

On Convergence of Gradient Expected Sarsa(λ)

Long Yang¹, Gang Zheng¹, Yu Zhang¹, Qian Zheng², Pengfei Li¹, Gang Pan^{1*}

¹College of Computer Science and Technology, Zhejiang University, China.

²School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore.

¹{yanglong, gang_zheng, hzzhangyu, pfl, gpan}@zju.edu.cn; ²zhengqian@ntu.edu.sg

Abstract

We study the convergence of Expected Sarsa(λ) with linear function approximation. We show that applying the off-line estimate (multi-step bootstrapping) to Expected Sarsa(λ) is unstable for off-policy learning. Furthermore, based on convex-concave saddle-point framework, we propose a convergent Gradient Expected Sarsa(λ) (GES(λ)) algorithm. The theoretical analysis shows that our GES(λ) converges to the optimal solution at a linear convergence rate, which is comparable to extensive existing state-of-the-art gradient temporal difference learning algorithms. Furthermore, we develop a Lyapunov function technique to investigate how the step-size influences finite-time performance of GES(λ), such technique of Lyapunov function can be potentially generalized to other GTD algorithms. Finally, we conduct experiments to verify the effectiveness of our GES(λ).

Introduction

Tabular Expected Sarsa(λ) (with importance sampling) is one of the widely used methods for off-policy evaluation in reinforcement learning (RL), whose goal is to estimate the value function of a given target policy via the data that is generated from a behavior policy. Due to the high-dimensional state space, instead of tabular learning, a standard approach is to estimate the value function with a linear function (Sutton and Barto 2018). There is very little literature to study Expected Sarsa(λ) with function approximation for off-policy learning. To our best knowledge, Sutton and Barto (2018) (section 12.9) firstly extend *off-line* estimate (multi-step bootstrapping) to Expected Sarsa(λ) with linear function approximation.

Unfortunately, as pointed out by Sutton and Barto (2018) intuitively, their off-line approach may be unstable, i.e., their way to extend Expected Sarsa(λ) with linear function approximation still lacks a provable convergence guarantee, which is undesirable for RL. It is critical to find the inherent essence of the above unstable appearance in Expected Sarsa(λ), which not only makes a complement for existing off-policy learning methods but also provides some inspirations to design a stable algorithm. Thus, extending Expected Sarsa(λ) with linear function approximation

for off-policy evaluation is a fundamental theoretical topic in RL, including: 1) how to character the instability of off-line Expected Sarsa(λ) with linear function approximation; 2) how to derive a convergent algorithm; 3) what convergence rate does Expected Sarsa(λ) with linear function approximation can reach. We focus on these questions in this paper.

Our Main Works To address the above problems, we propose Theorem 1, which characters a sufficient and necessary condition that presents stability criteria of off-line update Expected Sarsa(λ) with linear function approximation. Theorem 1 requires the *key* matrix (that has been defined in (10)) keeps the negative real components. Unfortunately, due to the discrepancy between behavior policy and target policy, off-line Expected Sarsa(λ) that is suggested by Sutton and Barto (2018) may not satisfy the condition appears in Theorem 1, i.e., their scheme maybe unstable. Then, we use a classic counterexample to verify the above instability lies in the off-line update Expected Sarsa(λ) with linear function approximation, see Example 1.

Furthermore, to get a stable algorithm, we derive an *on-line* gradient Expected Sarsa(λ) (GES(λ)) algorithm. Theorem 2 shows that the proposed GES(λ) learns the optimal solution at a linear convergence rate, which is comparable to extensive existing state-of-the-art gradient temporal difference learning algorithms. Although Xu et al. (2019) prove TDC (Sutton et al. 2009a) also converges at a linear convergence rate, they require a projection step that is unpractical in practice. Besides, the fussy blockwise step-size appears in (Xu et al. 2019) is more complicated than our step-size condition. Additionally, the results of (Dalal et al. 2018a; Lakshminarayanan et al. 2018) require an i.i.d assumption of function parameters, our proof removes this condition and achieve a better result than theirs. A more detailed comparison and an adequate discussion are provided in Table 1.

Finally, inspired by Srikant and Ying (2019), Wang et al. (2019), Gupta et al. (2019), we develop a Lyapunov function technique to establish Theorem 3, which illustrates the relationship between the finite-time performance of GES(λ) and step-size. Result shows that the upper-bounded error consists of two different parts: the first error depends on both step-size and the size of samples, and such error decays geometrically as the samples increase; while the second error is only determined by the step-size and it is independent of samples. Additionally, the technique of proving Theorem 3

*Corresponding author.

can be potentially generalized to other GTD algorithms.

Notations

We use $\text{Spec}(A)$ to denote the eigenvalues of the matrix $A \in \mathbb{C}^{p \times p}$, i.e., $\text{Spec}(A) = \{\lambda_1, \dots, \lambda_p\}$, where λ_i is the root of the characteristic equation $p(\lambda) = \det(A - \lambda I)$. We use \mathbb{C}_- to denote the collection that contains the complex numbers with negative real components, i.e.,

$$\mathbb{C}_- = \{c \in \mathbb{C}; \text{Re}(c) < 0\}.$$

Let $\lambda_{\min}(A)$ and $\lambda_{\max}(A)$ be the minimum and maximum eigenvalue of the matrix A correspondingly. We use $\|A\|_{\text{op}}$ to denote the operator norm of matrix A ; furthermore, if A is a symmetric real matrix, then $\|A\|_{\text{op}} = \max_{1 \leq i \leq p} \{|\lambda_i|\}$. $\kappa(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$ is the condition number of matrix A . We use $A \succ 0$ to denote a positive definite matrix A .

For a function $f(x) : \mathbb{R}^p \rightarrow \mathbb{R}$, let $\nabla^2 f(x)$ denote its Hessian matrix, and its convex conjugate function $f^*(y) : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined as $f^*(y) = \sup_{x \in \mathbb{R}^p} \{y^\top x - f(x)\}$.

Fact 1 ((Rockafellar 1970; Kakade et al. 2009)). *Let $f(\cdot)$ be α -strongly convex and β -smooth, i.e., $f(\cdot) - \frac{\alpha}{2}\|\cdot\|_2^2$ is convex and $\|f(u) - f(v)\| \leq \beta\|u - v\|$. If $0 \leq \alpha \leq \beta$, then the following fact holds,*

- (I) f^* is $\frac{1}{\alpha}$ -smooth and $\frac{1}{\beta}$ -strongly convex.
- (II) $\nabla f = (\nabla f^*)^{-1}$ and $\nabla f^* = (\nabla f)^{-1}$.

