# Differential TD Learning for Value Function Approximation

Adithya M. Devraj and Sean P. Meyn*

December 27, 2018

## Abstract

Value functions arise as a component of algorithms as well as performance metrics in statistics and engineering applications. Computation of the associated Bellman equations is numerically challenging in all but a few special cases.

A popular approximation technique is known as Temporal Difference (TD) learning. The algorithm introduced in this paper is intended to resolve two well-known problems with this approach: In the discounted-cost setting, the variance of the algorithm diverges as the discount factor approaches unity. Second, for the average cost setting, unbiased algorithms exist only in special cases.

It is shown that the gradient of any of these value functions admits a representation that lends itself to algorithm design. Based on this result, the new *differential TD* method is obtained for Markovian models on Euclidean space with smooth dynamics.

Numerical examples show remarkable improvements in performance. In application to speed scaling, variance is reduced by two orders of magnitude.

**Keywords:** Reinforcement learning, Approximate dynamic programming, Poisson's equation, stochastic optimal control

**2000 AMS Subject Classification:** 93E20, 93E35, 60J20

## 1  Introduction

The value functions considered in this paper are based on a Markov chain $\boldsymbol{X} = \{X(t) : t = 0, 1, 2, \ldots\}$, taking values in $\mathbb{R}^d$, and an associated cost function $c : \mathbb{R}^d \to \mathbb{R}$. A critical modeling assumption is the evolution equation,

$$X(t+1) = a(X(t), N(t+1)), \quad t \in \mathbb{Z}_+, \tag{1}$$

in which $\boldsymbol{N} = \{N(t) : t = 0, 1, 2, \ldots\}$ is an $m$-dimensional i.i.d. disturbance sequence, and $a : \mathbb{R}^{d \times m} \to \mathbb{R}^d$ is continuous. Under these assumptions, $X(t+1)$ is a continuous function of its initial condition $X(0)$; this observation is the starting point for the construction of algorithms for value function approximation.

We begin with some familiar background.

## 1.1 Value functions in control and statistics

A common performance metric in stochastic control and finance is the total discounted cost:

$$h_\alpha(x) = \sum_{t=0}^{\infty} \alpha^t \mathsf{E}[c(X(t)) \mid X(0) = x], \tag{2}$$

where $\alpha \in (0,1)$ is the discount factor. The average cost is defined as the ergodic limit,

$$\bar{c} = \lim_{n \to \infty} \frac{1}{n} \sum_{t=0}^{n-1} \mathsf{E}[c(X(t)) \mid X(0) = x], \tag{3}$$

which is of interest in many areas beyond control engineering. The following *relative value function* is central to analysis of the average cost:

$$h(x) = \sum_{t=0}^{\infty} \mathsf{E}[\tilde{c}(X(t)) \mid X(0) = x], \tag{4}$$

where $\tilde{c} = c - \bar{c}$. In particular, under general conditions, the asymptotic variance (the variance appearing in the Central Limit Theorem for the ergodic average (3)) can be expressed in terms of the relative value function [1].

Under the assumptions imposed in this paper, the average cost is deterministic and independent of $X(0) = x$. Moreover, the relative value function solves *Poisson's equation*:

$$\mathsf{E}[h(X(t+1)) - h(X(t)) \mid X(t) = x] = -\tilde{c}(x). \tag{5}$$

Significant applications include,

*Optimal control:* The policy iteration algorithm is used to compute an optimal policy based on two steps. Given a policy, it is first necessary to obtain the associated value function. The second step is to update the policy based on this value function [2]. This approach is used for both discounted and average-cost optimal control problems.

Poisson's equation finds application in many other fields:

*Variance reduction:* The control variate method is intended to reduce variance for various Monte-Carlo methods; a version of this technique involves the construction of an approximate solution to Poisson's equation [1,3].

*Nonlinear filtering:* A recent approach to approximate nonlinear filtering requires the solution to Poisson's equation to obtain the innovation gain [4]. Approximations are required for efficient implementation of this method.

We next recall the basic ideas surrounding TD-learning algorithms for value function approximation. The discussion is restricted to discounted-cost value functions.

## 1.2 TD-learning and value function approximation

Closed-form expressions for any of the value functions (2) or (5) is impossible in all but a few special cases, such as linear systems with quadratic cost. One approach to approximation is a simulation based algorithm known as Temporal Difference (TD) learning [5,6].

The goal of TD-learning is to approximate the value function $h_\alpha$ using a parameterized family of functions $\{h_\alpha^\theta : \theta \in \mathbb{R}^\ell\}$. Throughout most of the paper we restrict to a linear parameterization of the form

$$h_\alpha^\theta = \sum_{j=1}^{\ell} \theta_j \psi_j, \tag{6}$$

1

where $\psi \colon \mathbb{R}^d \to \mathbb{R}^\ell$ is continuously differentiable. The optimal parameter $\theta^*$ is the solution to a minimum-norm problem,

$$
\begin{aligned}
\theta^* &= \arg\min_\theta \|h_\alpha^\theta - h_\alpha\|_\pi^2 \\
&= \arg\min_\theta \mathsf{E}[(h_\alpha^\theta(X) - h_\alpha(X))^2],
\end{aligned}
\tag{7}
$$

where the expectation is with respect to $X \sim \pi$, the steady-state distribution; see Section 2.1 for details.

Theory for TD learning in the discounted cost setting is largely complete, in the sense that criteria for convergence are well-understood, and the asymptotic variance of the algorithm is computable based on standard theory from stochastic approximation theory [7]. Theory and algorithms for the average cost setting is more fragmented. The analog of (7) with $h_\alpha$ replaced by the relative value function can be solved using TD-learning techniques only for Markovian models that regenerate: there exists a single state $x^\bullet$ that is visited infinitely often [1,8].

Regeneration is not a restrictive assumption in many cases. However, the asymptotic variance of the algorithms introduced in [1] grows with the variance of inter-regeneration times. The variance can be massive in simple examples such as the M/M/1 queue. High variance is also predominantly observed in the discounted cost case when the discounting factor is close to 1.

The *differential* TD-learning algorithms introduced in this paper are designed to resolve these issues. The main idea is to estimate the *gradient* of the value function directly. Under the conditions imposed, the variance remains uniformly bounded over $0 < \alpha < 1$, and is also applicable for approximating the solution to Poisson's equation.

## 1.3 Differential TD-learning

In $\nabla$-TD learning algorithms, the gradient of the value function is approximated rather than the function itself.

Consider again the discounted-cost setting, and suppose that both $h_\alpha$ and each of the potential approximations $\{h_\alpha^\theta : \theta \in \mathbb{R}^\ell\}$ are continuously differentiable ($C^1$) as a function of $x$. In most of the paper, algorithms and analysis are restricted to the linear parameterization (6), so that

$$
\nabla h_\alpha^\theta = \sum_{j=1}^\ell \theta_j \nabla \psi_j.
\tag{8}
$$

The $\nabla$-TD learning algorithm is designed to compute the solution to the following nonlinear program:

$$
\theta^* = \arg\min_\theta \mathsf{E}[\|\nabla h_\alpha^\theta(X) - \nabla h_\alpha(X)\|^2], \quad X \sim \pi.
\tag{9}
$$

Once the optimal parameter has been obtained, the approximate value function requires the addition of a constant,

$$
h_\alpha^{\theta^*} = \sum_{j=1}^\ell \theta_j \psi_j + \kappa(\theta^*).
\tag{10}
$$

The optimal choice of $\kappa(\theta^*)$ is also obtained in the algorithms described in this paper.

In summary, the contributions of this work are,

1(a) The new $\nabla$-TD algorithms are applicable for either discounted- and average-cost.

(b) For a linear parameterization, the $\nabla$-LSTD algorithm solves the quadratic program (9).

