

# Fast LSTD Using Stochastic Approximation: Finite Time Analysis and Application to Traffic Control

L.A. Prashanth<sup>1</sup>, Nathaniel Korda<sup>2</sup>, and Rémi Munos<sup>1</sup>

<sup>1</sup> INRIA Lille - Nord Europe, Team SequeL, France  
{prashanth.la, remi.munos}@inria.fr

<sup>2</sup> Oxford University, United Kingdom  
nathaniel.korda@eng.ox.ac.uk

**Abstract.** We propose a stochastic approximation based method with randomisation of samples for policy evaluation using the least squares temporal difference (LSTD) algorithm. Our method results in an  $O(d)$  improvement in complexity in comparison to regular LSTD, where  $d$  is the dimension of the data. We provide convergence rate results for our proposed method, both in high probability and in expectation. Moreover, we also establish that using our scheme in place of LSTD does not impact the rate of convergence of the approximate value function to the true value function. This result coupled with the low complexity of our method makes it attractive for implementation in *big data* settings, where  $d$  is large. Further, we also analyse a similar low-complexity alternative for least squares regression and provide finite-time bounds there. We demonstrate the practicality of our method for LSTD empirically by combining it with the LSPI algorithm in a traffic signal control application.

## 1 Introduction

Several machine learning problems involve solving a linear system of equations from a given set of training data. In this paper we consider the problem of policy evaluation in reinforcement learning (RL) using the method of temporal differences (TD). Given a fixed training data set, one popular temporal difference algorithm for policy evaluation is LSTD [4]. However, LSTD is computationally expensive as it requires  $O(d^2)$  computations. We propose a stochastic approximation (SA) based algorithm that draws data samples from a uniform distribution on the training set. From the finite time analyses that we provide, we observe our algorithm converges at the optimal rate, in high probability as well as in expectation. Moreover, using our scheme in place of LSTD does not impact the rate of convergence of the approximate value function to the true value function. This finding coupled with the significant decrease in the computational cost of our algorithm, makes it appealing in the canonical *big data* settings.

The problem considered here is to estimate the value function  $V^\pi$  of a given policy  $\pi$ . Temporal difference (TD) methods are well-known in this context, and they are known to converge to the fixed point  $V^\pi = \mathcal{T}^\pi(V^\pi)$ , where  $\mathcal{T}^\pi$  is the Bellman operator (see Section 3.1 for a precise definition). A popular approach to overcome the curse of dimensionality associated with large state spaces is to parameterize the value function using a linear function approximation architecture. For every  $s$  in the state space  $\mathcal{S}$ ,

we approximate  $V^\pi(s) \approx \theta^\top \phi(s)$ , where  $\phi(\cdot)$  is a  $d$ -dimensional feature vector with  $d \ll |\mathcal{S}|$ , and  $\theta$  is a tunable parameter. The function approximation variant of TD [23] is known to converge to the fixed point of  $\Phi\theta = \Pi\mathcal{T}^\pi(\Phi\theta)$ , where  $\Pi$  is the orthogonal projection onto the space within which we approximate the value function, and  $\Phi$  is the feature matrix that characterises this space.

LSTD estimates the fixed point of  $\Pi\mathcal{T}^\pi$  using empirical data  $\mathcal{D} := \{(s_i, r_i, s'_i), i = 1, \dots, T\}$  obtained by simulating the Markov decision process (MDP) with the underlying policy  $\pi$ . For every  $i = 1, \dots, T$ , the 3-tuple  $(s_i, r_i, s'_i)$  corresponds to a transition from state  $s_i$  to  $s'_i$  under action  $\pi(s_i)$  and the resulting reward is denoted by  $r_i$ . The LSTD estimate is given as the solution to  $\hat{\theta}_T = \bar{A}_T^{-1}\bar{b}_T$ , where  $\bar{A}_T = \frac{1}{T} \sum_{i=1}^T \phi(s_i)(\phi(s_i) - \beta\phi(s'_i))^\top$ , and  $\bar{b}_T = \frac{1}{T} \sum_{i=1}^T r_i\phi(s_i)$ .

Computing the inverse of the matrix  $\bar{A}_T$  is computationally expensive, especially when  $d$  is large. Indeed, assuming that the features  $\phi(s_i)$  evolve in a compact subset of  $\mathbb{R}^d$ , the complexity of the above approach is  $O(d^2T)$ , where  $\bar{A}_T^{-1}$  is computed iteratively using the Sherman-Morrison lemma. On the other hand, if we employ the Strassen algorithm or the Coppersmith-Winograd algorithm for computing  $\bar{A}_T^{-1}$ , the complexity is of the order  $O(d^{2.807})$  and  $O(d^{2.375})$ , respectively, in addition to  $O(d^2T)$  complexity for computing  $\bar{A}_T$ .

A common trick, in practice, to alleviate this problem in high dimensions, is to replace the inversion of the  $\bar{A}_T$  matrix by an iterative procedure that performs a fixed point iteration. From a theoretical standpoint, this comes under the purview of stochastic approximation (SA), and one requires that the samples be chosen randomly to ensure convergence. In this paper, we analyse such an SA based scheme and show that it converges to the LSTD solution. The advantage is that the SA based scheme incurs lower computational cost in comparison to the approaches mentioned above. We also analyse a similar low-complexity alternative for the classic least squares parameter estimation problem.

We provide convergence rate results for our proposed method, both in high probability and in expectation. In particular, we show that, with probability  $1 - \delta$ , the SA based scheme constructs an  $\epsilon$ -approximation of the corresponding LSTD solution with  $O(d \ln(1/\delta)/\epsilon^2)$  complexity, irrespective of the number of samples  $T$ . Moreover, we also establish that using the SA based scheme in place of LSTD does not impact the rate of convergence of the approximate value function to the true value function (see Theorem 2).

The rate results coupled with the low complexity of our scheme make it more amenable to practical implementation in the canonical *big data* settings, where both  $d$  and  $T$  are large. Further, we provide explicit constants in the high probability bounds and we believe this opens several avenues for the use of SA based low complexity alternatives in higher level decision making procedures, for instance, least squares policy iteration (LSPI) [11] and linear bandit [5] algorithms. We demonstrate the practicality of our solution scheme for LSTD empirically by using it as a subroutine in the LSPI algorithm for adaptive traffic signal control<sup>1</sup>. In particular, for the experiments we

<sup>1</sup> See [16] for another set of experiments that combines the SA based low-complexity variant for least squares regression with the LinUCB algorithm for contextual bandits, using the large scale news recommendation dataset from Yahoo [24].

employ step-sizes that were used to derive the finite-time bounds (see Corollary 1). We demonstrate that this choice results in rapid convergence of our SA based scheme in the experiments and also that the performance of the SA variant of LSPI is comparable to that of LSPI.

The rest of the paper is organized as follows: In Section 2, we review relevant previous works and relevant literature. In Section 3 we present the fast LSTD algorithm based on stochastic approximation and in Section 4 we provide the non-asymptotic bounds for this algorithm. In Section 5, we outline the variants of our algorithm to incorporate regularization and iterate averaging, while in Section 7, we provide extensions to solve the problem of least squares regression. Next, in Section 6, we provide outlines for the proof and derivation of rates. In Section 8, we provide experiments on a traffic signal control application. Finally, in Section 9 we provide the concluding remarks.

