^0 for the original problem, and given μ^k , it finds R_{μ^k} satisfying $R_{\mu^k} = F_{\mu^k} R_{\mu^k}$, where F_{μ} is the mapping defined by

$$(F_{\mu}R)(x) = \sum_{i=1}^n d_{xi} \sum_{j=1}^n p_{ij}(\mu(i)) \left(g(i, \mu(i), j) + \alpha \sum_{y \in \mathcal{A}} \phi_{jy} R_{\mu}(y) \right), \quad x \in \mathcal{A}, \quad (4.2)$$

(this is the policy evaluation step). It then generates μ^{k+1} by

$$\mu^{k+1}(i) \in \arg \min_{u \in U(i)} \sum_{j=1}^n p_{ij}(u) \left(g(i, u, j) + \alpha \sum_{y \in \mathcal{A}} \phi_{jy} R_{\mu^k}(y) \right), \quad \forall i, \quad (4.3)$$

† A quick proof is to observe that F is the composition $F = DT\Phi$, where T is the usual DP mapping, and D and Φ are the matrices with rows the disaggregation and aggregation distributions, respectively. Since T is a contraction with respect to the sup-norm $\|\cdot\|_{\infty}$, and D and Φ are sup-norm nonexpansive in the sense

$$\|Dx\|_{\infty} \leq \|x\|_{\infty}, \quad \forall x \in \mathfrak{R}^n, \quad \|\Phi y\|_{\infty} \leq \|y\|_{\infty}, \quad \forall y \in \mathfrak{R}^s,$$

it follows that F is a sup-norm contraction.

(this is the policy improvement step). Using the analysis of Section 3.6 and Prop. 3.1 (with W equal to the monotone mapping ΦD), it can be shown that this policy iteration algorithm terminates finitely to some policy corresponding to the unique fixed point of F . The key fact here is that F and F_μ are not only sup-norm contractions, but also have the monotonicity property of DP mappings, which is used in an essential way in the convergence proof of ordinary policy iteration.

We may implement the preceding policy iteration method by simulation, using either a matrix inversion or an iterative approach. For a given policy μ , the aggregate version of Bellman's equation, $R = F_\mu R$, is linear of the form [cf. Eq. (4.2)]

$$R = DT_\mu(\Phi R),$$

where T_μ is the DP mapping associated with μ . We can thus write this equation as

$$ER = f,$$

where

$$E = I - \alpha DP\Phi, \quad f = Dg, \quad (4.4)$$

in analogy with the corresponding matrix and vector for the projected equation.

We may use low-dimensional simulation to approximate E and f based on a given number of samples, similar to Section 2 [cf. Eq. (2.12)]. In particular, a sample sequence $\{(i_0, j_0), (i_1, j_1), \dots\}$ is obtained by first generating a sequence of states $\{i_0, i_1, \dots\}$ by sampling according to a distribution $\{\xi_i \mid i = 1, \dots, n\}$ (with $\xi_i > 0$ for all i), and then by generating for each t the column index j_t using sampling according to the distribution $\{p_{i_t j} \mid j = 1, \dots, n\}$. Given the first $k + 1$ samples, we form the matrix \hat{E}_k and vector \hat{f}_k given by

$$\hat{E}_k = I - \frac{\alpha}{k+1} \sum_{t=0}^k \frac{1}{\xi_{i_t}} d(i_t) \phi(j_t)', \quad \hat{f}_k = \frac{1}{k+1} \sum_{t=0}^k \frac{1}{\xi_{i_t}} d(i_t) g(i_t, \mu(i_t), j_t), \quad (4.5)$$

where $d(i)$ is the i th column of D and $\phi(j)'$ is the j th row of Φ . The convergence $\hat{E}_k \rightarrow E$ and $\hat{f}_k \rightarrow f$ follows from the expressions

$$E = I - \alpha \sum_{i=1}^n \sum_{j=1}^n p_{ij}(\mu(i)) d(i) \phi(j)', \quad f = \sum_{i=1}^n \sum_{j=1}^n p_{ij}(\mu(i)) d(i) g(i, \mu(i), j),$$

the relation

$$\lim_{k \rightarrow \infty} \sum_{t=0}^k \frac{\delta(i_t = i, j_t = j)}{k+1} = \xi_i p_{ij},$$

and law of large numbers arguments (cf. Section 2).

It is important to note that the sampling probabilities ξ_i are restricted to be positive, but are otherwise arbitrary and need not depend on the current policy. Moreover, their choice does not affect the obtained approximate solution of the equation $ER = f$. Because of this possibility, the problem of exploration is less acute in the context of policy iteration when aggregation is used for policy evaluation. This is in contrast with the projected equation approach, where the choice of ξ_i affects the projection norm and the solution of the projected equation, as well as the contraction properties of the mapping ΠT .