2

(c) Extensions to nonlinear parameterizations are obtained.

In the discounted-cost setting, the algorithms also compute the optimal constant $\kappa(\theta^*)$ appearing in (10).

2. The new algorithms are applicable for models that do not have regeneration, and under general conditions the variance is uniformly bounded over all $0 < \alpha < 1$.

These algorithms do have limitations. First, they are only applicable in settings where the gradient is a meaningful concept. However, in numerical experiments we find that a pseudo-gradient can be defined even for a model with discrete state space, and the resulting ad-hoc algorithm is remarkably effective.

Also, in its current formulation, the algorithms introduced here fall into the category of *Approximate Dynamic Programming* (ADP), rather than *Reinforcement Learning* (RL): The algorithm relies on simulating a model of the system, rather than estimating a value function based on observations of a physical system. This distinction is not absolute: For example, in the applications to speed-scaling presented in Section 5, the $\nabla$-LSTD learning algorithm does fall into the class of RL algorithms.

The remainder of the paper is organized as follows: basic definitions and value function representations are presented in Section 2. The $\nabla$-LSTD learning algorithm is introduced in Section 3. Results from numerical experiments are surveyed in Section 5, and conclusions are contained in Section 6.

## 2 Representations and approximations

We begin with assumptions, and representations for value functions and their gradients.

### 2.1 Markovian model and problem formulation

The evolution equations (1) define a Markov chain $\boldsymbol{X}$ with transition semigroup defined for $t \geq 0$, $x \in \mathbb{R}^d$, and $A \in \mathcal{B}(\mathbb{R}^d)$, via

$$P^t(x, A) := \mathsf{P}_x\{X(t) \in A\} := \Pr\{X(t) \in A \,|\, X(0) = x\}.$$

For $t = 1$ we write $P = P^1$, which has the following form:

$$P(x, A) = \Pr\{a(x, N(1)) \in A\}.$$

The first set of assumptions ensures that each of the value functions exists. Fix a continuous function $v \colon \mathbb{R}^d \to [1, \infty)$ that serves as a weighting function. For any measurable function $f \colon \mathbb{R}^d \to \mathbb{R}$, the $v$-norm is denoted by,

$$\|f\|_v := \sup_x \frac{|f(x)|}{v(x)}.$$

The set of all measurable functions for which $\|f\|_v$ is finite is denoted $L_\infty^v$.

**Assumption A1:**

**A1.1**: The Markov chain is *v-uniformly ergodic*: There exists a unique invariant probability measure $\pi$, $b_0 < \infty$, and $0 < \rho_0 < 1$, such that for each function $f \in L_\infty^v$,

$$\left| \mathsf{E}_x[f(X(t))] - \pi(f) \right| \leq b_0 \rho_0^t \|f\|_v v(x), \quad t \geq 0, \tag{11}$$

where $\pi(f)$ denotes the steady-state mean of $f$. □

It is well known that (A1) is equivalent to the existence of a Lyapunov function: drift condition (V4) of [1]. The following consequences are immediate:

**Proposition 2.1.** *The following hold under (A1), and the bound $\|c\|_v < \infty$: The limit $\bar{c}$ in (3) exists, with $\bar{c} = \pi(c) < \infty$, and is independent of the initial condition $x$. Moreover, there exists $b_c < \infty$ such that:*

$$|h_\alpha(x)| \leq b_c\big(v(x) + (1 - \alpha)^{-1}\big)$$
$$|h_\alpha(x) - h_\alpha(y)| \leq b_c\big(v(x) + v(y)\big)$$
$$and \qquad |h(x)| \leq b_c v(x), \qquad x, y \in \mathbb{R}^d$$

*where $h_\alpha$ is defined in (2), and $h$ is defined in (4).* ∎

The following operator-theoretic notation will simplify exposition. For any measurable function $f \colon \mathbb{R}^d \to \mathbb{R}$ the new function $P^t f$ is defined as the conditional expectation

$$P^t f(x) = \mathsf{E}_x[f(X(t))] := \mathsf{E}[f(X(t)) \mid X(0) = x].$$

The *resolvent kernel* is the "$z$-transform" of the semi-group,

$$R_\alpha := \sum_{t=0}^{\infty} \alpha^t P^t, \quad 0 < \alpha < 1. \tag{12}$$

Under the assumptions of Prop. 2.1 we have

$$h_\alpha = R_\alpha c. \tag{13}$$

The solution to Poisson's equation has a similar representation under these assumptions – an operator-theoretic representation of (4) [9].

The representation (13) is valuable in deriving the TD-learning algorithm [1, 6]. We seek a similar representation for the gradient $\nabla h_\alpha = \nabla[R_\alpha c]$.

## 2.2 Representation for the gradient of a value function

The goal here is to obtain an operator $\Omega_\alpha$ for which the following holds:

$$\nabla h_\alpha = \Omega_\alpha \nabla c. \tag{14}$$

This requires assumptions on the model and the cost function. In the following, a heuristic construction of $\Omega_\alpha$ is presented, with justifications collected together in Section 3.2. For complete details, the reader is referred to [10].

The representation requires additional assumptions:

**Assumption A2:**

**A2.1**: The disturbance process $\boldsymbol{N}$ does not depend upon the initial condition $X(0)$.

**A2.2**: The function $a$ is continuously differentiable in its first variable, with

$$\sup_{x,n} \|\nabla a(x, n)\| < \infty$$

in which $\| \cdot \|$ is any matrix norm, and the $i$th column of the $d \times d$ matrix $\nabla a$ is equal to $\nabla a_i$. □

The first assumption A2.1 is critical so that the initial state $X(0) = x$ can be regarded as a variable, with $X(t)$ being a continuous function of $x$. Assumption A2 allows us to define the *sensitivity process* $\mathcal{S}(t)$:

$$\mathcal{S}_{i,j}(t) := \frac{\partial X_i(t)}{\partial X_j(0)}, \quad 1 \le i, j \le d. \tag{15}$$

From (1), the sensitivity process evolves according to the random linear system

$$\mathcal{S}(t+1) = \mathcal{A}^T(t+1)\mathcal{S}(t), \quad \mathcal{S}(0) = I, \tag{16}$$

where $\mathcal{A}(t) := \nabla a\left(X(t-1), N(t)\right)$.

The operator $\widetilde{\nabla}$ is defined for any $C^1$ function $f$ via

$$\widetilde{\nabla} f(X(t)) := \mathcal{S}^T(t) \nabla f(X(t)). \tag{17}$$

It follows from the chain rule that this coincides with the gradient of $f(X(t))$ with respect to the initial conditions.

The interpretation of (17) motivates the introduction of a semi-group $\{Q^t\}$ of operators, whose domain includes functions $g : \mathbb{R}^d \to \mathbb{R}^d$ for which $g_i \in L_\infty^v$ for each $i$. For $t = 0$, it is the identity operator, and for $t \ge 1$,

$$Q^t g(x) := \mathsf{E}_x\left[\mathcal{S}^T(t) g(X(t))\right]. \tag{18}$$

Provided we can exchange the gradient and the expectation,

$$\frac{\partial}{\partial x_i} \mathsf{E}[f(X(t))] = \mathsf{E}\left[[\widetilde{\nabla} f(X(t))]_i\right]$$

which implies that $\nabla P^t f(x) = \mathsf{E}_x[\widetilde{\nabla} f(X(t))] = Q^t \nabla f(x)$. Under further assumptions we obtain

$$\nabla h_\alpha = \sum_{t=0}^{\infty} \alpha^t \nabla P^t c = \sum_{t=0}^{\infty} \alpha^t Q^t \nabla c.$$

We then obtain (14), with

$$\Omega_\alpha g := \sum_{t=0}^{\infty} \alpha^t Q^t g. \tag{19}$$

The representation (14) is the basis of the $\nabla$-LSTD learning algorithms developed in this paper.