Preliminary

Reinforcement learning (RL) is formalized as Markov decision processes (MDP) which considers the tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma)$; \mathcal{S} is the space of states, \mathcal{A} is the space of actions; $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$, $p_{ss'}^a = P(S_t = s' | S_{t-1} = s, A_{t-1} = a)$ is the probability of state transition from s to s' under playing the action a ; $R(\cdot, \cdot) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^1$ is the expected reward function; $\gamma \in (0, 1)$ is the discount factor.

The policy π is a probability distribution on $\mathcal{S} \times \mathcal{A}$, we use $\pi(a|s)$ to denote the probability of playing a under the state s . Let $\{S_t, A_t, R_{t+1}\}_{t \geq 0}$ be generated by a given policy π , its *state-action value function* $q^\pi(s, a) = \mathbb{E}_\pi[\sum_{t=0}^{\infty} \gamma^t R_{t+1} | S_0 = s, A_0 = a]$, where $\mathbb{E}_\pi[\cdot]$ is the conditional expectation on the actions selected according to π . Let $\mathcal{B}^\pi : \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|} \rightarrow \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$ be the *Bellman operator*:

$$\mathcal{B}^\pi : q \mapsto R^\pi + \gamma P^\pi q, \quad (1)$$

where $P^\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$, $R^\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$, their corresponding elements are: $[P^\pi]_{s,s'} = \sum_{a \in \mathcal{A}} \pi(a|s) p_{ss'}^a$, $[R^\pi]_{s,a} = R(s, a)$. It is well-known that q^π is the unique fixed point of \mathcal{B}^π , i.e., $\mathcal{B}^\pi q^\pi = q^\pi$, which is known as Bellman equation.

Off-Policy Evaluation Let's consider the following trajectory \mathcal{T} generated by a policy μ :

$$\mathcal{T} = \{S_0 = s, A_0 = a, \dots, S_t, A_t, R_{t+1}, \dots\},$$

where $A_t \sim \mu(\cdot | S_t)$, $S_{t+1} \sim P(\cdot | S_t, A_t)$. In RL, the task of off-policy evaluation is to estimate the value function of the target policy π via the data that is generated by an another policy μ (that is called *behavior policy*), where $\mu \neq \pi$.

Assumption 1. *The Markov chain induced by behavior policy μ is ergodic, i.e., there exists a stationary distribution $\xi(\cdot, \cdot)$ over $\mathcal{S} \times \mathcal{A}$: for $\forall (S_0, A_0) \in \mathcal{S} \times \mathcal{A}$,*

$$\frac{1}{n} \sum_{k=1}^n P(S_k = s, A_k = a | S_0, A_0) \xrightarrow{n \rightarrow \infty} \xi(s, a) > 0. \quad (2)$$

The ergodicity of behavior policy μ is a standard assumption in off-policy learning (Bertsekas 2012), and it implies each-action pair can be visited under this behavior policy μ . In this paper, we use Ξ to denote a diagonal matrix whose diagonal element is $\xi(s, a)$, i.e.,

$$\Xi = \text{diag}\{\dots, \xi(s, a), \dots\}.$$

Temporal Difference (TD) Learning TD learning updates value function as follows, $\forall t \geq 0$,

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha_t \delta_t, \quad (3)$$

where $Q(\cdot, \cdot)$ is an estimate of state-action value function, α_t is step-size and δ_t is TD error. Let $Q_t = Q(S_t, A_t)$, if δ_t is expected TD error:

$$\delta_t^{\text{ES}} = R_{t+1} + \gamma \sum_{a \in \mathcal{A}} \pi(a | S_{t+1}) Q(S_{t+1}, a) - Q_t, \quad (4)$$

then update (3) is Expected Sarsa.

Expected Sarsa(λ) Sutton and Barto (2018)¹ propose a multi-step TD learning that extends Expected Sarsa to λ -return version: for each $t \geq 0$,

$$G_t^\lambda = Q_t + \sum_{k=t}^{\infty} (\gamma \lambda)^{k-t} \left(\prod_{i=t+1}^k \frac{\pi(A_i | S_i)}{\mu(A_i | S_i)} \right) \delta_k^{\text{ES}}, \quad (5)$$

where δ_k^{ES} is expected TD error. For the convenience, we set $\prod_{i=t}^k \rho_i = \prod_{i=t}^k \frac{\pi(A_i | S_i)}{\mu(A_i | S_i)} = \rho_{t:k}$, and $\rho_{t:t+1} = 1$.

Finally, we introduce λ -operator \mathcal{B}_λ^π that is a high level view of iteration (5):

$$\begin{aligned} \mathcal{B}_\lambda^\pi : q \mapsto q + \mathbb{E}_\mu \left[\sum_{k=t}^{\infty} (\lambda \gamma)^{k-t} \delta_k^{\text{ES}} \rho_{t+1:k} \right] \\ = q + (I - \lambda \gamma P^\pi)^{-1} (\mathcal{B}^\pi q - q), \end{aligned} \quad (6)$$

where \mathcal{B}^π is defined in (1). For the limitation of space, we provide the derivation of (7) from (6) in Appendix A.2.

Linear Function Approximation

TD learning (3) requires a very huge table to store the estimate value function $Q(\cdot, \cdot)$ when $|\mathcal{S}|$ is very large, which implies tabular TD learning is considerably expensive for high-dimensional RL. We often use a parametric function $Q_\theta(\cdot, \cdot)$ to approximate $q^\pi(s, a)$, i.e.,

$$q^\pi(s, a) \approx \phi^\top(s, a) \theta =: Q_\theta(s, a),$$

where $\theta \in \mathbb{R}^p$ is the parameter that needs to be learned, $\phi(s, a) = (\varphi_1(s, a), \varphi_2(s, a), \dots, \varphi_p(s, a))^\top$, and each $\varphi_i : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. Furthermore, Q_θ can be rewritten as a version of matrix $Q_\theta = \Phi \theta \approx q^\pi$, where Φ is a matrix whose rows are the state-action feature vectors $\phi^\top(s, a)$. In this paper, we mainly consider extending λ -return of Expected Sarsa (5) with linear function approximation.

¹It is noteworthy that the λ -return version of Expected Sarsa appears in section 12.9 of (Sutton and Barto 2018) is limited in the case of linear function case, Eq.(5) extends it to be a general case.

Off-Line Gradient Expected Sarsa(λ)

In this section, we use a counterexample to show the way to extend Expected Sarsa(λ) with linear function approximation via off-line estimate is unstable for off-policy learning.