## 2 Literature Review

Our algorithms are based on the well-known stochastic approximation technique, originally proposed for finding zeroes of a nonlinear function in [17]. The reader is referred to [10] for a textbook introduction to SA. Iterate averaging is a standard approach to accelerate the convergence of SA schemes and was proposed independently in [18] and [13]. Non asymptotic bounds for Robbins Monro schemes have been provided in [7] and extended to incorporate iterate averaging in [6].

In the context of the problem of prediction in RL, temporal difference (TD) learning is a well-known algorithm. See [3,20] for a textbook introduction and [23] for an asymptotic analysis. LSTD [4] is a popular batch algorithm that converges asymptotically to the TD solution. Finite time analysis of LSTD is provided in [12] and we extend it to the case when LSTD solution is replaced by a SA iterate.

A popular line of research in RL is on improving the complexity of TD-like algorithms (cf. GTD [21], GTD2 [22], iLSTD [8] and the references therein). The popular Computer Go with dimension  $d = 10^6$  [19] and several practical applications (e.g. transportation, networks) involve high-feature dimensions. Moreover, considering that linear function approximation is effective with a large number of features, our  $O(d)$  improvement in complexity of LSTD by employing SA is meaningful.

In comparison to previous work, we would like to point out that there is no finite time analysis of GTD-type algorithms. While iLSTD is an efficient approximation to LSTD, analysis in [8] requires that the feature matrix be sparse. In contrast, we provide finite-time bounds and do not make any sparsity assumption. To the best of our knowledge, efficient SA algorithms that approximate LSTD without impacting its rate of convergence to true value function, have not been proposed before in the literature. The high probability bounds that we derive for the SA based scheme do not directly follow from earlier work on LSTD algorithms. Further, unlike [7], we provide explicit constants in the bounds that we derive (see Corollary 1) and we employ these in our experiments as well.

Stochastic gradient descent (SGD) is a well-known method for optimising a function given only noisy observations. In the context of machine learning, finite time analysis of such methods have been provided in [1]. While the bounds in [1] are given in expectation, many machine learning applications require high probability bounds, which

we provide for our case. Regret bounds for online SGD techniques have been given in [25,9]: the gradient descent algorithm in [25] is in the setting of optimising the average of convex loss functions whose gradients are available, while that in [9] is for strongly convex loss functions.

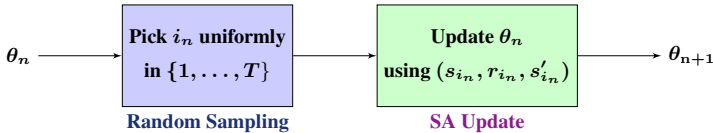
In comparison to previous work w.r.t. least squares regression, we highlight the following differences: **(i)** Earlier works on least squares regression (cf. [9]) require the knowledge of the strong convexity constant in deciding the step-size, while we average the iterates to get rid of this dependency. **(ii)** Our analysis is much simpler (since we work directly with least squares problems) and we make all the constants explicit for the problems considered.

### 3 Fast LSTD Using Stochastic Approximation (fLSTD-SA)

We propose here a stochastic approximation variant of the least squares temporal difference (LSTD) algorithm, whose iterates converge to the same fixed point as the regular LSTD algorithm, while incurring much smaller overall computational cost.

The algorithm, which we call fast LSTD through Stochastic Approximation (fLSTD-SA), is a simple stochastic approximation scheme with randomised samples. The results that we present establish that fLSTD-SA computes an  $\epsilon$ -approximation to the LSTD solution  $\hat{\theta}_T$  with probability  $1 - \delta$ , while incurring a complexity of the order  $O(d \ln(1/\delta)/\epsilon^2)$ , irrespective of the number of samples  $T$ . In turn, this enables us to give a performance bound for the approximate value function computed by fLSTD-SA. A schema of fLSTD-SA is given in Figure 1.

Although our analysis for fLSTD-SA depends on a strong convexity assumption that may not hold in all situations, we present also a variant of fLSTD-SA employing iterate averaging for which error bounds can be given without resorting to a strong convexity assumption.



**Fig. 1.** Overall flow of the fLSTD-SA algorithm

#### 3.1 Background for LSTD

Consider an MDP with (finite) state space  $\mathcal{S}$ , (finite) action space  $\mathcal{A}$  and transition probabilities  $p(s, a, s')$ ,  $s, s' \in \mathcal{S}$ ,  $a \in \mathcal{A}$ . For a given stationary policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ , the value function  $V^\pi$  is defined by

$$V^\pi(s) := E \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, \pi(s_t)) \mid s_0 = s \right], \quad (1)$$

where  $s_t$  denotes the state of the MDP at time  $t$ ,  $\beta \in (0, 1)$  is the discount factor, and  $r(s, a)$  denotes the instantaneous rewards obtained in state  $s$  with action  $a$ . The value function  $V^\pi$  can be expressed as the fixed point of the Bellman operator  $\mathcal{T}^\pi$  defined by

$$\mathcal{T}^\pi(V)(s) := r(s, \pi(s)) + \beta \sum_{s'} p(s, \pi(s), s') V(s'), \quad (2)$$

When the cardinality of  $\mathcal{S}$  is huge and in the absence of knowledge of the transition dynamics, a popular approach is to parameterize the value function using a linear function approximation architecture, i.e., for every  $s \in \mathcal{S}$ , we approximate  $V^\pi(s) \approx \phi(s)^\top \theta$ , where  $\phi(s)$  is a  $d$ -dimensional feature vector with  $d \ll |\mathcal{S}|$ , and  $\theta$  is a tunable parameter. The well-known TD learning algorithm [3] attempts to find the fixed point of the operator  $\Pi \mathcal{T}^\pi$  given by

$$\Phi \theta = \Pi \mathcal{T}^\pi(\Phi \theta), \quad (3)$$

where  $\mathcal{B} = \{\Phi \theta \mid \theta \in \mathbb{R}^d\}$  is the space within which we want to approximate the value function  $V^\pi$ ,  $\Pi$  is the orthogonal projection onto  $\mathcal{B}$ , and  $\Phi$  is the feature matrix with rows  $\phi(s)^\top, \forall s \in \mathcal{S}$  denoting the features corresponding to state  $s \in \mathcal{S}$ . Let  $\theta^*$  denote the solution to (3),  $P$  the transition probability matrix with components  $p(s, \pi(s), s')$  and  $\Psi$  the stationary distribution (assuming it exists) of the Markov chain for the underlying policy  $\pi$ . Then,  $\theta^*$  can be written as the solution to the following system of equations (cf. [2, Section 6.3])

$$A \theta^* = b, \text{ where } A = \Phi^\top \Psi (I - \beta P) \Phi \text{ and } b = \Phi^\top \Psi r. \quad (4)$$

The LSTD approach is to approximate  $A$  and  $b$  using  $T$  samples  $\{(s_i, r_i, s'_i), i = 1, \dots, T\}$  obtained by simulating the MDP with the underlying policy  $\pi$ .