Under the conditions of Prop. 3.4 the operator $\Omega_\alpha$ is well defined on a large domain of functions, with uniform bound over all $0 \le \alpha \le 1$. In the special case of $\alpha = 1$ we denote $\Omega = \Omega_1$, which under these conditions provides the representation $\nabla h = \Omega \nabla c$ for the gradient of the relative value function.

## 3   Differential TD learning

Algorithms are developed here for the Markov model (1), subject to Assumptions A1 and A2. The algorithms are presented first, with supporting theory postponed to Section 3.2.

Until Section 3.3 we restrict to a linear parameterization, in which $h_\alpha^\theta = \theta^T \psi$.

## 3.1 LSTD algorithms

We begin with a review of the standard algorithm, which is defined by the following recursion.

**Least squares TD-learning algorithm**

$$
\begin{aligned}
\varphi(t) &= \alpha\varphi(t-1) + \psi(X(t)) \\
b(t) &= (1-\gamma_t)b(t-1) + \gamma_t\varphi(t)c(X(t)) \\
M(t) &= (1-\gamma_t)M(t-1) + \gamma_t\psi(X(t))\psi^T(X(t)),
\end{aligned}
\tag{20}
$$

and obtain $\theta(t) = M^{-1}(t)b(t)$. The algorithm is initialized with $b(0)$, $\varphi(0) \in \mathbb{R}^\ell$, and $M(0) > 0$ a positive-definite $\ell \times \ell$ matrix [1,6].

Throughout the paper the gain sequence appearing in (20) and elsewhere is taken to be $\gamma_t = 1/t$, $t \geq 1$.

To simplify discussion we restrict to a stationary setting in the convergence results in this paper:

**Proposition 3.1.** *Suppose that (A1) holds, and that $c^2$ and $\|\psi\|^2$ are in $L_\infty^v$. Suppose moreover that the matrix $M = \mathsf{E}_\pi[\psi(X)\psi(X)^T]$ is full rank, where $X \sim \pi$. Then there is a version of the pair process $(\boldsymbol{X}, \boldsymbol{\varphi})$ that is stationary. For any initial conditions $b(0) \in \mathbb{R}^\ell$ and $M(0) > 0$, the algorithm is consistent:*

$$
\theta^* = \lim_{t\to\infty} M^{-1}(t)b(t).
$$

*Proof.* The existence of a stationary solution $\boldsymbol{X}$ follows directly from $v$-uniform ergodicity, and we then define

$$
\varphi(t) = \sum_{i=0}^{\infty} \alpha^i \psi(X(t-i)).
$$

It is known that the optimal parameter can be expressed $\theta^* = M^{-1}b$ in which $b = \mathsf{E}_\pi[\varphi(t)c(X(t))]$, so the result follows from the Law of Large Numbers for stationary processes. □

In the construction of the LSTD algorithm, the optimization problem (7) is cast in the Hilbert space,

$$
L_2^\pi = \left\{ \text{measurable } h \colon \mathbb{R}^d \to \mathbb{R} \ : \ \|h\|_\pi^2 = \langle h, h\rangle_\pi < \infty \right\}
$$

with $\langle f, g\rangle_\pi := \int f(x)g(x)\pi(dx)$. The $\nabla$-LSTD algorithm is based on a different inner product to define the norm in the approximation error.

For $C^1$ functions $f$, $g$ for which $\|\nabla f\|, \|\nabla g\| \in L_2^\pi$, define the inner product

$$
\langle f, g\rangle_{\pi,1} = \int \nabla f(x)^T \nabla g(x)\pi(dx),
$$

with the associated norm $\|f\|_{\pi,1} = \sqrt{\langle f, f\rangle_{\pi,1}}$. The nonlinear program (9) can be recast as

$$
\theta^* = \arg\min_\theta \|h_\alpha^\theta - h_\alpha\|_{\pi,1}.
\tag{21}
$$

Consider the linear parameterization (8) in which $\psi \colon \mathbb{R}^d \to \mathbb{R}^d$ is continuously differentiable, and assume as well that $c$ is continuously differentiable. The $\nabla$-LSTD learning algorithm is then defined by the following recursion:

**Differential least squares TD-learning algorithm**

$$
\begin{aligned}
\varphi(t) &= \alpha\mathcal{A}^T(t)\varphi(t-1) + \nabla\psi(X(t)) \\
b(t) &= (1-\gamma_t)b(t-1) + \gamma_t\varphi(t)^T\nabla c(X(t)) \\
M(t) &= (1-\gamma_t)M(t-1) + \gamma_t\nabla\psi(X(t))\nabla\psi(X(t))^T
\end{aligned}
\tag{22}
$$

where $\nabla \psi (x)$ denotes the $d \times \ell$ matrix

$$[\nabla \psi (x)]_{i,j} = \frac{\partial}{\partial x_i} \psi_j(x), \quad x \in \mathbb{R}^d, \tag{23}$$

and the parameters are obtained as, $\theta(t) = M^{-1}(t)b(t)$. Once again, $M(0) > 0$ is an arbitrary $\ell \times \ell$ positive-definite matrix, $b(0) \in \mathbb{R}^\ell$, $\varphi(0) \in \mathbb{R}^{d \times \ell}$ are arbitrary initializations.

Two more steps are required to obtain an estimate of $h_\alpha$. To ensure that $\pi(h_\alpha) = \pi(h_\alpha^\theta)$ we take $h_\alpha^\theta = \theta^T \psi + \kappa(\theta)$ where the constant is

$$\kappa(\theta) = -\pi(h_\alpha^\theta) + \pi(c)/(1-\alpha).$$

The two means can be estimated recursively:

$$\bar{h}_\alpha(t) = (1-\gamma_t)\bar{h}_\alpha(t-1) + \gamma_t h_\alpha^{\theta(t)} \tag{24}$$
$$\bar{c}(t) = (1-\gamma_t)\bar{c}(t-1) + \gamma_t c(X(t)). \tag{25}$$

It is immediate that $\bar{c}(t) \to \bar{c}$ as $t \to \infty$ by the Law of Large Numbers for $v$-uniformly ergodic Markov chains [9]. Convergence of $\bar{h}_\alpha(t)$ to $\pi(h_\alpha^{\theta^*})$ requires further assumptions.

This completes the description of the $\nabla$-LSTD learning algorithm.

## 3.2 Derivation and analysis

For a linear parameterization, the optimal parameter is the minimum of a quadratic. The proof of Prop. 3.2 follows immediately from the definition of the norm $\| \cdot \|_{\pi,1}$.

**Proposition 3.2.** *The norm appearing in* (21) *is a quadratic form,*

$$\|h_\alpha^\theta - h_\alpha\|_{\pi,1}^2 = \theta^T M \theta - 2b^T \theta + k, \tag{26}$$

*in which for each* $1 \le i, j \le j$,

$$M_{i,j} = \langle \psi_i, \psi_j \rangle_{\pi,1}, \quad b_i = \langle \psi_i, h_\alpha \rangle_{\pi,1}, \tag{27}$$

*and* $k = \langle h_\alpha, h_\alpha \rangle_{\pi,1}$. *Consequently, the optimizer* (21) *is any solution to*

$$M\theta^* = b. \tag{28}$$

∎

The matrix $M$ can be expressed in more compact notation:

$$M = \mathsf{E}_\pi[(\nabla \psi(X))^T \nabla \psi(X)]. \tag{29}$$

where $\nabla \psi$ is defined in (23). Similarly,

$$b = \mathsf{E}_\pi \big[ [\nabla \psi(X)]^T \nabla h_\alpha(X) \big]. \tag{30}$$

As in the standard TD learning algorithm, the vector $b$ is represented using the function $h_\alpha$, which is unknown. An alternative representation can be obtained whenever (14) is valid, and this is the basis of the $\nabla$-LSTD algorithm.