Note also that instead of using the probabilities ξ_i to sample original system states, we may alternatively sample the aggregate states x according to a distribution $\{\zeta_x \mid x \in \mathcal{A}\}$, generate a sequence of aggregate

states $\{x_0, x_1, \dots\}$, and then generate a state sequence $\{i_0, i_1, \dots\}$ using the disaggregation probabilities. In this case $\xi_i = \sum_{x \in \mathcal{A}} \zeta_x d_{xi}$ and Eqs. (4.5) should be modified as follows:

$$\hat{E}_k = I - \frac{\alpha}{k+1} \sum_{t=0}^k \frac{m}{d_{x_t i_t}} d(i_t) \phi(j_t)', \quad \hat{f}_k = \frac{1}{k+1} \sum_{t=0}^k \frac{m}{d_{x_t i_t}} d(i_t) g(i_t, \mu(i_t), j_t),$$

where m is the number of aggregate states.

The corresponding LSTD-type method generates $\hat{R}_k = \hat{E}_k^{-1} \hat{f}_k$, and approximates the cost vector of μ with the vector $\Phi \hat{R}_k$:

$$\tilde{J}_\mu = \Phi \hat{R}_k.$$

There is also a regression-based version that is suitable for the case where \hat{E}_k is nearly singular (cf. Section 2), as well as an iterative regression-based version of LSTD, which may be viewed as a special case of the (scaled) LSPE-type method (2.31)-(2.32). The latter method takes the form

$$\hat{R}_{k+1} = (\hat{E}_k' \Sigma_k^{-1} \hat{E}_k + \beta I)^{-1} (\hat{E}_k' \Sigma_k^{-1} \hat{f}_k + \beta \hat{R}_k), \quad (4.6)$$

where $\beta > 0$ and Σ_k is a positive definite symmetric matrix. Note that contrary to the projected equation case, for a discount factor $\alpha \approx 1$, \hat{E}_k will always be nearly singular [since $DP\Phi$ is a transition probability matrix, cf. Eq. (4.4)]. However, the iteration (4.6) is valid even if \hat{E}_k is singular.

As noted earlier, the nonoptimistic version of the aggregation-based policy iteration method does not exhibit the oscillatory behavior of the one based on the projected equation approach (cf. Section 3.6). The generated policies converge since by aggregation, we are essentially solving exactly the DP problem associated with the aggregate system of Fig. 4.1. The convergence of policies in turn implies that the error bound (3.2) holds. This suggests a more regular behavior and an advantage in terms of approximation quality of aggregation-based approximate policy iteration relative to its projected equation-based counterpart.

A related property holds for optimistic policy iteration methods. They behave similar to optimistic policy iteration for the DP problem associated with the aggregate system of Fig. 4.1, and they do not exhibit the chattering phenomenon described in Section 3.5.

The preceding arguments indicate that the aggregation approach holds an advantage over the projected equation approach in terms of regularity of behavior, error guarantees, and exploration-related difficulties. Note, however, that the basis functions in the aggregation approach are restricted by the requirement that the rows of Φ must be probability distributions. For example in the case of a single basis function ($s = 1$), there is only one possible choice for Φ in the aggregation context, namely the matrix whose single column is the unit vector. In practice, this means that the number of columns of Φ may have to be much larger than in the projected equation approach.

4.2 Multistep Aggregation

The aggregation methodology of this section can be generalized by considering a multistep aggregation-based dynamical system. This system, illustrated in Fig. 4.2, is specified by disaggregation and aggregation probabilities as before, but involves $k > 1$ transitions between original system states in between transitions from and to aggregate states.

We introduce vectors $\tilde{J}_0, \tilde{J}_1, \dots, \tilde{J}_k$, and R^* where:

$R^*(x)$ is the optimal cost-to-go from aggregate state x .

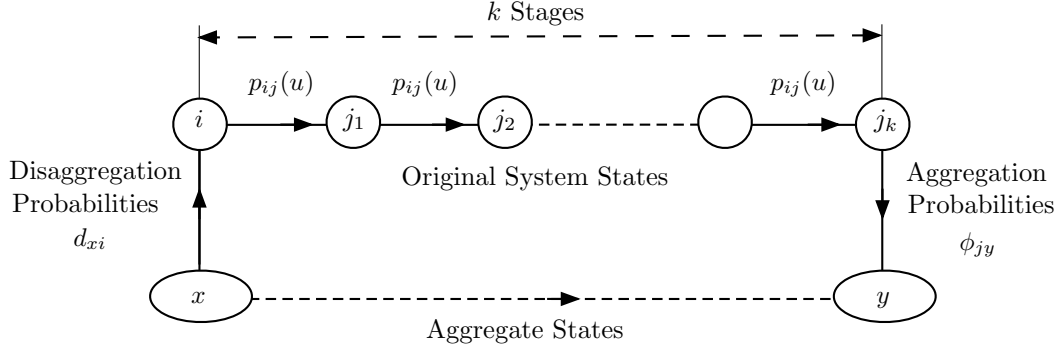


Figure 4.2. The transition mechanism for multistep aggregation. It is based on a dynamical system involving aggregate states, and k transitions between original system states in between transitions from and to aggregate states.

$\tilde{J}_0(i)$ is the optimal cost-to-go from original system state i that has just been generated from an aggregate state (left side of Fig. 4.2).

$\tilde{J}_1(j_1)$ is the optimal cost-to-go from original system state j that has just been generated from an original system state i .