An approximate solution to (4) is constructed as follows:

$$\hat{\theta}_T = \bar{A}_T^{-1} \bar{b}_T \quad (5)$$

where  $\bar{A}_T = T^{-1} \sum_{i=1}^T \phi(s_i)(\phi(s_i) - \beta \phi(s'_i))^\top$ , and  $\bar{b}_T = T^{-1} \sum_{i=1}^T r_i \phi(s_i)$ . Here  $\phi(s_i)$  is a  $d$ -dimensional feature vector corresponding to state  $s_i$ , for all  $i = 1, \dots, T$ . By invoking the strong law of large numbers, one can show that  $\bar{A}_T \rightarrow A$  and  $\bar{b}_T \rightarrow b$  as the number of samples  $T$  tends to infinity.

### 3.2 Update Rule for Flstd-SA

Starting with an arbitrary  $\theta_0$ , we update the parameter  $\theta_n$  as follows:

$$\theta_n = \theta_{n-1} + \gamma_n (r_{i_n} + \beta \theta_{n-1}^\top \phi(s'_{i_n}) - \theta_{n-1}^\top \phi(s_{i_n})) \phi(s_{i_n}), \quad (6)$$

where each  $i_n$  is chosen uniformly randomly from the set  $\{1, \dots, T\}$ . In other words, we pick a sample with uniform probability  $1/T$  from the set  $\mathcal{D} := \{(s_i, r_i, s'_i), i = 1, \dots, T\}$  and use it to perform a fixed point iteration in (6). The quantities  $\gamma_n$  above are *step sizes* that are chosen in advance and satisfy standard stochastic approximation conditions (see (A1) below). Notice that the above update is the usual TD update, except that the samples are drawn uniformly randomly from the sample set  $\mathcal{D}$ .

## 4 Main Results

### 4.1 Error Bounds

We make the following assumptions for the analysis fLSTD-SA:

(A1) The step sizes  $\gamma_n$  satisfy  $\sum_n \gamma_n = \infty$ , and  $\sum_n \gamma_n^2 < \infty$ .

(A2) Bounded features:  $\|\phi(s_i)\|_2 \leq 1$ , for  $i = 1, \dots, T$ .

(A3) Bounded rewards:  $|r_i| \leq R_{\max} < \infty$  for  $i = 1, \dots, T$  and bounded linear space, i.e.,  $-V_{\max} \leq \Phi\theta \leq V_{\max} < \infty$ .

(A4) Strong Convexity: Writing  $\Phi_T \triangleq (\phi(s_1)^\top; \dots; \phi(s_T)^\top)$ , the covariance matrix  $\frac{1}{T}\Phi_T^\top\Phi_T$  is positive definite and its smallest (positive) eigenvalue is at least  $\mu$ .

By working in a bounded linear space along with bounded rewards and features, along with step sizes that satisfy standard stochastic approximation conditions, we ensure that the parameter  $\theta$  remains stable, and hence that (6) converges.

To obtain high probability bounds on the error we consider separately the deviation of  $z_n$  from its mean (see (7) in Theorem 1), and the size of its mean itself (see (8) in Theorem 1). In this way the first quantity can be directly decomposed as a sum of martingale differences, and then a standard martingale concentration argument applied, while the second quantity can be analyzed by directly unrolling iteration (6) (a proof outline is provided in Section 6, while the detailed proofs are available in [16]).

**Theorem 1.** *Under (A1)-(A4), we have  $\forall \epsilon > 0$ ,*

$$P(\|\theta_n - \hat{\theta}_T\|_2 - \mathbb{E}\|\theta_n - \hat{\theta}_T\|_2 \geq \epsilon) \leq \exp\left(-\epsilon^2 / (2 \sum_{i=1}^n L_i^2)\right), \quad (7)$$

$$\begin{aligned} \mathbb{E}\|\theta_n - \hat{\theta}_T\|_2 &\leq \underbrace{\exp(-(1-\beta)\mu\Gamma_n)}_{\text{initial error}} \|\theta_0 - \hat{\theta}_T\|_2 \\ &\quad + \underbrace{\left(\sum_{k=1}^{n-1} H_\beta^2 \gamma_{k+1}^2 \exp(-2(1-\beta)\mu(\Gamma_n - \Gamma_{k+1}))\right)^{\frac{1}{2}}}_{\text{sampling error}}, \quad (8) \end{aligned}$$

where  $L_i := \gamma_i \prod_{j=i}^{n-1} (1 - 2\gamma_{j+1}\mu((1-\beta) - \beta(2-\beta)\gamma_{j+1}))^{1/2}$ ,  $\Gamma_n := \sum_{i=1}^n \gamma_i$  and  $H_\beta^2 := R_{\max}(R_{\max} + 2) + (1 + \beta)^2 V_{\max}^2$ .

The initial error depends on the initial point  $\theta_0$  of the algorithm. The sampling error arises out of a martingale difference sequence that depends on the random deviation of the stochastic update from the standard fixed point iteration, and is the dominant term in (8). Under a suitable choice of step-sizes (see Corollary 1), it can be shown that the initial error is forgotten faster than the sampling error.

The above theorem assumes no specific form for the step-sizes  $\gamma_n$ . Specifying the step-size sequence, we can merge the two claims above to deduce the following bounds on the approximation error  $z_n$  with explicit constants:

**Corollary 1 (Error Bound for iterates of fLSTD-SA).** *Under (A2)-(A4), choosing  $\gamma_n = \frac{(1-\beta)c}{2(c+n)}$  and  $c$  such that  $(1-\beta)^2\mu c \in (1.33, 2)$ , we have, for any  $\delta > 0$ ,*

$$\mathbb{E}\|\theta_n - \hat{\theta}_T\|_2 \leq \frac{K_1(n)}{\sqrt{n+c}} \text{ and } P\left(\|\theta_n - \hat{\theta}_T\|_2 \leq \frac{K_2(n)}{\sqrt{n+c}}\right) \geq 1 - \delta, \quad (9)$$

where  $K_1(n)$  and  $K_2(n)$  are functions of order  $O(1)$ , defined by:

$$K_1(n) = \frac{\sqrt{c}\|\theta_0 - \hat{\theta}_T\|_2}{n^{((1-\beta)^2\mu c - 1)/2}} + \frac{(1-\beta)cH_\beta}{2}, \quad K_2(n) = \frac{(1-\beta)c\sqrt{\log \delta^{-1}}}{2\sqrt{\left(\frac{4}{3}(1-\beta)^2\mu c - 1\right)}} + K_1(n).$$

*Remark 1.* We note that setting  $c$  such that  $(1-\beta)^2\mu c = \eta \in (1.33, 2)$  we can rewrite the constants in Corollary 1 as:

$$K_1(n) = \frac{\|\theta_0 - \hat{\theta}_T\|_2}{(1-\beta)\sqrt{\mu n^{(\eta-1)}}} + \frac{H_\beta}{2(1-\beta)\mu}, \quad K_2(n) = \frac{\sqrt{\log \delta^{-1}}}{2(1-\beta)\mu\sqrt{\left(\frac{4}{3}\eta - 1\right)}} + K_1(n).$$

So both the bounds in expectation and high probability have a linear dependence on the inverse of  $(1-\beta)\mu$ .