The following assumptions are used to justify this representation:

7

**Assumption A3:** For any $C^1$ functions $f, g$ satisfying $f^2, g^2 \in L_\infty^v$ and $\|\nabla f\|^2, \|\nabla g\|^2 \in L_\infty^v$, the following hold for the stationary version of the Markov chain:

$$\sum_{t=0}^{\infty} \left| \mathsf{E}_\pi \left[ \nabla f(X(t)) \right)^T \mathcal{S}(t) \nabla g(X(0)) \right] \right| < \infty \tag{31}$$

$$\sum_{t=0}^{\infty} \mathsf{E}_\pi \left[ \left| \nabla f(X(t)) \right)^T \mathcal{S}(t) \nabla g(X(0)) \right| \right] < \infty \tag{32}$$

Under (31) the right hand side of (14) is well defined a.e. $[\pi]$ when $c$ satisfies these conditions. General conditions for the validity of (31) are established in [10]. Theory to justify (32) is not as well developed. The condition is related to the existence of a negative Lyapunov exponent.

Under these assumptions we obtain a stationary solution for the pair $(\boldsymbol{X}, \boldsymbol{\varphi})$. The representation of $\boldsymbol{\varphi}$ requires the following shift-operator on sample space for a stationary version of $\boldsymbol{X}$: For a random variable of the form $Z = F(X(r), N(r), \ldots, X(s), N(s))$ with $r \leq s$ we denote, for any integer $k$,

$$\Theta^k Z = F(X(r+k), N(r+k), \ldots, X(s+k), N(s+k))$$

Consequently,

$$\Theta^k \mathcal{S}(t) = [\mathcal{A}(1+k)\mathcal{A}(2+k) \cdots \mathcal{A}(t+k)]^T. \tag{33}$$

Lemma 3.3 follows immediately from the assumptions: It follows from the definition (34) that this process follows the first recursion in (22).

**Lemma 3.3.** *Suppose that Assumptions A1–A3 hold, and that $\|\psi\|^2$ and $\|\nabla\psi\|^2$ are in $L_\infty^v$. Then there is a version of the pair process $(\boldsymbol{X}, \boldsymbol{\varphi})$ that is stationary, with*

$$\varphi(t) = \sum_{k=0}^{\infty} \alpha^k \left[ \Theta^{t-k} \mathcal{S}(k) \right] \nabla\psi(X(t-k)), \quad t \in \mathbb{Z}. \tag{34}$$

■

**Proposition 3.4.** *Suppose that Assumptions A1–A3 hold, and that $c^2$, $\|\nabla c\|$, $\|\psi\|^2$ and $\|\nabla\psi\|^2$ are in $L_\infty^v$. Suppose moreover that the matrix $M$ appearing in (29) is full rank. Then, for the stationary process $(\boldsymbol{X}, \boldsymbol{\varphi})$, the $\nabla$-LSTD learning algorithm is consistent: for any initial $b(0) \in \mathbb{R}^\ell$ and $M(0) > 0$,*

$$\theta^* = \lim_{t \to \infty} M^{-1}(t) b(t).$$

*Moreover, with probability one,*

$$\bar{c} = \lim_{t \to \infty} \bar{c}(t), \quad \pi(h_\alpha^{\theta^*}) = \lim_{t \to \infty} \bar{h}_\alpha(t),$$

*and hence $\lim_{t \to \infty} \{-\bar{h}_\alpha(t) + \bar{c}(t)/(1-\alpha)\} = \kappa(\theta^*)$.*

The remainder of this section consists of a proof of this proposition. We begin with a representation of $b$:

**Lemma 3.5.** *Under the assumptions of Prop. 3.4,*

$$b^T = \sum_{t=0}^{\infty} \alpha^t \mathsf{E}_\pi \left[ \left( \mathcal{S}^T(t) \nabla c(X(t)) \right)^T \nabla\psi(X(0)) \right] \tag{35}$$

$$= \mathsf{E} \left[ \left( \nabla c(X(0)) \right)^T \varphi(0) \right].$$

*in which $\boldsymbol{X}$ is stationary, with marginal $\pi$.*

*Proof.* The representation (14) is valid under (A3). Using this and (16) gives the first representation in (35):

$$b^T = \int \mathsf{E}_x\big[(\Omega_\alpha \nabla c(x))^T \nabla \psi(x)\big]\pi(dx)$$

$$= \sum_{t=0}^{\infty} \alpha^t \int \mathsf{E}_x\big[(\mathcal{S}^T(t)\nabla c(x))^T \nabla \psi(x)\big]\pi(dx) \tag{36}$$

$$= \sum_{t=0}^{\infty} \alpha^t \mathsf{E}_\pi\big[(\mathcal{S}^T(t)\nabla c(X(t)))^T \nabla\psi(X(0))\big]$$

Stationarity implies that for any $t, k \in \mathbb{Z}$,

$$\mathsf{E}_\pi\Big[\big(\mathcal{S}^T(t)\nabla c(X(t))\big)^T \nabla\psi(X(0))\Big]$$
$$= \mathsf{E}_x\Big[\big([\Theta^k \mathcal{S}^T(t)]\nabla c(X(t+k))\big)^T \nabla\psi(X(k))\Big].$$

Setting $k = -t$, the first representation in (35) becomes:

$$b^T = \sum_{t=0}^{\infty} \alpha^t \mathsf{E}_x\Big[\big(\nabla c(X(0))\big)^T \big(\Theta^{-t}\mathcal{S}(t)\big)\nabla\psi(X(-t))\Big]$$
$$= \mathsf{E}_x\Big[\big(\nabla c(X(0))\big)^T \Big(\sum_{t=0}^{\infty} \alpha^t\big(\Theta^{-t}\mathcal{S}(t)\big)\nabla\psi(X(-t))\Big)\Big].$$

The last equality is obtained under Assumption A3 by applying Fubini's theorem, and this completes the proof. □

**Proof of Prop. 3.4**  Lemma 3.5 combined with the stationarity assumption implies that

$$\lim_{T\to\infty} \frac{1}{T}b(t) = \lim_{T\to\infty} \frac{1}{T}\sum_{t=1}^{T} \varphi(t)^T \nabla c(X(t))$$
$$= \mathsf{E}[\varphi(0)^T \nabla c(X(0))] = b$$

Similarly, for each $T \geq 1$ we have

$$M(T) = M(0) + \sum_{t=1}^{T} (\nabla\psi(X(t)))^T \nabla\psi(X(t)),$$

and by the Law of Large Numbers we once again obtain

$$\lim_{T\to\infty} \frac{1}{T}M(T) = M.$$

Combining these results establishes $\theta^* = \lim_{t\to\infty} M^{-1}(t)b(t)$.

Convergence of $\{\bar{c}(t)\}$ is identical, and convergence of $\{\bar{h}_\alpha(t)\}$ also follows from the Law of Large Numbers since we have convergence of $\theta(t)$. ■

### 3.3 Extensions

**Extension to average-cost**

Nowhere in the proof of Prop. 3.4 do we use the assumption that $\alpha < 1$. It is not difficult to establish that under the conditions of the proposition, the $\nabla$-LSTD learning algorithm is convergent when $\alpha = 1$, and the limit solves the quadratic program

$$\theta^* = \arg\min_{\theta} \|h^{\theta} - h\|_{\pi,1}$$

in which $h$ is a solution to Poisson's equation.

**Nonlinear parameterization**

If the parameterized family $\{h_{\alpha}^{\theta}\}$ is nonlinear in $\theta$, then the optimization problem (21) may not be convex. It is possible to construct algorithms to compute a local minimum through stochastic gradient techniques.