Off-Line Update Sutton and Barto (2018) provide a way to extend (5) with linear function approximation as follows,

$$\begin{aligned}\theta_{t+1} &= \theta_t + \alpha_t (G_t^\lambda - Q_\theta(S_t, A_t)) \nabla Q_\theta(S_t, A_t) \\ &= \theta_t + \alpha_t \left(\sum_{k=t}^{\infty} (\gamma\lambda)^{k-t} \delta_{k,\theta}^{\text{ES}} \rho_{t+1:k} \right) \phi_t,\end{aligned}\quad (8)$$

where α_t is step-size, $\phi_t =: \phi(S_t, A_t)$, and

$$\delta_{k,\theta}^{\text{ES}} = R_{k+1} + \gamma \theta_t^\top \mathbb{E}_\pi[\phi(S_{k+1}, \cdot)] - \theta_t^\top \phi_k.$$

Furthermore, we can rewrite the expected parameter in (8):

$$\mathbb{E}_\mu[\theta_{t+1}] = \theta_t + \alpha_t (A\theta_t + b), \quad (9)$$

where

$$\begin{aligned}A &= \mathbb{E}_\mu \left(\sum_{k=t}^{\infty} (\gamma\lambda)^{k-t} \rho_{t+1:k} \phi_t (\gamma \mathbb{E}_\pi[\phi(S_k, \cdot)] - \phi_k)^\top \right) \\ &= \Phi^\top \Xi (I - \gamma \lambda P^\pi)^{-1} (\gamma P^\pi - I) \Phi,\end{aligned}\quad (10)$$

$$\begin{aligned}b &= \mathbb{E}_\mu \left(\sum_{k=t}^{\infty} (\gamma\lambda)^{k-t} \rho_{t+1:k} \phi_t R_{k+1} \right) \\ &= \Phi^\top \Xi (I - \gamma \lambda P^\pi)^{-1} R.\end{aligned}\quad (11)$$

If θ_t (9) converges to a certain point θ^* , then θ^* satisfies the following linear system:

$$A\theta^* + b = 0. \quad (12)$$

Such θ^* satisfies (12) is called *TD-fixed point*.

Stability Criteria According to Sutton et al. (2016); and Ghosh and Bellemare (2020), we formulate the stability of the iteration (9) as the next definition.

Definition 1 (Stability). *Update rule (9) is stable if θ_k converges to the point θ^* satisfies (12) for any initial θ_0 .*

Theorem 1 (Stability Criteria). *Under Assumption 1, the off-line update (9) is stable if and only if the eigenvalues of the matrix A (10) have negative real components, i.e.,*

$$\text{Spec}(A) \subset \mathbb{C}_-. \quad (13)$$

We provide its proof in Appendix B. Theorem 1 provides a sufficient and necessary condition (13) that guarantees the stability of iteration (8). Unfortunately, for off-policy learning, the matrix A (10) can not guarantee the condition (13) holds, which implies the iteration (8) may be divergent and unstable. Now, we use the following example (Touati et al. 2018) to illustrate the instability lies in the iteration (8).

Example 1. *For the MDP in Figure 1, we assign the features $\{(1, 0)^\top, (2, 0)^\top, (0, 1)^\top, (0, 2)^\top\}$ to the state-action pairs $\{(s_1, \text{right}), (s_2, \text{right}), (s_1, \text{left}), (s_2, \text{left})\}$. From the dynamic transition shown in Figure 1, we have*

$$P^\pi = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}, \Phi = \begin{pmatrix} 1 & 0 \\ 2 & 0 \\ 0 & 1 \\ 0 & 2 \end{pmatrix}, \Xi = \frac{1}{2} I_{4 \times 4}.$$

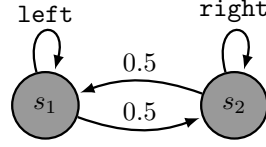


Figure 1: Counterexample, Two-State MDP: behavior policy $\mu(\text{right}|\cdot) = 0.5$ and target policy $\pi(\text{right}|\cdot) = 1$.

Then, according to (10), we have

$$A = \Phi^\top \Xi (I - \gamma \lambda P^\pi)^{-1} (\gamma P^\pi - I) \Phi = \begin{pmatrix} \frac{6\gamma - \gamma\lambda - 5}{2(1 - \gamma\lambda)} & 0 \\ \frac{3\gamma}{2} & -\frac{5}{2} \end{pmatrix},$$

and the eigenvalues of A are: $\frac{6\gamma - \gamma\lambda - 5}{2(1 - \gamma\lambda)}$ and $\frac{5}{2}$. For any initial $\theta_0 = (\theta_{0,1}, \theta_{0,2})^\top$, let $\mathbb{E}[\theta_{t+1}] =: (\theta_{t+1,1}, \theta_{t+1,2})^\top$ be the expectation of iteration (8), according to (9), the first component of $\mathbb{E}[\theta_{t+1}|\theta_t]$ is:

$$\theta_{t+1,1} = \theta_{0,1} \prod_{i=0}^t \left(1 + \alpha_i \frac{6\gamma - \gamma\lambda - 5}{2(1 - \gamma\lambda)} \right). \quad (14)$$

For any $\lambda \in (0, 1)$, if $\gamma \in (\frac{5}{6-\lambda}, 1)$, then $\frac{6\gamma - \gamma\lambda - 5}{2(1 - \gamma\lambda)}$ is a positive scalar, which implies A can not be a negative matrix. Furthermore, if step size $\alpha_t : \sum_{i \geq 0} \alpha_t = \infty$, we have ²

$$|\theta_{t+1,1}| = |\theta_{0,1}| \prod_{i=0}^t \left(1 + \alpha_i \frac{6\gamma - \gamma\lambda - 5}{2(1 - \gamma\lambda)} \right) \rightarrow +\infty, \quad (15)$$

which implies the way (8) to extend Expected Sarsa(λ) with linear function approximation via off-line estimate is unstable for off-policy learning.

On-Line Gradient Expected Sarsa(λ)

The above discussion of the instability for off-policy learning shows that we should abandon the off-line update (8). In this section, we provide a convergent on-line algorithm: Gradient Expected Sarsa(λ) (GES(λ)), which is based on the popular TD fixed point method.

The TD fixed point method (Sutton et al. 2009a; Bertsekas 2011; Dann et al. 2014) is widely used for policy evaluation and it focuses on finding the value function satisfies

$$\Phi\theta = \Pi \mathcal{B}_\lambda^\pi \Phi\theta, \quad (16)$$

where $\Pi = \Phi(\Phi^\top \Xi \Phi)^{-1} \Phi^\top \Xi$. It has been shown that if the *projected Bellman operator* $\Pi \mathcal{B}_\lambda^\pi$ has a fixed point θ^* , then it is unique (Lagoudakis and Parr 2003; Bertsekas 2011), and such a fixed-point θ^* also satisfies the linear system (12).