## 4.2 Performance Bound

Let  $\tilde{v}_T := \Phi\theta_T$  denote the approximate value function obtained from  $T$  steps of fLSTD-SA, and let  $v$  denote the true value function, evaluated at the states  $s_1, \dots, s_T$ . Then the following lower bound on the performance of  $\tilde{v}_T$  can be deduced from Corollary 1 in conjunction with Theorem 1 of [12]:

**Theorem 2.** *Under conditions of Corollary 1, for any  $\delta > 0$ , with probability  $1 - \delta$ ,*

$$\|v - \tilde{v}_T\|_T \leq \underbrace{\frac{\|v - \Pi v\|_T}{\sqrt{1-\beta^2}}}_{\text{residual error}} + \underbrace{O\left(\sqrt{\frac{d}{(1-\beta)^2\mu T}}\right)}_{\text{estimation error}} + \underbrace{O\left(\sqrt{\frac{1}{(1-\beta)\mu T} \ln \frac{1}{\delta}}\right)}_{\text{approximation error}},$$

where  $\|f\|_T^2 := T^{-1} \sum_{i=1}^T f(s_i)^2$ , for any function  $f$ .

The residual and estimation errors (first and second terms in the RHS above) are artifacts of function approximation and least squares methods, respectively. The third term, of order  $O(1/\sqrt{T})$ , is a consequence of using fLSTD-SA in place of the LSTD. From the above theorem, we observe that using our scheme in place of LSTD does not impact the rate of convergence of the approximate value function  $\tilde{v}_T$  to the true value function  $v$ . This finding coupled with the fact that our scheme is of low complexity makes it attractive for implementation in *big data* settings, where the feature dimension  $d$  is large.

## 5 Variants

To obtain the best performance from fLSTD-SA we need to know the value of  $\mu$ . However with minor adjustments to the analysis we can provide two variants of fLSTD-SA for which it is not necessary to know the value of  $\mu$  to obtain the (optimal) approximation error of order  $O(n^{-1/2})$  and explicit constants.

### 5.1 Regularization

A popular approach is to search not for the LSTD solution, but instead for a regularized LSTD solution defined as follows:

$$\hat{\theta}_T^{reg} = (\bar{A}_T + \mu I)^{-1} \bar{b}_T \quad (10)$$

where  $\mu$  is now a constant set in advance. The update rule for this variant is

$$\theta_n^{reg} = (1 - \gamma_n \mu) \theta_{n-1} + \gamma_n (r_{i_n} + \beta \theta_{n-1}^\top \phi(s'_{i_n}) - \theta_{n-1}^\top \phi(s_{i_n})) \phi(s_{i_n}). \quad (11)$$

This algorithm retains all the properties of the non-regularized fLSTD-SA algorithm, except that it converges to the solution of (10) rather than to that of (5). In particular the conclusions of Theorem 1, and of Corollary 1 hold without requiring assumption (A4), but measuring  $\theta_n - \hat{\theta}_T^{reg}$ , the error to the regularized fixed point  $\hat{\theta}_T^{reg}$ .

### 5.2 Iterate Averaging

Another well-known approach is to employ the Polyak-Ruppert scheme of averaging the iterates, together with choosing larger step-sizes. In particular, we fix the step-size  $\gamma_n := \frac{(1-\beta)}{2} \left(\frac{c}{c+n}\right)^\alpha$ , and then use the averaged iterate  $\bar{\theta}_{n+1} := (\theta_1 + \dots + \theta_n)/n$  to approximate the LSTD solution. Here the quantities  $\theta_n$  are just the iterates of the fLSTD-SA presented earlier. An analogue of Corollary 1 for iterate averaging is as follows (see [16] for a detailed proof):

**Corollary 2.** *Under (A2)-(A3), choosing  $\gamma_n = \frac{(1-\beta)}{2} \left(\frac{c}{c+n}\right)^\alpha$ , with  $\alpha \in (1/2, 1)$  and  $c \in (1.33, 2)$ , we have, for any  $\delta > 0$ ,*

$$\mathbb{E} \|\bar{\theta}_n - \hat{\theta}_T\|_2 \leq \frac{K_1^{IA}(n)}{(n+c)^{\alpha/2}} \text{ and } P \left( \|\bar{\theta}_n - \hat{\theta}_T\|_2 \leq \frac{K_2^{IA}(n)}{(n+c)^{\alpha/2}} \right) \geq 1 - \delta, \quad (12)$$

where, writing  $C = \sum_{n=1}^{\infty} \exp(-\mu c n^{1-\alpha}) (< \infty)$ ,

$$K_1^{IA}(n) := \frac{C \|\theta_0 - \hat{\theta}_T\|_2}{(n+c)^{(1-\alpha)/2}} + \frac{H_\beta c^\alpha (1-\beta)}{(\mu c^\alpha (1-\beta)^2)^{\alpha \frac{1+2\alpha}{2(1-\alpha)}}}, \text{ and}$$

$$K_2^{IA}(n) := \frac{\sqrt{\log \delta^{-1}}}{\mu(1-\beta)} \left[ 3^\alpha + \left[ \frac{2\alpha}{\mu c^\alpha (1-\beta)^2} + \frac{2^\alpha}{\alpha} \right]^2 \right] \frac{1}{(n+c)^{(1-\alpha)/2}} + K_1^{IA}(n).$$

Thus, it is possible to remove the dependency on the knowledge of  $\mu$  for the choice of  $c$  through averaging of the iterates, at the cost of  $(1-\alpha)/2$  in the rate. However, choosing  $\alpha$  close to 1 causes a sampling error blowup. As suggested by earlier works on stochastic approximation, it is preferred to average after a few iterations since the initial error is not forgotten exponentially faster than the sampling error with averaging.



## 6 Outline of Analysis

In this section we give outline proofs of the main results concerning the fLSTD-SA algorithm. We split these into two sections: first, we sketch the martingale analysis that leads to the proof of Theorem 1 and which forms the template for the proof for extension to least squares regression (see Appendix C in [16]) and the regularized and iterate averaged variants of fLSTD-SA (see Corollary 2); second, we give the derivation of the rates when the step sizes are chosen in specific forms.

### 6.1 Outline of Theorem 1 Proof

Denote the approximation error by  $z_n := \theta_n - \hat{\theta}_T$ . Recall that Theorem 1 decomposes the problem of bounding  $z_n$  into bounding the deviation from its mean in high probability and then the mean of  $z_n$  itself. In the following, we first provide a sketch of the proof of high probability bound and later outline the proof for the bound in expectation. For the former, we employ a proof technique similar to that used in [7]. However, our analysis is much simpler and we make all the constants explicit for the problem at hand. Moreover, in order to eliminate a possible exponential dependence of the constants in the resulting bound on the inverse of  $(1 - \beta)\mu$ , we depart from the argument in [7].