The basis should be chosen so that $h^{\theta}$ and $\nabla h_{\alpha}^{\theta}$ are continuous function of $x$, and continuously differentiable in $\theta$. On denoting $\psi_i^{\theta} = $ the partial derivative of $h_{\alpha}^{\theta}$ with respect to $\theta_i$, the first order condition for optimality of (21) is,

$$0 = \mathsf{E}_{\pi}[(\nabla\psi^{\theta^*}(X))^T(\nabla h_{\alpha}^{\theta^*}(X) - \nabla h_{\alpha}(X))],$$

where the $i$th column of the gradient matrix $\nabla\psi^{\theta^*}$ is equal to $\nabla\psi_i^{\theta^*}$.

The inner product $\langle h_{\alpha}, \psi_i^{\theta^*}\rangle_{\pi,1}$ depends on the unknown function $h_{\alpha}$, but this can be transformed into a practical algorithm. For example, a gradient descent like stochastic approximation algorithm is defined through the recursions,

$$\varphi(t) = \alpha\varphi(t-1) + \psi^{\theta(t-1)}(X(t))$$
$$d(t) = \nabla\psi^{\theta(t-1)}(X(t))^T\nabla h_{\alpha}^{\theta(t-1)}(X(t)) - \varphi(t)^T\nabla c(X(t))$$
$$\theta(t) = \theta(t-1) - \gamma_t\nabla\psi^{\theta(t-1)}(X(t))d(t).$$

Stability analysis of these algorithms will be the subject of future research.

Extensions to Markov models in continuous time are contained in the online version of the paper [10].

## 4 Continuous time

To highlight the main ideas, we restrict to a one-dimensional diffusion on $\mathbb{R}$,

$$dX(t) = a(X(t))\,dt + \sigma(X(t))\,dB(t), \tag{37}$$

in which $\boldsymbol{B}$ is standard Brownian motion, and the function $a\colon \mathbb{R} \to \mathbb{R}$ is Lipschitz continuous. For simplicity we also take $\sigma(x) \equiv 1$. Its semigroup is denoted $\{P^t\}$, and its differential generator is defined for $C^2$ functions $f\colon \mathbb{R} \to \mathbb{R}$ via, $\mathcal{D}f = af' + \frac{1}{2}f''$. Assumption A1 is imposed, where again the ergodic limit (11) is equivalent to the existence of a Lyapunov function, defined with respect to the differential generator.

The discounted cost is based on a discount rate $\gamma > 0$, with definition similar to (2):

$$h_{\gamma}(x) = \int_t^{\infty} e^{-\gamma t}\mathsf{E}_x[c(X(t))].$$

The *resolvent kernel* $R_\gamma$ is the Laplace transform,

$$R_\gamma := \int_0^\infty e^{-\gamma t} P^t \, dt, \quad \gamma > 0, \tag{38}$$

so that the value function can be expressed $h_\gamma = R_\gamma c$.

If the value function is $C^2$, then it solves the dynamic programming equation,

$$\mathcal{D}h_\gamma = \gamma h_\gamma - c \tag{39}$$

A representation for the derivative $h'_\gamma$ is obtained in the following subsection, from which we obtain a continuous-time analog of the $\nabla$-LSTD algorithm.

## 4.1 Gradient representation

The dynamic programming equation (39) suggests the inverse formula $R_\gamma = [I_\gamma - \mathcal{D}]^{-1}$. This formula is valid on a suitable domain, and with suitable interpretation [11].

The representation for $h'_\gamma$ makes use of the the generalized resolvent kernel [10–12]: For a measurable function $G \colon \mathbb{R} \to \mathbb{R}$, and measurable functions $f$ in some domain,

$$R_G f(x) := \int_0^\infty \mathsf{E}_x\Big[\exp\Big(-\int_0^t G(X(s)) \, ds\Big) f(X(t))\Big] \, dt. \tag{40}$$

In [11, 12] it is assumed that $G > 0$ everywhere. These conditions are relaxed in [10, 13].

To apply these concepts, differentiate with respect to $x$ each side of (39) to obtain

$$a'h'_\gamma + ah''_\gamma + \tfrac{1}{2}h'''_\gamma = \frac{d}{dx}(\mathcal{D}h_\gamma) = \gamma h'_\gamma - c'$$

Rearranging terms gives $[I_G - \mathcal{D}]h'_\gamma = c'$, with $G = -a' + \gamma$. Provided the inverse exists, we obtain

$$h'_\gamma = [I_G - \mathcal{D}]^{-1} c' = R_G c'. \tag{41}$$

These steps can be justified subject to a growth condition on $c'$, and a Lyapunov drift condition for the diffusion [14].

## 4.2 $\nabla$-LSTD-learning

The goals are unchanged in this continuous time setting: We seek the parameter $\theta^*$ that solves

$$\theta^* = \arg\min_\theta \|h_\gamma^\theta - h_\gamma\|_{\pi,1}^2 \tag{42}$$

The $\nabla$-LSTD-learning algorithm designed to solve this problem is defined by these ODEs:

$$\frac{d}{dt}\varphi(t) = [a'(X(t)) - \gamma]\varphi(t) + \psi'(X(t)) \tag{43a}$$

$$\frac{d}{dt}b(t) = \varphi(t)c'(X(t)) \tag{43b}$$

$$\frac{d}{dt}M(t) = \psi'(X(t))\psi'^T(X(t)) \tag{43c}$$

generating estimates of $\theta^*$ as before via $\theta(t) = M(t)^{-1}b(t)$.

The construction and analysis of this algorithm is based on the characterization of $\theta^*$ for the linearly parameterized approximation. We close this section with an overview of the main ideas.

As in the discrete time case, the objective is quadratic in $\theta$, and as in Prop. 3.2 we obtain $\theta^* = M^{-1}b$ with $M$ defined in (29), and $b = \mathsf{E}_\pi\big[\psi'(X)h'_\gamma(X)\big]$.

The main difference in the continuous time case is the alternate representation for $b$. We again require assumptions to ensure the existence of a steady-state solution to (43a), of the form

$$\varphi(t) = \int_{-\infty}^{t} \exp\Big(-\int_{r}^{t} G(X(s))\,ds\Big)\psi'(X(r))\,dr, \quad t \in \mathbb{R}.$$

with $G(x) = \gamma - a'(x)$, $x \in \mathbb{R}$. This requires a version of Assumption A3 in the continuous time setting; in [14] it is shown that this holds under a Lyapunov drift condition.

When these steps are justified we can conclude that $b = \mathsf{E}_\pi\big[[\varphi(t)c'(X(t))]\big]$, and convergence of the $\nabla$-LSTD algorithm then follows as in the discrete-time setting.

## 5   Simulation Results

This section contains a survey of numerical experiments to illustrate the general theory, and suggest possible extensions of the algorithm.

Common elements in all of our experiments are a linear parameterization for the value function, and the implementation of the $\nabla$-LSTD algorithm. Comparisons with other approaches include the standard LSTD algorithm for discounted cost, and the regenerative LSTD algorithm of [1, 8] for average cost applications where there is regeneration. The standard TD($\lambda$) algorithm was also considered, but in each example the variance was found to be several orders of magnitude greater than alternatives. A matrix gain variant called TD-$K(\lambda)$ is introduced to obtain a better algorithm for comparison. For this linearly parameterized setting, the matrix gain algorithm is essentially equivalent to the LSTD($\lambda$) algorithm of [15].

The asymptotic covariance is used to compare these algorithms:

$$\Sigma = \lim_{t \to \infty} \frac{1}{t}\mathsf{E}[\tilde{\theta}(t)\tilde{\theta}(t)^T] \tag{44}$$

where $\tilde{\theta}(t) = \theta(t) - \theta^*$. Under general conditions, the asymptotic covariance coincides with the covariance in the Central Limit Theorem. It is estimated by observing a histogram following multiple runs of each algorithm.