$\tilde{J}_m(j_m)$, $m = 2, \dots, k$, is the optimal cost-to-go from original system state j_m that has just been generated from an original system state j_{m-1} .

These vectors satisfy the following set of Bellman equations:

$$R^*(x) = \sum_{i=1}^n d_{xi} \tilde{J}_0(i), \quad x \in \mathcal{A},$$

$$\tilde{J}_0(i) = \min_{u \in U(i)} \sum_{j_1=1}^n p_{ij_1}(u) (g(i, u, j_1) + \alpha \tilde{J}_1(j_1)), \quad i = 1, \dots, n, \quad (4.7)$$

$$\tilde{J}_m(j_m) = \min_{u \in U(j_m)} \sum_{j_{m+1}=1}^n p_{j_m j_{m+1}}(u) (g(j_m, u, j_{m+1}) + \alpha \tilde{J}_{m+1}(j_{m+1})), \quad j_m = 1, \dots, n, \quad m = 1, \dots, k-1, \quad (4.8)$$

$$\tilde{J}_k(j_k) = \sum_{y \in \mathcal{A}} \phi_{j_k y} R^*(y), \quad j_k = 1, \dots, n. \quad (4.9)$$

By combining these equations, we obtain an equation for R^* :

$$R^*(x) = (FR^*)(x), \quad x \in \mathcal{A},$$

where F is the mapping defined by

$$FR = DT^k(\Phi R),$$

where T is the usual DP mapping of the problem. As earlier, it can be seen that F is a sup-norm contraction, but its contraction modulus is α^k rather than α .

There is a similar mapping corresponding to a fixed policy and it can be used to implement a policy iteration algorithm, which evaluates a policy through calculation of a corresponding vector R and then

improves it. However, there is a major difference from the single-step aggregation case: a policy involves a set of k control functions $\{\mu_0, \dots, \mu_{k-1}\}$, and while a known policy can be easily simulated, its improvement involves multistep lookahead using the minimizations of Eqs. (4.7)-(4.9), and may be costly. Thus multistep aggregation is a useful idea only for problems where the cost of this multistep lookahead minimization (from a single given starting state) is not prohibitive. By contrast, policy improvement using projected equation-based multistep policy evaluation methods such as TD(λ), does not have this limitation. On the other hand, note that from the theoretical point of view, a multistep scheme provides a means of better approximation of the true optimal cost vector J^* , independent of the use of a large number of aggregate states. This can be seen from Eqs. (4.7)-(4.9), which based on classical value iteration convergence results, show that $\tilde{J}_0(i) \rightarrow J^*(i)$ as $k \rightarrow \infty$, regardless of the choice of aggregate states.

4.3 Distributed Asynchronous Aggregation

We now consider the distributed solution of discounted DP problems using aggregation. We envision a network of s processors/agents, with each processor updating asynchronously a local policy and a local cost function, defined on a subset of the state space. Aggregate estimates of the computed values are communicated asynchronously between processors, and are used to perform local cost updates and possibly local policy updates. In extreme cases, the aggregate estimates available to the agents may be greatly outdated, or otherwise may carry marginally useful information, thereby approaching a multi-agent system model where each agent makes decisions based solely on local information.

In a synchronous value iteration method, each processor $a = 1, \dots, s$, maintains/updates a (local) cost $J_a(i)$ for every state i in a subset of states I_a , and also an aggregate cost

$$R_a = \sum_{i \in I_a} d_{ai} J_a(i),$$

where $\{d_{ai} \mid i \in I_a\}$ is a probability distribution. We assume that the subsets I_a , $a = 1, \dots, s$, form a partition of the state space (as in hard aggregation, cf. Section 3.7), and we generically denote by J and R the vectors with components $J_a(i)$, $i \in I_a$, and R_a , $a = 1, \dots, s$, respectively. Note that J is an n -dimensional vector, the Cartesian product $J = \prod_{a=1, \dots, s} J_a$. For state i , $J(i)$ is the i th component of the vector J_a corresponding to the processor a for which $i \in I_a$.

We first discuss a value iteration algorithm, which involves updating both R and J according to

$$R_a^t = \sum_{i \in I_a} d_{ai} J_a^t(i), \quad \forall a = 1, \dots, s, \quad (4.10)$$

$$J_a^{t+1}(i) = \min_{u \in U(i)} H_a(i, u, J^t, R^t), \quad \forall i \in I_a, a = 1, \dots, s, \quad (4.11)$$

where the mapping H_a is defined for all $a = 1, \dots, s$, $i \in I_a$, $u \in U(i)$, and $J \in \mathfrak{R}^n$, $R \in \mathfrak{R}^s$, by

$$H_a(i, u, J, R) = \sum_{j=1}^n p_{ij}(u) g(i, u, j) + \alpha \sum_{j \in I_a} p_{ij}(u) J(j) + \alpha \sum_{j \notin I_a} p_{ij}(u) R_{a(j)}, \quad (4.12)$$

and where for each original system state j , we denote by $a(j)$ the subset to which j belongs [i.e., $j \in I_{a(j)}$]. Thus the iteration (4.11) is the same as ordinary value iteration, except that the aggregate costs $R_{a(j)}$ are used for the states $j \notin I_a$, whose costs are updated by other processors. It is possible to show that the

iteration (4.10)-(4.11) involves a sup-norm contraction mapping of modulus α , so it converges to the unique solution of the following system of equations in (J, R)

$$R_a = \sum_{i \in I_a} d_{ai} J_a(i), \quad \forall a = 1, \dots, s, \quad J_a(i) = \min_{u \in U(i)} H_a(i, u, J, R), \quad \forall i \in I_a. \quad (4.13)$$

(The contraction property follows from the fact that $\{d_{ai} \mid i \in I_a\}$ is a probability distribution.)