*Proof (High probability bound).* (Sketch) Recall that  $z_n := \theta_n - \hat{\theta}_T$ . We rewrite  $\|z_n\|_2^2 - \mathbb{E}\|z_n\|_2^2$  as a telescoping sum of martingale differences:

$$\|z_n\|_2 - \mathbb{E}\|z_n\|_2 = \sum_{i=1}^n g_i - \mathbb{E}[g_i | \mathcal{F}_{i-1}] = \sum_{i=1}^n D_i,$$

where  $D_i := g_i - \mathbb{E}[g_i | \mathcal{F}_{i-1}]$ ,  $g_i := \mathbb{E}[\|z_n\|_2 | \theta_i]$ , and  $\mathcal{F}_i$  denotes the sigma algebra generated by the random variables  $\{i_1, \dots, i_n\}$ .

The next step is to show that the functions  $g_i$  are Lipschitz continuous in the rewards, with Lipschitz constants  $L_i$ . In order to obtain constants with no exponential dependence on the inverse of  $(1 - \beta)\mu$  we depart from the general scheme of [7], and use our knowledge of the form of the update function  $f_i$  to eliminate the noise due to the rewards between time  $i + 1$  and time  $n$ . Specifically, letting  $\Theta_j^i(\theta)$  denote the mapping that returns the value of the iterate  $\theta_j$  at instant  $j$ , given that  $\theta_i = \theta$ , we show that

$$\begin{aligned} \mathbb{E} [\|\Theta_n^i(\theta) - \Theta_n^i(\theta')\|_2^2] &= \mathbb{E} [\mathbb{E} ([I - \gamma_n \phi(s_{i_n}) \phi(s_{i_n})^\top - \beta \phi(s_{i_n}) \phi(s'_{i_n})^\top]) \\ &\quad \cdot (\Theta_{n-1}^i(\theta) - \Theta_{n-1}^i(\theta')) \mid \Theta_{n-1}^i(\theta), \Theta_{n-1}^i(\theta')])] \\ &\leq (1 - \gamma_n \mu (1 - \beta - \gamma_n \beta (2 - \beta))) \mathbb{E} [\|\Theta_{n-1}^i(\theta) - \Theta_{n-1}^i(\theta')\|_2^2], \end{aligned}$$

where we used the specific form of  $f_i$  in obtaining the equality, and have applied assumption (A4) to obtain the inequality. Unrolling this iteration then yields the new Lipschitz constants.

Now we can invoke a standard martingale concentration bound: Using the  $L_i$ -Lipschitz property of the  $g_i$  functions and the assumption (A3) we find that

$$P(\|z_n\|_2 - \mathbb{E}\|z_n\|_2 \geq \epsilon) = P\left(\sum_{i=1}^n D_i \geq \epsilon\right) \leq \exp(-\lambda\epsilon) \exp\left(\frac{\alpha\lambda^2}{2} \sum_{i=1}^n L_i^2\right).$$

The claim follows by optimizing the above over  $\lambda$ . The full proof is available in [16].

*Proof (Bound in expectation).* (Sketch) First we extract a martingale difference from the update rule (6): Recall that  $z_n := \theta_n - \hat{\theta}_T$ . Let  $f_n(\theta) := (\theta^\top x_{i_n} - (r_{i_n} + \beta \theta^\top x'_{i_n})) x_{i_n}$  and let  $F(\theta) := \mathbb{E}_{i_n}(f_n(\theta))$ . Then, we have

$$z_n = \theta_n - \hat{\theta}_T = \theta_{n-1} - \hat{\theta}_T - \gamma_n (F(\theta_{n-1}) - \Delta M_n),$$

where  $\Delta M_{n+1}(\theta) = F_n(\theta) - f_n(\theta)$  is a martingale difference. Now since  $\hat{\theta}_T$  is the LSTD solution,  $F(\hat{\theta}_T) = 0$ . Moreover,  $F(\cdot)$  is linear, and so we obtain

$$z_n = z_{n-1} - \gamma_n (z_{n-1} \bar{A}_n - \Delta M_n) = \Pi_n z_0 - \sum_{k=1}^n \gamma_k \Pi_n \Pi_k^{-1} \Delta M_k,$$

where  $\bar{A}_n = \frac{1}{n} \sum_{i=1}^n x_i (x_i - \beta x'_i)^\top$  and  $\Pi_n := \prod_{k=1}^n (I - \gamma_k \bar{A}_k)$ .

By Jensen's inequality, we obtain

$$\mathbb{E}(\|z_n\|_2) \leq (\mathbb{E}(\langle z_n, z_n \rangle))^{\frac{1}{2}} = \left( \mathbb{E} \|\Pi_n z_0\|_2^2 + \sum_{k=1}^n \gamma_k^2 \mathbb{E} \|\Pi_n \Pi_k^{-1} \Delta M_k\|_2^2 \right)^{\frac{1}{2}} \quad (13)$$

The rest of the proof amounts to bounding the martingale difference  $\Delta M_n$  as follows:

$$\mathbb{E}[\|\Delta M_n\|_2^2] \leq \mathbb{E}_{i_t} \langle f_{i_t}(\theta_{t-1}), f_{i_t}(\theta_{t-1}) \rangle \leq R_{\max}(R_{\max} + 2) + (1 + \beta)^2 \|\theta_{t-1}\|_2^2 \leq H_\beta^2.$$

## 6.2 Derivation of Rates

Now we give the proof of Corollary 1, which gives explicitly the rate of convergence of the approximation error in high probability for the specific choice of step sizes:

*Proof (Proof of Corollary 1).* Note that when  $\gamma_n = \frac{(1-\beta)c}{2(c+n)}$ ,

$$\begin{aligned} \sum_{i=1}^n L_i^2 &= \sum_{i=1}^n \frac{(1-\beta)^2 c^2}{4(c+i)^2} \prod_{j=i}^n \left( 1 - 2\mu \frac{(1-\beta)c}{2(c+n)} \left( (1-\beta) - \beta(2-\beta) \frac{(1-\beta)c}{2(c+n)} \right) \right) \\ &\leq \sum_{i=1}^n \frac{(1-\beta)^2 c^2}{4(c+i)^2} \exp \left( -\frac{3}{4} (1-\beta)^2 \mu c \sum_{j=i}^n \frac{1}{(c+n)} \right) \\ &\leq \frac{(1-\beta)^2 c^2}{4(n+c)^{\frac{3}{4}(1-\beta)^2 \mu c}} \sum_{i=1}^n (i+c)^{-(2-\frac{3}{4}(1-\beta)^2 \mu c)}. \end{aligned}$$

We now find three regimes for the rate of convergence, based on the choice of  $c$ :

- (i)  $\sum_{i=1}^n L_i^2 = O \left( (n+c)^{\frac{3}{4}(1-\beta)^2 \mu c} \right)$  when  $\frac{3}{4}(1-\beta)^2 \mu c \in (0, 1)$ ,
- (ii)  $\sum_{i=1}^n L_i^2 = O(n^{-1} \ln n)$  when  $\frac{3}{4}(1-\beta)^2 \mu c = 1$ , and
- (iii)  $\sum_{i=1}^n L_i^2 = \frac{(1-\beta)^2 c^2}{4(\frac{3}{4}(1-\beta)^2 \mu c - 1)} (n+c)^{-1}$  when  $\frac{3}{4}(1-\beta)^2 \mu c \in (1, 2)$ .