Instead of using the method of value iteration according to the projected Bellman operator $\Pi \mathcal{B}_\lambda^\pi$, we derive the algorithm on the mean square projected Bellman equation (MSPBE) (Sutton et al. 2009a) as follows,

$$\begin{aligned}\min_{\theta} \text{MSPBE}(\theta, \lambda) &=: \min_{\theta} \frac{1}{2} \|\Phi\theta - \Pi \mathcal{B}_\lambda^\pi(\Phi\theta)\|_{\Xi}^2 \\ &= \min_{\theta} \frac{1}{2} \|A\theta + b\|_{M^{-1}}^2,\end{aligned}\quad (17)$$

²Eq.(15) is a direct result of the following conclusion that could be found in any calculus textbook. Let $p_i = 1 + a_i$, where $a_i > 0$, if $\sum_{i=1}^{\infty} a_i = +\infty$, then $\prod_{i=1}^{\infty} p_i = \prod_{i=1}^{\infty} (1 + a_i) = +\infty$.

Algorithm 1 Gradient Expected Sarsa(λ) (GES(λ))

1: **Initialization:** $\omega_0 = 0, \theta_0 = 0, \alpha_0 > 0, \beta_0 > 0, T \in \mathbb{N}$.
2: $e_{-1} = 0$
3: **for** $t = 0$ **to** T **do**
4: Observe $\{S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}\} \sim \mu$
5: $\rho_t = \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)}$
6: $e_t = \lambda\gamma\rho_t e_{t-1} + \phi_t$
7: $\delta_t = R_{t+1} + \gamma\theta_t^\top \mathbb{E}_\pi \phi(S_{t+1}, \cdot) - \theta_t^\top \phi_t$
8: $\omega_{t+1} = \omega_t + \beta_t(e_t \delta_t - \phi_t \phi_t^\top \omega_t)$
9: $\theta_{t+1} = \theta_t - \alpha_t(\gamma \mathbb{E}_\pi[\phi(S_{t+1}, \cdot)] - \phi_t) e_t^\top \omega_t$
10: **end for**
11: **Output:** $\{\theta_t, \omega_t\}_{t=1}^T$

where $\|x\|_\Xi = x^\top \Xi x$ is a weighted norm, and $M = \Phi^\top \Xi \Phi$. We provide the derivation of (17) in Appendix C.1.

Since the computational complexity of the invertible matrix M^{-1} is very large, it is too expensive to use stochastic gradient method to solve the problem (17) directly. Let

$$g(\omega) = \frac{1}{2} \|\omega\|_M^2 - b^\top \omega \quad (18)$$

$$\Psi(\theta, \omega) = (A\theta + b)^\top \omega - \frac{1}{2} \|\omega\|_M^2 = \theta^\top A\omega - g(\omega).$$

According to Liu et al. (2015), the original problem (17) is equivalent to the convex-concave saddle-point problem

$$\min_{\theta} \max_{\omega} \{\Psi(\theta, \omega)\}. \quad (19)$$

Proposition 1. *If (θ^*, ω^*) is the solution of problem (19), then θ^* is the solution of original problem (17), i.e.,*

$$\theta^* = \arg \min_{\theta} \text{MSPBE}(\theta, \lambda).$$

We provide the proof of Proposition 1 in Appendix C.2. Proposition 1 illustrates that the solution of (17) is contained in the problem (19). Gradient update is a natural way to solve problem (19) (ascending in ω and descending in θ):

$$\omega_{t+1} = \omega_t + \beta_t(A\theta_t + b - M\omega_t), \quad (20)$$

$$\theta_{t+1} = \theta_t - \alpha_t A^\top \omega_t, \quad (21)$$

where α_t, β_t is step-size, $t \geq 0$.

However, since A, b , and M are versions of expectations, we can not get the transition probability in practice. A practical way is to find the unbiased estimators of them. Let $e_0 = 0, \rho_t = \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)}, e_t = \lambda\gamma\rho_t e_{t-1} + \phi_t, \hat{b}_t = R_{t+1} e_t, \hat{A}_t = e_t(\gamma \mathbb{E}_\pi[\phi(S_{t+1}, \cdot)] - \phi_t)^\top, \hat{M}_t = \phi_t \phi_t^\top$. According to the Theorem 9 of (Maei 2011), we have

$$\mathbb{E}_\mu[\hat{A}_t] = A, \mathbb{E}_\mu[\hat{b}_t] = b, \mathbb{E}_\mu[\hat{M}_t] = M. \quad (22)$$

Replacing the expectations in (20) and (21) by corresponding unbiased estimates, we define the stochastic on-line implementation of (20) and (21) as follows,

$$\omega_{t+1} = \omega_t + \beta_t(\hat{A}_t \theta_t + \hat{b}_t - \hat{M}_t \omega_t), \quad (23)$$

$$\theta_{t+1} = \theta_t - \alpha_t \hat{A}_t^\top \omega_t. \quad (24)$$

We provide more details in Algorithm 1.

Finite-Time Performance Analysis

In this section, we mainly focus on the finite-time performance of GES(λ). Theorem 2 shows the proposed GES(λ) converges at a linear rate under a concrete step-size, which is comparable to extensive existing state-of-the-art GTD algorithms. We have provided an adequate comparison in Table 1 and a comprehensive discussion in the section of related works. Furthermore, to investigate how the step-size influences finite-time performance, we establish Theorem 3. We develop a technique of Lyapunov function to prove Theorem 3, before we present the main result, we provide the motivation and some necessary details of Lyapunov function.

Throughout this paper, we make two additional standard assumptions, which are widely used in reinforcement learning (Wang et al. 2017; Bhandari et al. 2018; Xu et al. 2019).

Assumption 2 (Boundedness of Feature Map, Reward). *The features $\{\phi_t\}_{t \geq 0}$ is uniformly bounded by ϕ_{\max} . The reward function is uniformly bounded by R_{\max} . The importance sampling ρ_t is uniformly bounded by ρ_{\max} .*

Assumption 3 (Solvability of Problem). *The matrix A is non-singular and $\text{rank}(\Phi) = p$.*

As claimed by Xu et al., (2019), Assumption 2 can be ensured by normalizing the feature maps $\{\phi_t\}_{t \geq 1}$ and when $\mu(\cdot|s)$ is non-degenerate for all $s \in \mathcal{S}$. Besides, Assumption 2 implies the boundedness of the estimators \hat{A}_t, \hat{M}_t and \hat{b}_t . For the limitation of space, we provide more details and discussions in Remark 5 (see Appendix D.1).

Assumption 3 requires the non-singularity of the matrix A , which implies the optimal parameter $\theta^* = -A^{-1}b$ is well defined. The feature matrix Φ has linearly independent columns implies the matrix M is non-singular.