Two extensions are considered for a specific example: the approximation of relative value function for the speed-scaling model of [16]. First, for this reflected process evolving on $\mathbb{R}_+$, it is shown that the sensitivity process can be defined, subject to conditions on the dynamics near the boundary. Second, the algorithm is tested for an example with *discrete state space*. There is no apparent justification for this approach, but it worked well in the examples considered.

### 5.1   Linear stochastic process

A scalar linear model is ideal for illustrating the difference between $\nabla$-LSTD learning and alternative approaches. The dynamics are given by the recursion

$$X(t+1) = aX(t) + N(t+1)$$

in which $a \in (0, 1)$ and $\boldsymbol{N}$ is Gaussian $\mathcal{N}(0, 1)$.

In all of the numerical results surveyed here, the cost function is defined to be the quadratic, $c(x) = x^2$, $a = 0.7$, and we restrict to the discounted-cost.

The relative value function and the discounted-cost value functions are quadratic in this case, and also symmetric: $h(x) = h(-x)$. Consequently, the function class obtained using the basis $\psi(x) = (1, x^2)^T$ includes the true value function. On taking the gradient, the function space collapses from two dimensions to one, and the growth rate of the functions is reduced from quadratic to linear: $\nabla \psi(x) = (0, 2x)^T$.

For this linear model, the first recursion for the $\nabla$-LSTD algorithm defined in Section 3 becomes

$$\varphi(t) = \alpha a \varphi(t-1) + \nabla \psi(X(t)), \tag{45}$$

while the corresponding equation in the standard LSTD algorithm is
$$\varphi(t) = \alpha \varphi(t-1) + \psi(X(t)). \tag{46}$$

Both these algorithms are consistent. However, two differences suggest that the asymptotic covariance is much smaller when using the $\nabla$-LSTD algorithm. First is the additional discounting factor $a$ appearing in (45), but absent in (46). This is why the LSTD asymptotic covariance grows without bound as $\alpha$ tends to 1. A second advantage of the $\nabla$-LSTD algorithm is that the gradients reduce the growth rate of each function of $x$. In this case, reducing the quadratic growth of $c$ and $\psi$ to the linear growth of their gradients.

Experiments were run for two different discounting factors, $\alpha = 0.9$ and $\alpha = 0.99$. Variance estimates were obtained by conducting $10^3$ independent simulations for each set of parameters tested.
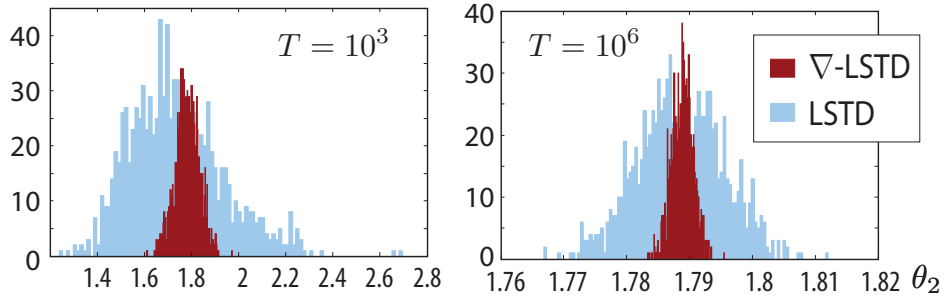


Figure 1: Histogram of $\theta_2(T)$ using both TD-learning and $\nabla$TD-learning, $\alpha = 0.9$.

The optimal parameter is $\theta^* = (16.1, 1.79)^T$ when $\alpha = 0.9$. Figure 1 shows the resulting histograms for $\theta_2(T)$ (the coefficient of $\psi_2(x) = x^2$) for two time horizons, $T = 10^3$ and $10^6$.

It was found that convergence of the $\nabla$-LSTD-learning algorithm is about 10 times faster than the LSTD algorithm. That is, for a given time $T$, the variance of $\theta_2$ estimated using $\nabla$TD-learning algorithm is about the same as that of the TD-learning algorithm which has run for 10 times longer.

The difference in performance of the two algorithms is greater as $\alpha$ is increased. The case $\alpha = 0.99$ is considered next, for which $\theta^* = (192.27, 1.9421)^T$. Figure 2 contains histograms for the same two time horizons.

In conclusion, for this example, the asymptotic covariance of the $\nabla$-LSTD algorithm is bounded uniformly over $0 < \alpha < 1$, and it can also be used to estimate the solution to Poisson's equation.

## 5.2 Dynamic speed scaling

Dynamic speed scaling is a popular approach to power management in computer system design. The goal is to control the processing speed so as to optimally balance energy and delay costs; this
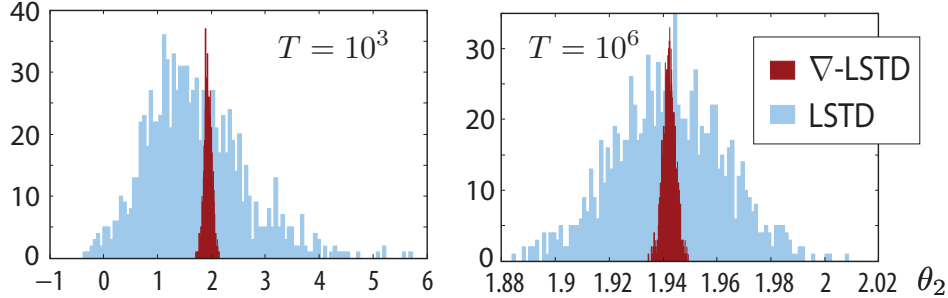
13

Figure 2: Histogram of $\theta_2(T)$ using both TD-learning and $\nabla$TD-learning, $\alpha = 0.99$.

can be done by reducing (increasing) the processor speed at times when the workload is small (large). For the purposes of this paper, speed scaling is a simple stochastic control problem: a single server queue with a controllable service rate.

A regenerative form of LSTD learning was applied in [16] for this example to approximate the solution to the average-cost optimality equation. Approximate policy iteration algorithm was implemented, in which the LSTD algorithm provided an approximate relative value function at each iteration of the algorithm.

The discrete time MDP (Markov Decision Process) model is described as follows: At each time $t$, the state $X(t)$ is interpreted either as the queue length, or the workload in the system; $N(t) \geq 0$ denotes the number of job arrivals, and $U(t)$ the service completion at time $t$. This is subject to the constraint $0 \leq U(t) \leq X(t)$. The evolution equation is the controlled random walk:

$$X(t+1) = X(t) - U(t) + N(t+1), \quad t \geq 0. \tag{47}$$

Under the assumption that $\boldsymbol{N}$ is i.i.d., and $\boldsymbol{U}$ is obtained using a state feedback policy $U(t) = \mathrm{f}(X(t))$, the controlled model is a Markov chain of the form (1).

In the experiments that follow we focus exclusively on the average-cost setting, with $c(x,u) = x + u^2/2$, and

$$\mathrm{f}(x) = \min(x, 1 + \varepsilon\sqrt{x}), \tag{48}$$

with $\varepsilon > 0$. This is similar in form to the optimal average cost policy calculated in [16]. It is shown in [16] that the value function is approximated by the function $h^\theta(x) = \theta^T \psi(x)$ for some $\theta \in \mathbb{R}_+^2$, with $\psi(x) = (x^{3/2}, x)^T$. As in the linear example, the gradient $\nabla\psi(x) = (\frac{3}{2}x^{1/2}, 1)^T$ has much slower growth as a function of $x$.