(We have used comparisons with integrals to bound the summations.) Thus, setting  $2/((1-\beta)^2\mu) > c > 1/((1-\beta)^2\mu)$ , the high probability bound from Theorem 1 gives

$$P(\|\theta_n - \hat{\theta}_T\|_2 - \mathbb{E}\|\theta_n - \hat{\theta}_T\|_2 \geq \epsilon) \leq \exp\left(-\frac{\epsilon^2(n+c)}{2K_{\mu,c,\beta}}\right) \quad (14)$$

where  $K_{\mu,c,\beta} := \frac{(1-\beta)^2c^2}{4((1-\beta)^2\mu c - 1)}$ .

Under the same choice of step-size, the bound in expectation in Theorem 1 we have:

$$\begin{aligned} & \sum_{k=1}^{n-1} H_\beta^2 \gamma_{k+1}^2 \exp(-2(1-\beta)\mu(\Gamma_n - \Gamma_{k+1})) \\ & \leq \frac{(1-\beta)^2c^2H_\beta^2}{4(n+c)(1-\beta)^2\mu c} \sum_{k=1}^n (c+k)^{-(2-(1-\beta)^2\mu c)} \leq \frac{(1-\beta)^2c^2H_\beta^2}{4(n+c)} \end{aligned}$$

we in the last inequality we have again compared the sum with an integral. Similarly

$$\exp(-(1-\beta)\mu\Gamma_n) \leq \left(\frac{c}{n+c}\right)^{\frac{(1-\beta)^2\mu c}{2}} \leq \left(\frac{c}{n+c}\right)^{\frac{1}{2}}.$$

So we have

$$\mathbb{E}\|\theta_n - \hat{\theta}_T\|_2 \leq \left(\sqrt{c}\|\theta_0 - \theta^*\|_2 + \frac{(1-\beta)cH_\beta}{2}\right) (c+n)^{-\frac{1}{2}}, \quad (15)$$

and the result now follows.

## 7 Extension to Least Squares Regression

In this section, we describe the classic parameter estimation problem using the method of least squares, the standard approach to solve this problem and a low-complexity alternative using stochastic approximation.

In this setting, we are given a set of samples  $\mathcal{D} := \{(x_i, y_i), i = 1, \dots, T\}$  with the underlying observation model  $y_i = x_i^\top \theta^* + \xi_i$  ( $\xi_i$  is zero mean and variance bounded by  $\sigma < \infty$ , and  $\theta^*$  is an unknown parameter). The least squares estimate  $\hat{\theta}_T$  minimizes  $\sum_{i=1}^T (y_i - \theta^\top x_i)^2$ . It can be shown that  $\hat{\theta}_T = \bar{A}_T^{-1} \bar{b}_T$ , where  $\bar{A}_T = T^{-1} \sum_{i=1}^T x_i x_i^\top$  and  $\bar{b}_T = T^{-1} \sum_{i=1}^T x_i y_i$ .

Notice that, unlike the RL setting,  $\hat{\theta}_T$  here is the minimizer of an empirical loss function. However, as in the case of LSTD, the computational cost for a Sherman-Morrison lemma based approach for solving the above would be of the order  $O(d^2T)$ . Similarly to the case of the fLSTD-SA algorithm, we update the iterate  $\theta_n$  using a SA scheme as follows (starting with an arbitrary  $\theta_0$ ),

$$\theta_n = \theta_{n-1} + \gamma_n (y_{i_n} - \theta_{n-1}^\top x_{i_n}) x_{i_n}, \quad (16)$$

where, as before, each  $i_n$  is chosen uniformly randomly from the sample set  $\mathcal{D}$  and  $\gamma_n$  are step-sizes.

Unlike fLSTD-SA which is a fixed point iteration, the above is a stochastic gradient descent procedure. Nevertheless, using the same proof template as for fLSTD-SA earlier, we can derive bounds on the approximation error, i.e., the distance between  $\theta_n$  and least squares solution  $\hat{\theta}_T$ , both in high probability as well as expectation.

**Results.** As in the case of fLSTD-SA, we assume that the features are bounded, the noise is i.i.d, zero-mean and bounded and the matrix  $\bar{A}_T$  is positive definite, with smallest eigenvalue at least  $\mu > 0$ . An analogue of Corollary 1 for this setting is as follows (See Appendix C in [16] for a detailed proof.):

**Corollary 3.** *Choosing  $\gamma_n = \frac{c}{2(c+n)}$  and  $c$  such that  $\mu c \in (1.33, 2)$ , for any  $\delta > 0$ ,*

$$\mathbb{E}\|\theta_n - \hat{\theta}_T\|_2 \leq \frac{K_1^{LS}}{\sqrt{n+c}} \text{ and } P\left(\|\theta_n - \hat{\theta}_T\|_2 \leq \frac{K_2^{LS}}{\sqrt{n+c}}\right) \geq 1 - \delta,$$

where, defining  $h(n) := c\left[\left(\sigma + 2\|\theta_0 - \hat{\theta}_T\|_2^2\right) + 4\|\theta_0 - \hat{\theta}_T\|_2 \ln n + 2\ln^2 n\right]$ ,

$$K_1^{LS}(n) := \frac{\sqrt{c}\|\theta_0 - \hat{\theta}_T\|_2}{(n+c)^{(\mu c-1)/2}} + \frac{h(n)}{2}, \quad K_2^{LS}(n) := \frac{\sqrt{c}}{\sqrt{((\mu c)/2-1)}} \sqrt{\log \frac{1}{\delta}} + K_1(n).$$

## 8 Traffic Control Application

LSPI [11] is a well-known algorithm for control based on the policy iteration procedure for MDPs. It performs policy evaluation and policy improvement in tandem. For the purpose of policy evaluation, LSPI uses a LSTD-like algorithm called LSTDQ, which learns the state-action value function. In contrast, LSTD learns the state value function.

We now briefly describe LSTDQ and its fast SA variant fLSTDQ-SA: We are given a set of samples  $\mathcal{D} := \{(s_i, a_i, r_i, s'_i), i = 1, \dots, T\}$ , where each sample  $i$  denotes a one-step transition of the MDP from state  $s_i$  to  $s'_i$  under action  $a_i$ , while resulting in a reward  $r_i$ . LSTDQ attempts to approximate the Q-value function for any policy  $\pi$  by solving the linear system  $\hat{\theta}_T = \bar{A}_T^{-1} \bar{b}_T$ , where  $\bar{A}_T = T^{-1} \sum_{i=1}^T \phi(s_i, a_i)(\phi(s_i, a_i) - \beta \phi(s'_i, \pi(s'_i)))^\top$ , and  $\bar{b}_T = T^{-1} \sum_{i=1}^T r_i \phi(s_i, a_i)$ . fLSTDQ-SA approximates LSTDQ by an iterative update scheme as follows (starting with an arbitrary  $\theta_0$ ):

$$\theta_k = \theta_{k-1} + \gamma_k \left( r_{i_k} + \beta \theta_{k-1}^\top \phi(s'_{i_k}, \pi_n(s'_{i_k})) - \theta_{k-1}^\top \phi(s_{i_k}, a_{i_k}) \right) \phi(s_{i_k}, a_{i_k}) \quad (17)$$

From Section 3, it is evident that the claims in Theorem 1 and Corollary 1 hold for the above scheme as well.