Linear Convergence Rate

We consider the first-order optimality condition of the problem (19), i.e., the optimal solution (θ^*, ω^*) satisfies

$$\begin{cases} \nabla_{\theta} \Psi(\theta^*, \omega^*) = A^\top \omega^* = 0, \\ \nabla_{\omega} \Psi(\theta^*, \omega^*) = -\nabla g(\omega^*) + A\theta^* = 0. \end{cases} \quad (25)$$

According to the Fact 1 and the condition (25), we have

$$\omega^* = (\nabla g)^{-1}(A\theta^*) = \nabla g^*(A\theta^*),$$

which implies ω^* can be represented by θ^* , thus, we mainly focus on the performance of $\{\theta_t\}_{t \geq 1}$.

Theorem 2. *$\{(\theta_t, \omega_t)\}_{t \geq 0}$ is generated by Algorithm 1. Let $\Delta_{\theta_t} = \|\theta_t - \theta^*\|_2^2, \Delta_{\omega_t} = \|\omega_t - \nabla g^*(A\theta_t)\|_2^2, \nu = \frac{2\kappa^2(A)\kappa(M)\lambda_{\max}(A)}{\lambda_{\min}(M)}$, and $D_t = \nu\Delta_{\theta_t} + \Delta_{\omega_t}$. If we choose*

step-size $\alpha = \frac{\lambda_{\min}(M)}{(\lambda_{\max}(M) + \lambda_{\min}(M)) \left(\frac{\lambda_{\max}^2(A)}{\lambda_{\min}(M)} + \nu\lambda_{\max}(A) \right)}, \beta = \frac{2}{\lambda_{\max}(M) + \lambda_{\min}(M)}$, under Assumption 1-3, we have

$$\mathbb{E}[D_{t+1}] \leq \left(1 - \frac{1}{12\kappa^3(M)\kappa^4(A)}\right) \mathbb{E}[D_t]. \quad (26)$$

Furthermore, we have

$$\mathbb{E}[\|\theta_t - \theta^*\|_2^2] \leq \frac{1}{\nu} \left(1 - \frac{1}{12\kappa^3(M)\kappa^4(A)}\right)^t \mathbb{E}[D_0]. \quad (27)$$

Remark 1. We provide its proof in Appendix C.3. Recall the fact $D_t = \nu \Delta_{\theta_t} + \Delta_{\omega_t} \geq \nu \Delta_{\theta_t}$, which implies the inequality

$$\mathbb{E}[\|\theta_t - \theta^*\|_2^2] \leq \frac{\mathbb{E}[D_t]}{\nu} \stackrel{(26)}{\leq} \frac{1}{\nu} \left(1 - \frac{1}{12 \kappa^3(M) \kappa^4(A)}\right)^t \mathbb{E}[D_0].$$

The term $(1 - \frac{1}{12 \kappa^3(M) \kappa^4(A)}) \in (0, 1)$ implies $\text{GES}(\lambda)$ produces the sequence $\{\theta_t\}_{t \geq 0}$ converges to the optimal solution at a linear convergence rate. After some simple algebra, with a computational cost of

$$\mathcal{O}\left(\max\left\{1, \frac{\lambda_{\max}(A)}{\lambda_{\max}(M)\nu}\right\} \left(1 - \frac{1}{12 \kappa^3(M) \kappa^4(A)}\right) \log \frac{1}{\delta}\right),$$

the output of Algorithm 1 closes to (θ^*, ω^*) as follows,

$$\mathbb{E}[\|\theta_t - \theta^*\|_2^2] \leq \delta^2, \quad \mathbb{E}[\|\omega_t - \omega^*\|_2^2] \leq \delta^2.$$

Remark 2. Theorem 2 provides a concrete step-size that depends on the some unknown parameters. As suggested by Du et al. (2017), Touati et al. (2018), and Voloshin et al. (2019), we can use Monte Carlo method to estimate the unknown parameters.

A Further Analysis via Lyapunov Function.

In this section, we propose Theorem 3 that illustrates a relationship between the performance of $\text{GES}(\lambda)$ and step-size.

The proof of Theorem 3 involves a novel Lyapunov function technique, we start by presenting the motivation behind such Lyapunov function. Let

$$H = -A^\top M^{-1}A, L = 2A.$$

Let Q_1 and Q_2 be the solutions of the following equations

$$\begin{cases} -H^\top Q_1 - Q_1 H &= I \\ M^\top Q_2 + Q_2 M &= I. \end{cases} \quad (28)$$

Since both M and $-H$ are Hurwitz matrix, the solution of the linear system (28) always exists (Lakshminarayanan et al. 2018; Srikant and Ying 2019). Let Q be a matrix as follows,

$$Q = \frac{1}{p_1 + p_2} \begin{pmatrix} p_1 & 0 \\ 0 & p_2 \end{pmatrix},$$

where $p_1 = \|Q_1 A^\top\|_{op} Q_1$, $p_2 = \|Q_2 M^{-1} A L\|_{op} Q_2$. Finally, we define ϱ_t and z_t as follows,

$$\varrho_t = \omega_t - M^{-1} A \theta_t, \quad z_t = (\theta_t - \theta^*, \varrho_t - \varrho^*)^\top, \quad (29)$$

where $\varrho^* = \omega^* - M^{-1} A \theta^*$. Lyapunov function $L(z_t)$ is:

$$L(z_t) = z_t^\top Q z_t. \quad (30)$$

Motivation of Lyapunov Function We consider the expected difference of iteration (23)-(24) as follows

$$\frac{1}{\alpha} \mathbb{E}[\theta_{t+1} - \theta_t | Y_{t-\tau}] = \mathbb{E}[-\hat{A}_t \omega_t | Y_{t-\tau}] \quad (31)$$

$$\frac{\alpha}{\beta} \mathbb{E}[\omega_{t+1} - \omega_t | Y_{t-\tau}] = \mathbb{E}[\hat{A}_t \theta_t + \hat{b}_t - \hat{M}_t \omega_t | Y_{t-\tau}], \quad (32)$$

where $Y_{t-\tau} = \{\theta_{t-\tau}, \omega_{t-\tau}, X_{t-\tau}\}$, and the sequence $X_t = \{S_0, A_0, S_1, A_1, \dots, S_t, A_t\}$ denotes a Markov chain according to Algorithm 1. The expectation is conditioned sufficiently in the past information of the underlying Markov

chain. Approximating the left parts of (31)-(32) by derivatives, then we have the ordinary differential equation (ODE):

$$\begin{cases} \dot{\theta}(t) = -A\omega(t), \\ \frac{\alpha}{\beta} \dot{\omega}(t) = A\theta(t) + b - M\omega(t), \end{cases} \quad (33)$$