In the algorithm (4.10)-(4.11), all processors a must be updating their aggregate costs R_a and local costs $J_a(i)$ synchronously, and communicate the aggregate costs to the other processors before a new iteration may begin. In a more practical asynchronous version of the method, the aggregate costs R_a may be outdated because of communication “delays” between processors. In this case, the iteration (4.10)-(4.11) is modified to take the form

$$J_a^{t+1}(i) = \min_{u \in U(i)} H_a(i, u, J^t, R_1^{\tau_{a1}(t)}, \dots, R_s^{\tau_{as}(t)}), \quad \forall i \in I_a, \quad a = 1, \dots, s, \quad (4.14)$$

with $0 \leq \tau_{aa'}(t) \leq t$ for $a, a' = 1, \dots, s$, and

$$R_a^t = \sum_{i \in I_a} d_{ai} J_a^t(i), \quad \forall a = 1, \dots, s.$$

The differences $t - \tau_{aa'}(t)$, $a, a' = 1, \dots, s$, in Eq. (4.14) may be viewed as “delays” between the current time k and the times $\tau_{aa'}(t)$ when the corresponding aggregate costs were computed by other processors. Moreover, the iteration (4.14) may be executed only at a subset \mathbf{T}_a of times for each processor a , whereas for the other times $t \notin \mathbf{T}_a$ the vector J_a remains unchanged, i.e., $J_a^{t+1}(i) = J_a^t(i)$ for all $i \in I_a$. The convergence of this algorithm to the unique solution of Eq. (4.13) can be shown using the theory of asynchronous distributed DP developed in Bertsekas [Ber82], and has been discussed recently in the paper by Bertsekas and Yu [BeY10b].

An asynchronous policy iteration method involving aggregation operates similarly, except that each processor a performs a policy evaluation iteration at a subset of times $t \in \mathbf{T}_a$, a policy improvement iteration at a subset of times $t \in \bar{\mathbf{T}}_a \subset \mathbf{T}_a$, and no update at the remaining times. The local policy improvement iteration is

$$\mu_a^{t+1}(i) = \begin{cases} \arg \min_{u \in U(i)} H_a(i, u, J^t, R_1^{\tau_{a1}(t)}, \dots, R_m^{\tau_{am}(t)}) & \text{if } t \in \bar{\mathbf{T}}_a, \\ \mu_a^t(i) & \text{if } t \notin \bar{\mathbf{T}}_a, \end{cases} \quad \forall i \in I_a. \quad (4.15)$$

The local policy evaluation is iterative/optimistic according to

$$J_a^{t+1}(i) = \begin{cases} H_a(i, \mu_a^t(i), J^t, R_1^{\tau_{a1}(t)}, \dots, R_s^{\tau_{as}(t)}) & \text{if } t \in \mathbf{T}_a, \\ J_a^t(i) & \text{if } t \notin \mathbf{T}_a, \end{cases} \quad \forall i \in I_a. \quad (4.16)$$

Unfortunately, this algorithm can fail even under favorable circumstances, as shown by Williams and Baird [WiB93]: a lookup table representation (one processor for each state), only one processor updating at a time, and a deterministic system. An alternative and almost as simple algorithm was proposed recently in [BeY10b], which was shown to converge to the optimum under the most general conditions. In this algorithm, the iteration (4.16) is replaced by

$$J_a^{t+1}(i) = \begin{cases} \min \left\{ V_a^t(i), H_a(i, \mu_a^t(i), J^t, R_1^{\tau_{a1}(t)}, \dots, R_s^{\tau_{as}(t)}) \right\} & \text{if } t \in \mathbf{T}_a, \\ J_a^t(i) & \text{if } t \notin \mathbf{T}_a, \end{cases} \quad \forall i \in I_a. \quad (4.17)$$

where $V_a^t(i)$ is the value of the minimum computed during the latest policy improvement iteration (4.15):

$$V_a^{t+1}(i) = \begin{cases} \min_{u \in U(i)} H_a(i, u, J^t, R_1^{\tau_{a1}(t)}, \dots, R_s^{\tau_{as}(t)}) & \text{if } t \in \bar{\mathbf{T}}_a, \\ V_a^t(i) & \text{if } t \notin \bar{\mathbf{T}}_a, \end{cases} \quad \forall i \in I_a. \quad (4.18)$$

Other related algorithms, similarly convergent to the unique solution of Eq. (4.13), have also been discussed in [BeY10b].