The idea behind the experimental setup is to study both LSPI and a variant of LSPI, referred to as fLSPI-SA, where we use fLSTDQ-SA as a subroutine to approximate the LSTDQ solution. Algorithm 1 provides the pseudo-code for the latter algorithm.

We consider a traffic signal control application for conducting the experiments. The problem here is to adaptively choose the sign configurations for the signalized intersections in the road network considered, in order to maximize the traffic flow in the long run. Let  $L$  be the total number of lanes in the road network considered. Further, let  $q_i(t), i = 1, \dots, L$  denote the queue lengths and  $t_i(t), i = 1, \dots, L$  the elapsed time (since signal turned to red) on the individual lanes of the road network. Following [14], the traffic signal control MDP is formulated as follows:

**Algorithm 1.** fLSPI-SA

**Input:** Sample set  $D := \{s_i, a_i, r_i, s'_i\}_{i=1}^T$ , obtained from an initial (arbitrary) policy

**Initialisation:**  $\epsilon, \tau$ , step-sizes  $\{\gamma_k\}_{k=1}^T$ , initial policy  $\pi_0$  (given as  $\theta_0$ )

$\pi \leftarrow \pi_0, \theta \leftarrow \theta_0$

**repeat**

**Policy Evaluation**

Approximate LSTDQ( $D, \pi$ ) using fLSTDQ-SA( $D, \pi$ ) as follows:

**for**  $k = 1 \dots \tau$  **do**

Get random sample index:  $i_k \sim U(\{1, \dots, T\})$

Update fLSTD-SA iterate  $\theta_k$  using (17)

**end for**

$\theta' \leftarrow \theta_\tau, \Delta = \|\theta - \theta'\|_2$

**Policy Improvement**

Obtain a greedy policy  $\pi'$  as follows:  $\pi'(s) = \arg \max_{a \in \mathcal{A}} \theta'^T \phi(s, a)$

$\theta \leftarrow \theta', \pi \leftarrow \pi'$

**until**  $\Delta < \epsilon$

**State**  $s_t = (q_1(t), \dots, q_L(t), t_1(t), \dots, t_L(t))$ ,

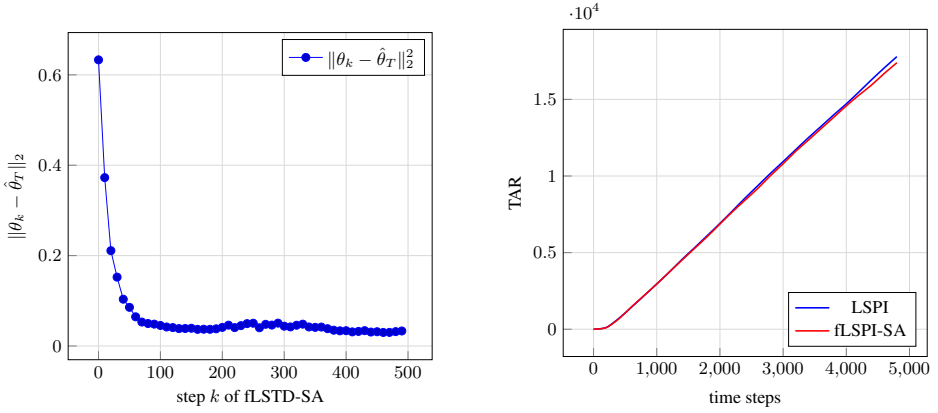
**Action**  $a_t$  belongs to the set of feasible sign configurations,

**Single-stage cost**  $h(s_t) = u_1 \left[ \sum_{i \in I_p} u_2 \cdot q_i(t) + \sum_{i \notin I_p} w_2 \cdot q_i(t) \right] + w_1 \left[ \sum_{i \in I_p} u_2 \cdot t_i(t) + \sum_{i \notin I_p} w_2 \cdot t_i(t) \right]$ , where  $u_i, w_i \geq 0$  such that  $u_i + w_i = 1$  for  $i = 1, 2$  and  $u_2 > w_2$ . Here, the set  $I_p$  is the set of prioritized lanes.

**Table 1.** Feature selection

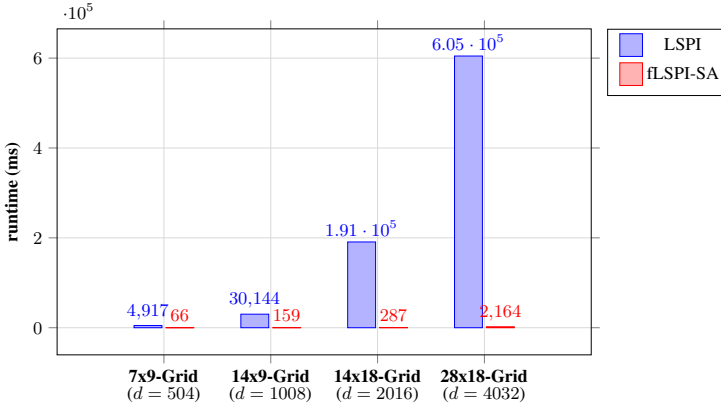
State	Action	Feature $\phi_i(s, a)$
$q_i < \mathcal{L}_1$ and $t_i < \mathcal{T}_1$	RED	0.01
	GREEN	0.06
$q_i < \mathcal{L}_1$ and $t_i \geq \mathcal{T}_1$	RED	0.02
	GREEN	0.05
$\mathcal{L}_1 \leq q_i < \mathcal{L}_2$ and $t_i < \mathcal{T}_1$	RED	0.03
	GREEN	0.04
$\mathcal{L}_1 \leq q_i < \mathcal{L}_2$ and $t_i \geq \mathcal{T}_1$	RED	0.04
	GREEN	0.03
$q_i \geq \mathcal{L}_2$ and $t_i < \mathcal{T}_1$	RED	0.05
	GREEN	0.02
$q_i \geq \mathcal{L}_2$ and $t_i \geq \mathcal{T}_1$	RED	0.06
	GREEN	0.01

Function approximation is a standard technique employed to handle high-dimensional state spaces (as is the case with the traffic signal control MDP on large road networks). We employ the feature selection scheme from [15], which is briefly described in the



(a) Norm difference on 7x9-grid network

(b) Throughput (TAR) on 7x9-grid network



(c) Run-times on four road networks

**Fig. 2.** Norm difference, throughput and runtime performance of LSPI and fLSPI-SA

following: The features  $\phi(s, a)$  corresponding to any state-action tuple  $(s, a)$  is a  $L$ -dimensional vector, with one bit for each line in the road network. The feature value  $\phi_i(s, a)$ ,  $i = 1, \dots, L$  corresponding to lane  $i$  is chosen as described in Table. 1, with  $q_i$  and  $t_i$  denoting the queue length and elapsed times for lane  $i$ . Thus, as the size of the network increases, the feature dimension scales in a linear fashion.

Note that the above feature selection scheme depends on certain thresholds  $\mathcal{L}_1$  and  $\mathcal{L}_2$  on the queue length and  $\mathcal{T}_1$  on the elapsed times. The motivation for using such graded thresholds is owing to the fact that queue lengths are difficult to measure precisely in practice. We set  $(\mathcal{L}_1, \mathcal{L}_2, \mathcal{T}_1) = (6, 14, 130)$  in all our experiments and this choice has been used, for instance, in [15].