Implementation of the $\nabla$-LSTD algorithm requires attention to the boundary of the state space. The sensitivity process $\{\mathcal{S}(t)\}$ as defined in (15) requires that the state space be open, and that the dynamics are smooth. Both of these assumptions are violated in this example. However, we do have a representation for the right derivative, which evolves according to the recursive equation,

$$\mathcal{S}(t+1) = \mathcal{A}^T(t+1)\mathcal{S}(t) = [1 - \tfrac{d^+}{dx}\mathrm{f}(X(t))]\mathcal{S}(t) \tag{49}$$

We begin with the case in which the marginal of $\boldsymbol{N}$ is exponential. In this case the right derivatives and ordinary derivatives coincide a.s..

**Exponential arrivals**

The marginal distribution of $\boldsymbol{N}$ was taken to be the unit-mean exponential.

14

When $\boldsymbol{X}$ evolves on $\mathbb{R}_+$ with policy f as defined above, the form of the $\nabla$-LSTD algorithm is unchanged from the definition given in Section 3. The recursion for $\varphi$ in (22) is implemented based on (49):

$$\mathcal{A}(t+1) = \mathbf{1}\{X(t) > \bar{\varepsilon}\}\big[1 - \tfrac{1}{2}\varepsilon X(t)^{-1/2}\big] \tag{50}$$

where $\bar{\varepsilon} = \frac{1}{2}(\varepsilon + \sqrt{\varepsilon^2 + 4})$.

The regenerative LSTD algorithm used in [16] is not directly applicable in this example. Various forms of the TD($\lambda$) algorithm were tested, but all appeared to have infinite asymptotic variance. The introduction of a matrix gain resulted in improved performance. The examples that follow compare the $\nabla$-LSTD algorithm with the best results we were able to obtain using other methods.

The matrix gain algorithm will be called TD-$K(\lambda)$. It is identical to the standard algorithm, except for the introduction of a matrix gain sequence $\{K_t\}$ in the following:

$$
\begin{aligned}
\theta(t+1) &= \theta(t) + \gamma_{t+1} K_t z(t) d(t+1)\\
d(t+1) &= \tilde{c}(t) + \big[\psi(X(t+1)) - \psi(X(t))\big]^T \theta(t)\\
\bar{c}(t+1) &= \bar{c}(t) + \gamma_{t+1}\big[-\bar{c}(t) + c(X(t+1))\big]\\
z(t+1) &= \lambda z(t) + \psi(X(t+1)).
\end{aligned}
\tag{51}
$$

$\tilde{c}(t) = c(X(t)) - \bar{c}(t)$. To optimize a constant gain $K_t \equiv K$ over all $\ell \times \ell$ matrices, the solution is obtained by considering the associated ODE, $\dot{\vartheta} = KA(\vartheta - \theta^*)$. The choice $K = -A^{-1}$ is known to be optimal. In this example we have,

$$A = \mathsf{E}\big[z(t)(\psi(X(t+1)) - \psi(X(t)))^T\big]$$

where the expectation is taken in steady state, with $z(t) = \sum_{k=0}^{\infty} \lambda^k \psi(X(t-k))$. This was estimated using

$$A_{t+1} = A_t + \gamma_{t+1}\big[-A_t + z(t)(\psi(X(t+1)) - \psi(X(t)))^T\big]$$

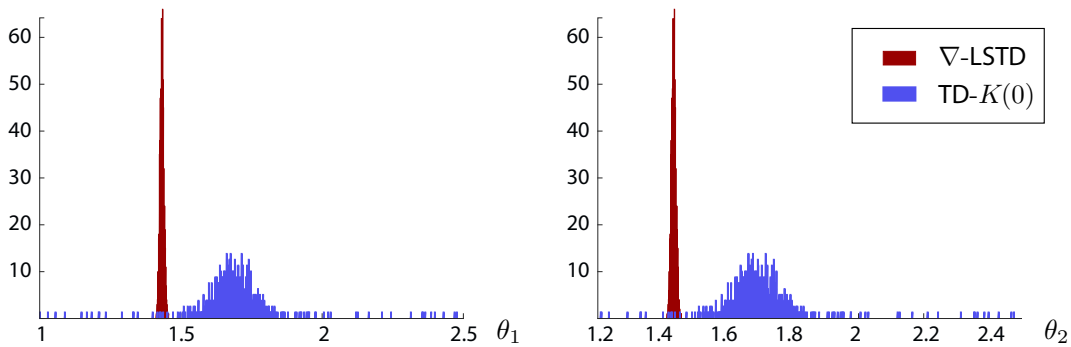and $K_{t+1} = -A_{t+1}^{-1}$ (i.e., Stochastic Newton Raphson).



Figure 3: Histogram of the parameters estimated using LSTD and $\nabla$-LSTD, for the speed scaling model

Figure 3 shows the histogram of the two parameters, estimated using both $\nabla$-LSTD-learning and TD-$K(0)$-learning, run for a duration of $T = 10^5$ time steps. The stationary policy used is as defined in (48) with $\varepsilon = 0.5$. Note that here again, the variance reduction obtained using the $\nabla$-LSTD-learning algorithm is remarkable.

## $\nabla$-LSTD and regenerative LSTD

In [16], the authors consider a discrete state space, with $N(t)$ geometrically distributed on an integer lattice $\{0, \Delta, 2\Delta, \dots\}$. In this case, the theory developed for the $\nabla$-LSTD-algorithm does not fit

the model since we have no convenient representation of a sensitivity process. Nevertheless, the algorithm can be run by replacing gradients with ratios of differences. In particular, in implementing the algorithm we substitute the definition (50) with $\mathcal{A}(t) = 1 - [\mathrm{f}(X(t) + \Delta) - \mathrm{f}(X(t))]/\Delta$, and $\nabla c$ was approximated similarly.

The geometric distribution gives $\mathsf{P}(N(t) = n\Delta) = (1 - p_A)^n p_A$; the values $p_A = 0.04$ and $\Delta = 1/24$ were chosen, so that $\mathsf{E}[N(t)] = 1$.

The sequence of steps followed in the regenerative LSTD-learning algorithm are similar to (20):

$$
\begin{aligned}
\varphi(t) &= \mathbf{1}\{X(t-1) \neq 0\}\varphi(t-1) + \tilde{\psi}(t) \\
b(t) &= (1 - \gamma_t)b(t-1) + \gamma_t \tilde{c}(t)\varphi(t) \\
M(t) &= (1 - \gamma_t)M(t-1) + \gamma_t \tilde{\psi}(t)\tilde{\psi}^T(t) \\
\bar{c}(t) &= (1 - \gamma_t)\bar{c}(t-1) + \gamma_t c(X(t)),
\end{aligned}
\tag{52}
$$

where $\tilde{c}(t) = c(X(t)) - \bar{c}(t)$, $\tilde{\psi}(t) := \psi(X(t)) - \eta_\psi(t)$, and

$$
\eta_\psi(t) = (1 - \gamma_t)\eta_\psi(t-1) + \gamma_t \psi(X(t)) \,.
$$

The parameter at time $t$ is obtained as $\theta(t) = M^{-1}(t)b(t)$.