5. CONCLUSIONS

We have surveyed some aspects of policy iteration methods with cost function approximation. From an analytical point of view, this is a subject with a rich theory and interesting algorithmic issues. From a practical point of view, this is a methodology that can address very large and difficult problems, and yet present major technical difficulties (such as exploration deficiencies and nearly singular equations) and exhibit unpredictable behaviors (such as policy oscillations and chattering), which are not fully understood at present and may seriously challenge practitioners.

We have contrasted two policy evaluation approaches: projected equation and aggregation. It appears that aggregation-based methods have more regular behavior, offer better error bound guarantees, and have less exploration-related difficulties than projected equation-based methods. On the other hand, aggregation methods are restricted in the choice of basis functions that they can use, and this can be a significant limitation for many problems.

6. REFERENCES

- [BBD10] Busoniu, L., Babuska, R., De Schutter, B., and Ernst, D., 2010. Reinforcement Learning and Dynamic Programming Using Function Approximators, CRC Press, N. Y.
- [BBN04] Bertsekas, D. P., Borkar, V. S., and Nedić, A., 2004. “Improved Temporal Difference Methods with Linear Function Approximation,” in *Learning and Approximate Dynamic Programming*, by J. Si, A. Barto, W. Powell, and D. Wunsch (Eds.), IEEE Press, N. Y.
- [BBS95] Barto, A. G., Bradtke, S. J., and Singh, S. P., 1995. “Real-Time Learning and Control Using Asynchronous Dynamic Programming,” *Artificial Intelligence*, Vol. 72, pp. 81-138.
- [BES10] Busoniu, L., Ernst, D., De Schutter, B., and Babuska, R., 2010. “Cross-Entropy Optimization of Control Policies with Adaptive Basis Functions,” *IEEE Trans. on Systems, Man, and Cybernetics - Part B: Cybernetics*, Vol. PP, pp. 1-14.
- [BKM05] de Boer, P. T., Kroese, D. P., Mannor, S., and Rubinstein, R. Y. 2005. “A Tutorial on the Cross-Entropy Method,” *Annals of Operations Research*, Vol. 134, pp. 19-67.
- [BSA83] Barto, A. G., Sutton, R. S., and Anderson, C. W., 1983. “Neuronlike Elements that Can Solve Difficult Learning Control Problems,” *IEEE Trans. on Systems, Man, and Cybernetics*, Vol. 13, pp. 835-846.
- [Bai93] Baird, L. C., 1993. “Advantage Updating,” Report WL-TR-93-1146, Wright Patterson AFB, OH.
- [BeG82] Bertsekas, D. P., and Gafni, E., 1982. “Projection Methods for Variational Inequalities with Applications to the Traffic Assignment Problem,” *Math. Progr. Studies*, Vol. 17, pp. 139-159.
- [BeI96] Bertsekas, D. P., and Ioffe, S., 1996. “Temporal Differences-Based Policy Iteration and Applications in Neuro-Dynamic Programming,” Lab. for Info. and Decision Systems Report LIDS-P-2349, MIT.
- [BeS78] Bertsekas, D. P., and Shreve, S. E., 1978. *Stochastic Optimal Control: The Discrete Time Case*, Academic Press, N. Y.; may be downloaded from <http://web.mit.edu/dimitrib/www/home.html>
- [BeT96] Bertsekas, D. P., and Tsitsiklis, J. N., 1996. *Neuro-Dynamic Programming*, Athena Scientific, Belmont, MA.
- [BeY07] Bertsekas, D. P., and Yu, H., 2007. “Solution of Large Systems of Equations Using Approximate Dynamic Programming Methods,” Lab. for Information and Decision Systems Report LIDS-P-2754, MIT.
- [BeY09] Bertsekas, D. P., and Yu, H., 2009. “Projected Equation Methods for Approximate Solution of Large Linear Systems,” *Journal of Computational and Applied Mathematics*, Vol. 227, pp. 27-50.