where $\theta(t), \omega(t) \in \mathbb{R}^p$ of (33) are the functions that are defined on the continuous time $(0, \infty)$. The update rule of (23)-(24) can be thought of as a discretization of the ODE (33) that is known as *singular perturbation ODE* (chapter 11 of (Khalil and Grizzle 2002)), our goal is to provide a non-asymptotic analysis of $\text{GES}(\lambda)$ according to the asymptotically stable equilibrium of ODE (33). According to Khalil and Grizzle (2002), the following Lyapunov function $L(t)$ is widely used as a stability criteria for the ODE (33),

$$\begin{aligned} L(a, t) &= a(\omega(t) - M^{-1} A \theta(t))^\top Q_2 (\omega(t) - M^{-1} A \theta(t)) \\ &\quad + (1 - a) \theta^\top(t) Q_1 \theta(t), \end{aligned} \quad (34)$$

where $a \in (0, 1)$. Our $L(z_t)$ (30) can be seen as a discretization of $L(t)$ (34) after a proper choice of a , which inspires us to conduct the Lyapunov function $L(z_t)$.

Lemma 1. Under Assumption 1-3, there exists a positive scalar τ such that: $t \geq \tau$,

$$\begin{aligned} \mathbb{E}[L(z_{t+1})] - \mathbb{E}[L(z_t)] &\leq -\alpha \left(\frac{1}{2} \varkappa_1 - \frac{\alpha}{\beta} \varkappa_2\right) \mathbb{E}[L(z_t)] \\ &\quad + 2\alpha^2 \zeta^2 \lambda_{\max}(Q) \tilde{c}_b^2, \end{aligned} \quad (35)$$

where the constants $\varkappa_1, \varkappa_2, \zeta =: C_1 + \frac{\beta}{\alpha} C_2, \tilde{c}_b$ are defined in Appendix D.

Finally, we know

$$\mathbb{E}[\|z_t\|_2^2] \leq (\lambda_{\min}(Q))^{-1} \mathbb{E}[L(z_t)],$$

applying the result of (35) recurrently, we have Theorem 3.

Theorem 3. Let $\eta_1 = 4\zeta^2 \tau^2 (\|z_0\|_2 + \tilde{c}_b)^2 + \|z_0\|_2^2$, $\eta_2 = \frac{2\kappa(Q)\zeta^2 \lambda_{\max}(Q) \tilde{c}_b^2}{\frac{1}{2} \varkappa_1 - \frac{\alpha}{\beta} \varkappa_2}$. Under Assumption 1-3, there exists a positive scalar τ such that: $t \geq \tau$,

$$\mathbb{E}[\|z_t\|_2^2] \leq \alpha^2 \eta_1 \left(1 - \frac{\alpha}{\lambda_{\max}(Q)} \left(\frac{1}{2} \varkappa_1 - \frac{\alpha}{\beta} \varkappa_2\right)\right)^{t-\tau} + \alpha \eta_2.$$

Remark 3. Recall $z_t = (\theta_t - \theta^*, \varrho_t - \varrho^*)^\top$, then $\mathbb{E}[\|\theta_t - \theta^*\|_2^2] \leq \mathbb{E}[\|z_t\|_2^2]$, which implies the expected mean square error of $\theta_t - \theta^*$ is also upper-bounded by the result of Theorem 3. Furthermore, after a total computational cost of

$$\tau + \mathcal{O}\left(\frac{1}{\delta} \log \frac{1}{\delta}\right),$$

the Algorithm 1 outputs θ_t closes to the optimal solution θ^* as follows,

$$\mathbb{E}[\|\theta_t - \theta^*\|_2^2] \leq \mathcal{O}(\tau \delta).$$

Remark 4. Theorem 3 shows that the upper-bounded error consists of two different parts: the first error bound depends on both step-size and the size of samples, and this error decays geometrically as the number of iteration increases; while the second part is only determined by the step-sizes and it is independent of the number of iterations.

Algorithm	Step-size	Convergence Rate	TD Fixed Point
TD(0) (Nathaniel et al. 2015)	$\alpha_t = \mathcal{O}(t^{-\eta})$ $\eta \in (0, 1)$	$\mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$	$\Phi^\top \Xi(\gamma P^\mu - I)\Phi\theta^* = -b$
TD(0) (Dalal et al. 2018a)	$\sum_{t=1}^{\infty} \alpha_t = \infty$	$\mathcal{O}\left(\frac{1}{T^\eta}\right)$ $\eta \in (0, 1)$	$\Phi^\top \Xi(\gamma P^\mu - I)\Phi\theta^* = -b$
TD(0) (Lakshminarayanan et al. 2018)	Constant	$\mathcal{O}\left(\frac{1}{T}\right)$	$\Phi^\top \Xi(\gamma P^\mu - I)\Phi\theta^* = -b$
GTD(0) (Dalal et al. 2018b)	$\sum_{t=1}^{\infty} \alpha_t = \infty, \frac{\beta_t}{\alpha_t} \rightarrow 0$	$\mathcal{O}\left(\left(\frac{1}{T}\right)^{\frac{1-\kappa}{3}}\right)$ $\kappa \in (0, 1)$	$\Phi^\top \Xi(\gamma P^\pi - I)\Phi\theta^* = -b$
GTD(0)/GTD2/TDC (Dalal et al. 2020)	$\alpha_t = \frac{1}{t^{\eta_1}}, \beta_t = \frac{1}{t^{\eta_2}}$ $0 < \eta_2 < \eta_1 < 1$	$\mathcal{O}\left(\frac{1}{T^{\eta_1}}\right)$	$\Phi^\top \Xi(\gamma P^\pi - I)\Phi\theta^* = -b$
GTB(λ) (Touati et al. 2018)	$\alpha_t, \beta_t = \mathcal{O}\left(\frac{1}{t}\right)$	$\mathcal{O}\left(\frac{1}{T}\right)$	$\Phi^\top \Xi(I - \gamma\lambda P^\mu)^{-1}(\gamma P^\pi - I)\Phi\theta^* = -b$
SARSA (Zou et al. 2019)	$\alpha_t = \mathcal{O}\left(\frac{1}{t}\right)$	$\mathcal{O}\left(\frac{\log^3(T)}{T}\right)$	$\Phi^\top \Xi(\gamma P^\mu - I)\Phi\theta^* = -b$
TDC (Xu et al. 2019)	$\max\{\alpha_t \log\left(\frac{1}{\alpha_t}\right), \alpha_t\}$ $\leq \min\left\{\frac{\ \theta_0 - \theta^*\ _2}{2^{\varepsilon-1}}, C\right\}$	Linear	$\Phi^\top \Xi(\gamma P^\pi - I)\Phi\theta^* = -b$
TDC (Kaledin et al. 2020)	$\alpha_t = \frac{1}{t^v}, \beta_t = \frac{1}{t}$ $v \in (0, 1)$	$\mathcal{O}\left(\frac{1}{T}\right)$	$\Phi^\top \Xi(\gamma P^\pi - I)\Phi\theta^* = -b$
GES(λ) Theorem 2	Constant	$\mathcal{O}\left(\left(1 - \frac{C}{12}\right)^T\right)$	$\Phi^\top \Xi(I - \gamma\lambda P^\pi)^{-1}(\gamma P^\pi - I)\Phi\theta^* = -b$

Table 1: Comparison of GTD family algorithms over performance measurement $\mathbb{E}\|\theta_T - \theta^*\|_2^2$.