We implement both LSPI as well as fLSPI-SA for the above problem. We collect  $T = 10000$  samples from an exploratory policy that picks the actions in a uniformly random

manner. For both LSPI and fLSPI-SA, we set  $\beta = 0.9$  and  $\epsilon = 0.1$ . For fLSPI-SA, we set  $\tau = 500$  steps. This choice is motivated by an experiment where we observed that at 500 steps, fLSTD-SA is already very close to LSTDQ and taking more steps did not result in any significant improvements for fLSPI-SA. We implement the regularized variant of LSTDQ, with regularization constant  $\mu$  set to 1. Motivated by Corollary 1, we set the step-size  $\gamma_k = (1 - \beta)c/(2(c + k))$ , with  $c = 1.33(1 - \beta)^{-2}$ .

**Results** We report the norm differences, total arrived road users (TAR) and run-times obtained from our experimental runs in Figs. 2a–2c. Norm difference measures the distance in  $\ell^2$  norm between the fLSTD-SA iterate  $\theta_k$ ,  $k = 1, \dots, \tau$  and LSTDQ solution  $\hat{\theta}_T$  in iteration 1 of fLSPI-SA. TAR is a throughput metric that denotes the total number of road users who have reached their destination. The choice 1 of the iteration in Fig 2a is arbitrary, as we observed that fLSTD-SA iterate  $\theta_\tau$  is close to the corresponding LSTDQ solution in each iteration of fLSPI-SA. The runtime reports in Fig. 2c are for four different road networks of increasing size and hence, increasing feature dimension.

From Fig. 2a, we observe that fLSTD-SA algorithm converges rapidly to the corresponding LSTDQ solution. Further, from the runtime plots (see Fig. 2c), we notice that fLSPI-SA is several orders of magnitude faster than regular LSPI. From a traffic application standpoint, we observe in Fig. 2b that fLSPI-SA results in a throughput (TAR) performance that is on par with LSPI.

## 9 Conclusions

We analysed a stochastic approximation based algorithm with randomised samples for policy evaluation by the method of LSTD. We provided convergence rate results for this algorithm, both in high probability and in expectation. Further, we also established that using this scheme in place of LSTD does not impact the rate of convergence of the approximate value function to the true value function. This result coupled with the fact that the SA based scheme possesses lower computational complexity in comparison to traditional techniques, makes it attractive for implementation in *big data* settings, where the feature dimension is large. On a traffic signal control application, we demonstrated the practicality of a low-complexity alternative to LSPI that uses our SA based scheme in place of LSTDQ for policy evaluation.

**Acknowledgments.** The first and third authors would like to thank the European Community’s Seventh Framework Programme (FP7/2007 – 2013) under grant agreement n<sup>o</sup> 270327 for funding the research leading to these results. The second author was gratefully supported by the EPSRC Autonomous Intelligent Systems project EP/I011587.

## References

1. Bach, F., Moulines, E.: Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In: NIPS (2011)
2. Bertsekas, D.P.: Dynamic Programming and Optimal Control, 4th edn. Approximate Dynamic Programming, vol. II (2012)

3. Bertsekas, D.P., Tsitsiklis, J.N.: Neuro-Dynamic Programming. Optimization and Neural Computation Series 3, vol. 7. Athena Scientific (1996)
4. Bradtke, S., Barto, A.: Linear least-squares algorithms for temporal difference learning. *Machine Learning* 22, 33–57 (1996)
5. Dani, V., Hayes, T.P., Kakade, S.M.: Stochastic linear optimization under bandit feedback. In: COLT, pp. 355–366 (2008)
6. Fathi, M., Frikha, N.: Transport-entropy inequalities and deviation estimates for stochastic approximation schemes. arXiv preprint arXiv:1301.7740 (2013)
7. Frikha, N., Menozzi, S.: Concentration Bounds for Stochastic Approximations. *Electron. Commun. Probab.* 17(47), 1–15 (2012)
8. Geramifard, A., Bowling, M., Zinkevich, M., Sutton, R.S.: iLSTD: Eligibility traces and convergence analysis. In: NIPS, vol. 19, p. 441 (2007)
9. Hazan, E., Kale, S.: Beyond the regret minimization barrier: an optimal algorithm for stochastic strongly-convex optimization, pp. 421–436 (2011)
10. Kushner, H.J., Yin, G.: Stochastic approximation and recursive algorithms and applications, vol. 35. Springer (2003)
11. Lagoudakis, M.G., Parr, R.: Least-squares policy iteration. *The Journal of Machine Learning Research* 4, 1107–1149 (2003)
12. Lazaric, A., Ghavamzadeh, M., Munos, R.: Finite-sample analysis of least-squares policy iteration. *Journal of Machine Learning Research* 13, 3041–3074 (2012)
13. Polyak, B.T., Juditsky, A.B.: Acceleration of stochastic approximation by averaging. *SIAM Journal on Control and Optimization* 30(4), 838–855 (1992)
14. Prashanth, L., Bhatnagar, S.: Reinforcement Learning with Function Approximation for Traffic Signal Control. *IEEE Transactions on Intelligent Transportation Systems* 12(2), 412–421 (2011)
15. Prashanth, L., Bhatnagar, S.: Threshold Tuning using Stochastic Optimization for Graded Signal Control. *IEEE Transactions on Vehicular Technology* 61(9), 3865–3880 (2012)
16. Prashanth, L., Korda, N., Munos, R.: Fast LSTD using stochastic approximation: Finite time analysis and application to traffic control. arXiv preprint arXiv:1306.2557v4 (2014)
17. Robbins, H., Monro, S.: A stochastic approximation method. In: *The Annals of Mathematical Statistics*, pp. 400–407 (1951)
18. Ruppert, D.: Stochastic approximation. In: *Handbook of Sequential Analysis*, pp. 503–529 (1991)
19. Silver, D., Sutton, R.S., Müller, M.: Reinforcement Learning of Local Shape in the Game of Go. In: IJCAI, vol. 7, pp. 1053–1058 (2007)
20. Sutton, R.S., Barto, A.G.: Reinforcement learning: An introduction, vol. 1. Cambridge Univ. Press (1998)
21. Sutton, R.S., Szepesvári, C., Maei, H.R.: A convergent  $O(n)$  algorithm for off-policy temporal-difference learning with linear function approximation, pp. 1609–1616 (2009)
22. Sutton, R.S., et al.: Fast gradient-descent methods for temporal-difference learning with linear function approximation. In: ICML, pp. 993–1000. ACM (2009)
23. Tsitsiklis, J.N., Van Roy, B.: An analysis of temporal-difference learning with function approximation. *IEEE Transactions on Automatic Control* 42(5), 674–690 (1997)
24. Webscope, Y.: Yahoo! Webscope dataset ydata-frontpage-todaymodule-clicks-v2\_0 (2011), “[http://research.yahoo.com/Academic\\_Relations](http://research.yahoo.com/Academic_Relations)”
25. Zinkevich, M.: Online convex programming and generalized infinitesimal gradient ascent. In: ICML, pp. 928–925 (2003)