- [BeY10a] Bertsekas, D. P., and Yu, H., 2010. “Q-Learning and Enhanced Policy Iteration in Discounted Dynamic Programming,” Lab. for Information and Decision Systems Report LIDS-P-2831, MIT.
- [BeY10b] Bertsekas, D. P., and Yu, H., 2010. “Asynchronous Distributed Policy Iteration in Dynamic Programming,” Proc. of Allerton Conf. on Information Sciences and Systems.
- [Ber77] Bertsekas, D. P., 1977. “Monotone Mappings with Application in Dynamic Programming,” SIAM J. on Control and Optimization, Vol. 15, pp. 438-464.
- [Ber82] Bertsekas, D. P., 1982. “Distributed Dynamic Programming,” IEEE Transactions on Aut. Control, Vol. AC-27, pp. 610-616.
- [Ber96] Bertsekas, D. P., 1996. Lecture at NSF Workshop on Reinforcement Learning, Hilltop House, Harper’s Ferry, N.Y.
- [Ber07] Bertsekas, D. P., 2007. Dynamic Programming and Optimal Control, 3rd Edition, Vol. II, Athena Scientific, Belmont, MA.
- [Ber09a] Bertsekas, D. P., 2009. “Projected Equations, Variational Inequalities, and Temporal Difference Methods,” Lab. for Information and Decision Systems Report LIDS-P-2808, MIT; IEEE Trans. on Aut. Control, to appear in 2011.
- [Ber09b] Bertsekas, D. P., 2009. Convex Optimization Theory, Athena Scientific, Belmont, MA; on-line Chapter 6 at <http://web.mit.edu/dimitrib/www/convexdualitychapter.pdf>.
- [Ber10a] Bertsekas, D. P., 2010. Approximate Dynamic Programming, on-line at <http://web.mit.edu/dimitrib/www/dpchapter.html>.
- [Ber10b] Bertsekas, D. P., 2010. “Pathologies of Temporal Difference Methods in Approximate Dynamic Programming,” Proc. 2010 IEEE Conference on Decision and Control.
- [Bor08] Borkar, V. S., 2008. Stochastic Approximation: A Dynamical Systems Viewpoint, Cambridge Univ. Press.
- [Bor09] Borkar, V. S., 2009. “Reinforcement Learning: A Bridge Between Numerical Methods and Monte Carlo,” World Scientific Review, Vol. 9, Chapter 4.
- [Boy02] Boyan, J. A., 2002. “Technical Update: Least-Squares Temporal Difference Learning,” Machine Learning, Vol. 49, pp. 1-15.
- [BrB96] Bradtke, S. J., and Barto, A. G., 1996. “Linear Least-Squares Algorithms for Temporal Difference Learning,” Machine Learning, Vol. 22, pp. 33-57.
- [CFH07] Chang, H. S., Fu, M. C., Hu, J., Marcus, S. I., 2007. Simulation-Based Algorithms for Markov Decision Processes, Springer, N. Y.
- [CPS92] Cottle, R. W., Pang, J-S., and Stone, R. E., 1992. The Linear Complementarity Problem, Academic Press, N. Y.; republished by SIAM in 2009.
- [Cao07] Cao, X. R., 2007. Stochastic Learning and Optimization: A Sensitivity-Based Approach, Springer, N. Y.
- [ChV06] Choi, D. S., and Van Roy, B., 2006. “A Generalized Kalman Filter for Fixed Point Approximation and Efficient Temporal-Difference Learning, Discrete Event Dynamic Systems: Theory and Applications, Vol. 16, pp. 207-239.
- [DFM09] Desai, V. V., Farias, V. F., and Moallemi, C. C., 2009. “Approximate Dynamic Programming via a Smoothed Approximate Linear Program, Advances in Neural Information Processing Systems 22, MIT Press, Cambridge, MA.
- [Den67] Denardo, E. V., 1967. “Contraction Mappings in the Theory Underlying Dynamic Programming,” SIAM Review, Vol. 9, pp. 165-177.
- [DiM10] Di Castro, D., Mannor, S., 2010. “Adaptive Bases for Reinforcement Learning,” <http://arxiv.org/abs/1005.0125>.
- [FaV06] Farias, V. F., and Van Roy, B., 2006. “Tetris: A Study of Randomized Constraint Sampling, in Probabilistic

and Randomized Methods for Design Under Uncertainty, Springer-Verlag.