Related Works

In this section, we review existing finite-time performance of GTD algorithms over $\|\theta_T - \theta^*\|_2^2$.

Although the asymptotic analysis of GTD family has been established in (Sutton et al. 2009a,b; Maei 2011), which holds only in the limit as the number of iterations increases to infinity, and we can not get the information of convergence rate from asymptotic results. This is the main reason why we focus on the finite-time performance over $\|\theta_T - \theta^*\|_2^2$. It is noteworthy that Liu et al. (2015) firstly introduce *primal-dual gap error* to measure the convergence of GTD algorithm, we provide the discussion of finite-time primal-dual gap error analysis in Appendix E.

Nathaniel et al. (2015) proves that TD(0) (Sutton 1988) converges at $\mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$ with the step-size $\alpha_t = \mathcal{O}\left(\frac{1}{t^\eta}\right)$, $\eta \in (0, 1)$. Later, Dalal et al. (2018a) further explore the property of TD(0), they prove the convergence rate of TD(0) achieves $\mathcal{O}\left(e^{-\frac{\lambda}{2}T^{1-\eta}} + \frac{1}{T^\eta}\right)$, but it never reaches $\mathcal{O}\left(\frac{1}{T}\right)$, where $\eta \in (0, 1)$, λ is the minimum eigenvalue of the matrix $A^\top + A$. Lakshminarayanan et al., (2018) show TD(0) converges at $\mathcal{O}\left(\frac{1}{T}\right)$ with a more relaxed step-size than the works of (Nathaniel et al. 2015; Dalal et al. 2018a), it only requires a constant step-size. Recently, Dalal et al. (2018b) proves GTD(0) family algorithm (Sutton et al. 2009b) converges at $\mathcal{O}\left(\left(\frac{1}{T}\right)^{\frac{1-\kappa}{3}}\right)$, but never reach $\mathcal{O}\left(\frac{1}{T}\right)$, where $\kappa \in (0, 1)$. A very similar convergence rate appears in (Dalal et al. 2020), which considers TDC and GTD2. Touati et al. (2018) propose GTB(λ)/GRetrace(λ), they prove the convergence rate of GTB(λ)/GRetrace(λ) reaches $\mathcal{O}\left(\frac{1}{T}\right)$. Zou et al. (2019) show SARSA with linear function approximation converges

at the rate of $\mathcal{O}\left(\frac{\log^3(T)}{T}\right)$. Recently, Kaledin et al. (2020) further develop two timescale stochastic approximation with Markovian noise, and they show that TDC converges at a rate of $\mathcal{O}\left(\frac{1}{T}\right)$ if $\alpha_t = \frac{1}{t^v}$, $\beta_t = \frac{1}{t}$, $v \in (0, 1)$.

Theorem 2 illustrates our GES(λ) achieves a linear convergence rate, thus GES(λ) converges faster than all above gradient TD learning algorithms theoretically. Although Xu et al., (2019) prove TDC also converges at a linear convergence rate, they require a fussy blockwise diminishing step-size condition: $\max\{\alpha_t \log\left(\frac{1}{\alpha_t}\right), \alpha_t\} \leq \{\min\left\{\frac{\|\theta_0 - \theta^*\|_2}{2^{\varepsilon-1}}, C\right\}$, $\alpha_t = \mathcal{O}\left(\frac{1}{(t+1)^{\eta_1}}\right)$, $\beta_t = \mathcal{O}\left(\frac{1}{(t+1)^{\eta_2}}\right)$, $0 < \eta_2 < \eta_1 < 1$, where C is a constant. Apparently, our Theorem 2 requires a simpler condition of step-size than Xu et al., (2019). It is noteworthy that our Theorem 2 does not require an additional projection step (that is unnecessary in practice) that appears in (Xu et al. 2019).

Significantly, the finite-time performances of (Dalal et al. 2018a; Lakshminarayanan et al. 2018) requires an additional assumption that all the samples required to update the function parameters are i.i.d. In this paper, we remove this condition and achieve a better result than theirs.

Experiments

In this section, we test the capacity of GES(λ) for off-policy evaluation in three typical domains: MountainCar, Baird Star (Baird 1995), Two-State MDP (Touati et al. 2018). We compare GES(λ) with three state-of-art algorithms: GQ(λ) (Maei and Sutton 2010), ABQ(ζ) (Mahmood et al. 2017b), GTB(λ) (Touati et al. 2018) over two typical measurements: MSPBE and mean square error (MSE). We choose those

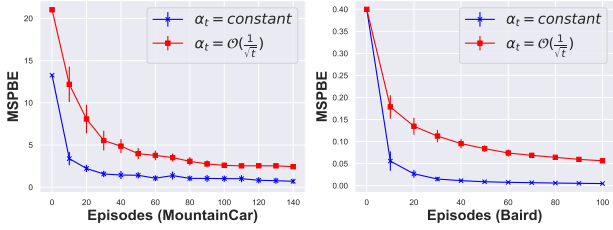


Figure 2: Comparison between a constant step-size and $\frac{1}{\sqrt{t}}$.

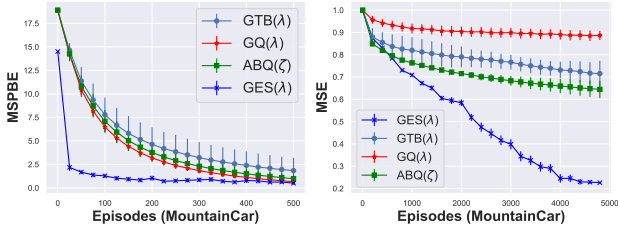


Figure 3: MSPBE and MSE comparison on MountainCar.

three algorithms as baselines since they are all learning via expected TD-error δ_t^{ES} , which is same as $\text{GES}(\lambda)$.