We replace $\psi$ with $\tilde{\psi}$ in (52) to restrict the growth rate of the eligibility vector $\varphi(t)$, which in turn reduces the variance of the estimates $\theta(t)$. This is justified because $h_a^\theta = \theta^T \tilde{\psi}$ differs from $h_b^\theta = \theta^T \psi$ by only a constant term.
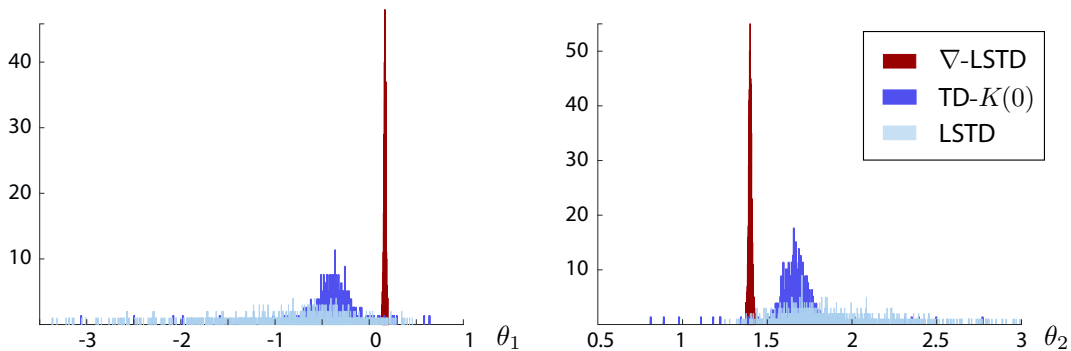


Figure 4: Histogram of the parameters estimated using the three algorithms.

For comparison purposes, we also implemented the TD-$K(0)$ algorithm, defined in (51). Figure 4 shows the histogram of $\theta(T)$ obtained using $\nabla$-LSTD-learning, regenerative LSTD-learning, and TD-$K(0)$ learning algorithms, with $T = 10^5$. The variance using $\nabla$-LSTD-learning algorithm is extremely small compared to the other two. Also, TD-$K(0)$ had the largest outliers.

It is also noticeable in Figure 4 that there is a difference in the values to which the algorithms have converged. To investigate the quality of the estimates through a different lens, the *Bellman error* was computed for each algorithm:

$$
\mathcal{E}_B(x) = [P - I]h(x) + \tilde{c}(x),
$$

where $P$ of course depends on the policy f, and $h = \bar{\theta}^T \psi$, where $\bar{\theta}$ is the mean of the $10^3$ parameter estimates obtained for each of the different algorithms.

Figure 5 shows plots of the Bellman error observed, for each of the three algorithms, for typical values of $\theta(T)$, with $T = 10^3$, $10^4$ and $10^5$. In each case the stationary policy (48) was used, with $\varepsilon = 0.5$.

16

The limit of the Bellman error is the same in each experiment. In the case of the $\nabla$-LSTD algorithm, the Bellman error is unchanged for $T \geq 10^3$. Achieving similar performance using either of the other algorithms required about $10^5$ samples.
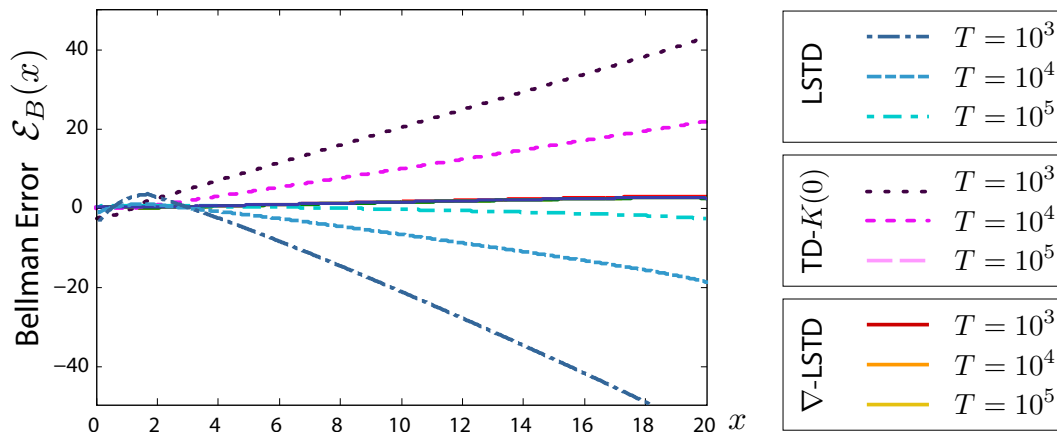


Figure 5: Bellman error corresponding to the estimates of $h$ based on the three algorithms. The convergence time for $\nabla$-LSTD algorithm is two orders of magnitude faster than the two other algorithms.

# 6 Conclusions

The new gradient based TD-learning algorithms for value function approximation introduced here show remarkable variance reduction in the examples considered. This is explained by the reduction in magnitude of the functions used as inputs to the algorithm, and also from the additional "discounting" that is inherent in the new algorithms.

The most interesting open problem is why the algorithm is so effective even in a discrete state-space setting in which there is no theory to justify its application.

# References

[1] S. P. Meyn, *Control Techniques for Complex Networks*. Cambridge University Press, 2007, pre-publication edition available online.

[2] D. Bertsekas and S. Shreve, *Stochastic Optimal Control: The Discrete-Time Case*. Athena Scientific, 1996.

[3] S. Asmussen and P. W. Glynn, *Stochastic Simulation: Algorithms and Analysis*, ser. Stochastic Modelling and Applied Probability. New York: Springer-Verlag, 2007, vol. 57.

[4] T. Yang, P. Mehta, and S. Meyn, "Feedback particle filter," *IEEE Trans. Automat. Control*, vol. 58, no. 10, pp. 2465–2480, Oct 2013.

[5] R. Sutton and A. Barto, *Reinforcement Learning: An Introduction*. Cambridge, MA: MIT Press. On-line edition at http://www.cs.ualberta.ca/~sutton/book/the-book.html, 1998.

[6] D. Bertsekas and J. N. Tsitsiklis, *Neuro-Dynamic Programming*. Cambridge, Mass: Atena Scientific, 1996.

[7] V. S. Borkar, *Stochastic Approximation: A Dynamical Systems Viewpoint.* Delhi, India and Cambridge, UK: Hindustan Book Agency and Cambridge University Press (jointly), 2008.

[8] D. Huang, W. Chen, P. Mehta, S. Meyn, and A. Surana, "Feature selection for neuro-dynamic programming," in *Reinforcement Learning and Approximate Dynamic Programming for Feedback Control*, F. Lewis, Ed. Wiley, 2011.

[9] S. P. Meyn and R. L. Tweedie, *Markov chains and stochastic stability*, 2nd ed. Cambridge: Cambridge University Press, 2009, published in the Cambridge Mathematical Library. 1993 edition online.

[10] A. Devraj, I. Kontoyiannis, and S. P. Meyn, "Exponential ergodicity and Lyapunov exponents Part I: Markov chains in discrete time," In preparation, March 2016.

[11] S. P. Meyn and R. L. Tweedie, "Generalized resolvents and Harris recurrence of Markov processes," *Contemporary Mathematics*, vol. 149, pp. 227–250, 1993.

[12] J. Neveu, "Potentiel Markovien récurrent des chaînes de Harris," *Ann. Inst. Fourier, Grenoble*, vol. 22, pp. 7–130, 1972.

[13] I. Kontoyiannis and S. P. Meyn, "Spectral theory and limit theorems for geometrically ergodic Markov processes," *Ann. Appl. Probab.*, vol. 13, pp. 304–362, 2003, presented at the INFORMS Applied Probability Conference, NYC, July, 2001.

[14] A. Devraj, I. Kontoyiannis, and S. P. Meyn, "Exponential ergodicity and Lyapunov exponents Part II: Markovian diffusions," In preparation., March 2016.

[15] J. A. Boyan, "Technical update: Least-squares temporal difference learning," *Mach. Learn.*, vol. 49, no. 2-3, pp. 233–246, 2002.

[16] W. Chen, D. Huang, A. A. Kulkarni, J. Unnikrishnan, Q. Zhu, P. Mehta, S. Meyn, and A. Wierman, "Approximate dynamic programming using fluid and diffusion approximations with applications to power management," in *Proc. of the 48th IEEE Conf. on Dec. and Control; held jointly with the 2009 28th Chinese Control Conference*, 2009, pp. 3575–3580.