- [Fle84] Fletcher, C. A. J., 1984. Computational Galerkin Methods, Springer-Verlag, N. Y.
- [GLH94] Gurvits, L., Lin, L. J., and Hanson, S. J., 1994. “Incremental Learning of Evaluation Functions for Absorbing Markov Chains: New Methods and Theorems,” Preprint.
- [Gos03] Gosavi, A., 2003. Simulation-Based Optimization Parametric Optimization Techniques and Reinforcement Learning, Springer-Verlag, N. Y.
- [HBK94] Harmon, M. E., Baird, L. C., and Klopff, A. H., 1994. “Advantage Updating Applied to a Differential Game,” unpublished report, presented at the 1994 Neural Information Processing Systems Conference, Denver, CO.
- [JJS94] Jaakkola, T., Jordan, M. I., and Singh, S. P., 1994. “On the Convergence of Stochastic Iterative Dynamic Programming Algorithms,” Neural Computation, Vol. 6, pp. 1185-1201.
- [Kak02] Kakade, S., 2002. “A Natural Policy Gradient,” in Advances in Neural Information Processing Systems, Vol. 14, MIT Press, Cambridge, MA.
- [Kon02] Konda, V. R., 2002. Actor-Critic Algorithms, Ph.D. Thesis, Dept. of EECS, M.I.T., Cambridge, MA.
- [Kra72] Krasnoselskii, M. A., et. al, 1972. Approximate Solution of Operator Equations, Translated by D. Louvish, Wolters-Noordhoff Pub., Groningen.
- [LLL08] Lewis, F. L., Lendaris, G. G., and Liu, D., 2008. Special Issue on Adaptive Dynamic Programming and Reinforcement Learning in Feedback Control, IEEE Transactions on Systems, Man, and Cybernetics, Part B, Vol. 38.
- [LaP03] Lagoudakis, M. G., and Parr, R., 2003. “Least-Squares Policy Iteration,” J. of Machine Learning Research, Vol. 4, pp. 1107-1149.
- [LeV09] Lewis, F. L., and Vrabie, D., 2009. “Reinforcement Learning and Adaptive Dynamic Programming for Feedback Control,” IEEE Circuits and Systems Magazine, 3rd Q. Issue.
- [MMS06] Menache, I., Mannor, S., and Shimkin, N., 2005. “Basis Function Adaptation in Temporal Difference Reinforcement Learning,” Ann. Oper. Res., Vol. 134, pp. 215-238.
- [MSB08] Maei, H. R., Szepesvari, C., Bhatnagar, S., Silver, D., Precup, D., and Sutton, R. S., 2009. “Convergent Temporal-Difference Learning with Arbitrary Smooth Function Approximation,” Proc. NIPS 22.
- [Mar70] Martinet, B., 1970. “Regularisation d’ Inequations Variationnelles par Approximations Successives,” Rev. Francaise Inf. Rech. Oper., pp. 154-159.
- [Mey07] Meyn, S., 2007. Control Techniques for Complex Networks, Cambridge University Press, N. Y.
- [Mun03] Munos, R., 2003. “Error Bounds for Approximate Policy Iteration,” Proc. 20th ICML, Washington DC, pp. 560567.
- [NeB03] Nedić, A., and Bertsekas, D. P., 2003. “Least Squares Policy Evaluation Algorithms with Linear Function Approximation,” Discrete Event Dynamic Systems: Theory and Applications, Vol. 13, pp. 79-110.
- [Pin97] Pineda, F., 1997. “Mean-Field Analysis for Batched TD(λ),” Neural Computation, pp. 1403-1419.
- [Pow07] Powell, W. B., 2007. Approximate Dynamic Programming: Solving the Curses of Dimensionality, Wiley, N. Y.
- [Put94] Puterman, M. L., 1994. Markov Decision Processes: Discrete Stochastic Dynamic Programming, J. Wiley, N. Y.
- [Roc76] Rockafellar, R. T., 1976. “Monotone Operators and the Proximal Point Algorithm,” SIAM J. on Control and Optimization, Vol. 14, pp. 877-898.
- [Rot79] Rothblum, U. G., 1979. “Iterated Successive Approximation for Sequential Decision Processes,” in Stochastic Control and Optimization, by J. W. B. van Overhagen and H. C. Tijms (eds), Vrije University, Amsterdam.
- [RuK08] Rubinstein, R. Y., and Kroese, D. P., 2008. Simulation and the Monte Carlo Method, 2nd Edition, J. Wiley,

N. Y.

- [SBP04] Si, J., Barto, A., Powell, W., and Wunsch, D., (Eds.), 2004. *Learning and Approximate Dynamic Programming*, IEEE Press, N. Y.
- [SJJ94] Singh, S. P., Jaakkola, T., and Jordan, M. I., 1994. “Learning without State-Estimation in Partially Observable Markovian Decision Processes,” *Proceedings of the Eleventh Machine Learning Conference*, pp. 284-292.
- [SJJ95] Singh, S. P., Jaakkola, T., and Jordan, M. I., 1995. “Reinforcement Learning with Soft State Aggregation,” in *Advances in Neural Information Processing Systems 7*, MIT Press, Cambridge, MA.
- [SMP09] Sutton, R. S., Maei, H. R., Precup, D., Bhatnagar, S., Silver, D., Szepesvari, C., and Wiewiora, E., 2009. “Fast Gradient-Descent Methods for Temporal-Difference Learning with Linear Function Approximation,” *Proc. of ICML*.
- [SSM08] Sutton, R. S., Szepesvari, C., and Maei, H. R., 2008. “A Convergent $O(n)$ Algorithm for Off-Policy Temporal-Difference Learning with Linear Function Approximation,” *Proc. of NIPS* 21.
- [Sam59] Samuel, A. L., 1959. “Some Studies in Machine Learning Using the Game of Checkers,” *IBM Journal of Research and Development*, pp. 210-229.
- [Sam67] Samuel, A. L., 1967. “Some Studies in Machine Learning Using the Game of Checkers. II – Recent Progress,” *IBM Journal of Research and Development*, pp. 601-617.
- [Sch10] Scherrer, B., 2010. “Should One compute the Temporal Difference Fixed Point or Minimize the Bellman Residual? The Unified Oblique Projection View,” *Proc. of 2010 ICML, Haifa, Israel*.
- [Sha53] Shapley, L. S., 1953. “Stochastic Games,” *Proc. Nat. Acad. Sci. U.S.A.*, Vol. 39.
- [SuB98] Sutton, R. S., and Barto, A. G., 1998. *Reinforcement Learning*, MIT Press, Cambridge, MA.
- [Sut88] Sutton, R. S., 1988. “Learning to Predict by the Methods of Temporal Differences,” *Machine Learning*, Vol. 3, pp. 9-44.
- [SzL06] Szita, I., and Lorinz, A., 2006. “Learning Tetris Using the Noisy Cross-Entropy Method,” *Neural Computation*, Vol. 18, pp. 2936-2941.
- [Sze09] Szepesvari, C., 2009. “Reinforcement Learning Algorithms for MDPs,” *Dept. of Computing Science Report TR09-13, University of Alberta, Ca.*
- [Tes92] Tesauro, G. J., 1992. “Practical Issues in Temporal Difference Learning,” *Machine Learning*, Vol. 8, pp. 257-277.
- [ThS09] Thiery, C., and Scherrer, B., 2009. “Improvements on Learning Tetris with Cross-Entropy,” *International Computer Games Association Journal*, Vol. 32, pp. 23-33.
- [ThS10] Thiery, C., and Scherrer, B., 2010. “Least-Squares Policy Iteration: Bias-Variance Trade-off in Control Problems,” *Proc. of 2010 ICML, Haifa, Israel*.
- [TsV96] Tsitsiklis, J. N., and Van Roy, B., 1996. “Feature-Based Methods for Large-Scale Dynamic Programming,” *Machine Learning*, Vol. 22, pp. 59-94.
- [TsV97] Tsitsiklis, J. N., and Van Roy, B., 1997. “An Analysis of Temporal-Difference Learning with Function Approximation,” *IEEE Transactions on Automatic Control*, Vol. 42, pp. 674-690.
- [TsV99a] Tsitsiklis, J. N., and Van Roy, B., 1999. “Average Cost Temporal-Difference Learning,” *Automatica*, Vol. 35, pp. 1799-1808.
- [TsV99b] Tsitsiklis, J. N., and Van Roy, B., 1999. “Optimal Stopping of Markov Processes: Hilbert Space Theory, Approximation Algorithms, and an Application to Pricing Financial Derivatives,” *IEEE Transactions on Automatic Control*, Vol. 44, pp. 1840-1851.
- [Van95] Van Roy, B., 1995. “Feature-Based Methods for Large Scale Dynamic Programming,” *Lab. for Info. and Decision Systems Report LIDS-TH-2289, Massachusetts Institute of Technology, Cambridge, MA.*

- [Van06] Van Roy, B., 2006. “Performance Loss Bounds for Approximate Value Iteration with State Aggregation,” *Mathematics of Operations Research*, Vol. 31, pp. 234-244.
- [WPB09] Wang, M., Polydorides, N., and Bertsekas, D. P., 2009. “Approximate Simulation-Based Solution of Large-Scale Least Squares Problems,” Lab. for Information and Decision Systems Report LIDS-P-2819, MIT.
- [Wer92a] Werbös, P. J., 1992. “Approximate Dynamic Programming for Real-Time Control and Neural Modeling,” D. A. White, D. A. Sofge, (eds.), *Handbook of Intelligent Control*, Van Nostrand, N. Y.
- [Wer92b] Werbös, P. J., 1992. “Neurocontrol and Supervised Learning: an Overview and Valuation,” D. A. White, D. A. Sofge, (eds.), *Handbook of Intelligent Control*, Van Nostrand, N. Y.
- [WhS92] White, D., and Sofge, D., 1992. *Handbook of Intelligent Control*, Van Nostrand Reinhold, N.Y.
- [WiB93] Williams, R. J., and Baird, L. C., 1993. “Analysis of Some Incremental Variants of Policy Iteration: First Steps Toward Understanding Actor-Critic Learning Systems,” Report NU-CCS-93-11, College of Computer Science, Northeastern University, Boston, MA.
- [YaL08] Yao, H., and Liu, Z.-Q., 2008. “Preconditioned Temporal Difference Learning,” *Proc. of the 25th ICML*, Helsinki, Finland.
- [YuB06] Yu, H., and Bertsekas, D. P., 2006. “Convergence Results for Some Temporal Difference Methods Based on Least Squares,” Lab. for Information and Decision Systems Report 2697, MIT; *IEEE Trans. on Aut. Control*, Vol. 54, 2009, pp. 1515-153.
- [YuB07] Yu, H., and Bertsekas, D. P., 2007. “A Least Squares Q-Learning Algorithm for Optimal Stopping Problems,” Lab. for Information and Decision Systems Report 2731, MIT; also in *Proc. European Control Conference 2007*, Kos, Greece.
- [YuB09] Yu, H., and Bertsekas, D. P., 2009. “Basis Function Adaptation Methods for Cost Approximation in MDP,” *Proceedings of 2009 IEEE Symposium on Approximate Dynamic Programming and Reinforcement Learning (ADPRL 2009)*, Nashville, Tenn.
- [YuB10] Yu, H., and Bertsekas, D. P., 2010. “Error Bounds for Approximations from Projected Linear Equations,” *Mathematics of Operations Research*, Vol. 35, 2010, pp. 306-329; a shorter version appears in *Proc. of European Workshop on Reinforcement Learning (EWRL’08)*, July 2008, Lille, France.
- [Yu10a] Yu, H., 2010. “Least Squares Temporal Difference Methods: An Analysis Under General Conditions,” Technical report C-2010-39, Dept. Computer Science, Univ. of Helsinki.
- [Yu10b] Yu, H., 2010. “Convergence of Least Squares Temporal Difference Methods Under General Conditions,” *Proc. of the 27th ICML*, Haifa, Israel.