Feature Map and Parameters Recall the states and actions of MountainCar: $\mathcal{S} = \{\text{Velocity, Position}\} = [-0.07, 0.07] \times [-1.2, 0.6]$, $\mathcal{A} = \{\text{left, neutral, right}\}$. In this example, if $\text{Velocity} > 0$, we use behavior policy $\mu = (\frac{1}{100}, \frac{1}{100}, \frac{98}{100})$, $\pi = (\frac{1}{10}, \frac{1}{10}, \frac{8}{10})$; else $\mu = (\frac{98}{100}, \frac{1}{100}, \frac{1}{100})$, $\pi = (\frac{8}{10}, \frac{1}{10}, \frac{1}{10})$. Since the state space is continuous, we use an open *tile coding*³ software to extract feature of states. We set the number of tilings to be 4, and there are no white noise features. The performance is an average 5 runs, and each run contains 5000 episodes. As suggested by Sutton and Barto (2018), we set all the initial parameters to be 0, which is optimistic about causing extensive exploration.

The Baird Star is an episodic seven states MDP with two actions: dashed action and solid action. In this example, we set the behavior policy $\mu(\cdot|\text{dashed}) = \frac{6}{7}$, $\mu(\cdot|\text{solid}) = \frac{1}{7}$ and target policy $\pi(\cdot|\text{solid}) = 1$. We choose the feature map matrix as follows $\Phi = \begin{pmatrix} 2I_{7 \times 7} & \mathbb{1}_{7 \times 1} & \mathbf{0}_{7 \times 8} \\ \mathbf{0}_{7 \times 8} & 2I_{7 \times 7} & \mathbb{1}_{7 \times 1} \end{pmatrix}$, where $\mathbf{0}$ denotes a matrix whose elements are all 0, and $\mathbb{1}_{7 \times 1}$ denotes a vector whose elements are all 1. The dynamics of Two-State MDP is presented in Example 1. We set $\lambda = 0.99$, $\gamma = 0.99$ in all the experiments. The MSPBE/MSE distribution is computed over the combination of step-size, $(\alpha_t, \frac{\beta_t}{\alpha_t}) \in [0.1 \times 2^j | j = -10, -9, \dots, -1, 0]^2$.

Effect of Step-size Figure 2 shows the comparison of the empirical MSPBE performance between a constant step-size and the decay step-size $\mathcal{O}(\frac{1}{\sqrt{t}})$. Result of Figure 2 illustrates that the $\text{GES}(\lambda)$ with a proper constant step-size converges

³<http://incompleteideas.net/rlai.cs.ualberta.ca/RLAI/RLtoolkit/tilecoding.html>

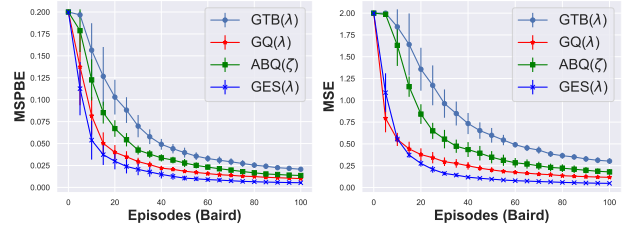


Figure 4: MSPBE and MSE comparison on Baird Star.

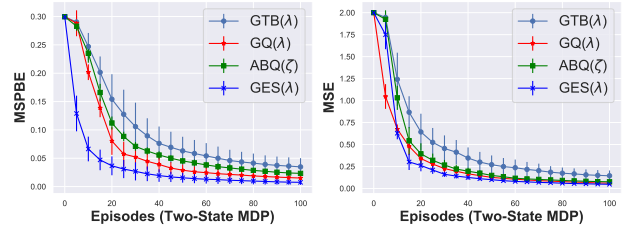


Figure 5: MSPBE and MSE comparison on Two-State MDP.

significantly faster than the learning with step-size $\mathcal{O}(\frac{1}{\sqrt{t}})$, which also supports Theorem 2: learning with a proper constant step-size can reach a very faster rate.

Comparison of Empirical MSPBE and MSE In this section, we use empirical $\text{MSPBE} = \frac{1}{2} \|\hat{b} + \hat{A}\theta\|_{\hat{M}^{-1}}^2$ to evaluate the performance, where we evaluate \hat{A} , \hat{b} , and \hat{M} according to their unbiased estimators by Monte Carlo method with 5000 episodes, and our implementation of MSPBE is inspired by (Touati et al. 2018). Besides, we also compare the performance over a common measurement empirical MSE: $\text{MSE} = \|\Phi\theta - q^\pi\|_{\Xi}^2 / \|q^\pi\|_{\Xi}^2$, where q^π is estimated by simulating the target policy π and averaging the discounted cumulative rewards over trajectories. The combination of step-size for MSE is the same as previous empirical MSPBE.

Results in Figure 3 to 5 show that our $\text{GES}(\lambda)$ learns faster with a better performance than $\text{GQ}(\lambda)$, $\text{ABQ}(\zeta)$ and $\text{GTB}(\lambda)$. Besides, $\text{GES}(\lambda)$ converges with a lower variance. In the Two-State MDP and Baird Star experiments, $\text{GES}(\lambda)$ outperforms the baselines slightly. This is because both Two-State MDP and Baird Star are relatively easy; many gradient TD learning could learn a convergent result. While the advantage of $\text{GES}(\lambda)$ over baselines becomes more significant in the MountainCar domain, which shows that $\text{GES}(\lambda)$ is more robust than baselines in the more difficult task.

Conclusion

We propose $\text{GES}(\lambda)$ that extends Expected Sarsa(λ) with linear function approximation. We prove $\text{GES}(\lambda)$ learns the optimal solution at a linear convergence rate, which is comparable to extensive GTD algorithms. The primal-dual gap error of $\text{GES}(\lambda)$ matches the best-known theoretical results, but we require a simpler condition of step-size. Finally, we conduct experiments to verify the effectiveness of $\text{GES}(\lambda)$.

References

- Baird, L. 1995. Residual algorithms: Reinforcement learning with function approximation. In *ICML*, 30–37.
- Bertsekas, D. P. 2011. Temporal difference methods for general projected equations. *IEEE Transactions on Automatic Control* 56(9): 2128–2139.
- Bertsekas, D. P. 2012. *Dynamic Programming and Optimal Control*, volume 2. Athena scientific Belmont, MA.
- Bhandari, J.; Russo, D.; Singal, R.; et al. 2018. A finite time analysis of temporal difference learning with linear function approximation. *COLT*.
- Dalal, G.; Szorenyi, B.; Thoppe, G.; and Mannor, S. 2018a. Finite sample analyses for td(0) with function approximation. In *AAAI2018*.
- Dalal, G.; Szorenyi, B.; Thoppe, G.; and Mannor, S. 2018b. Finite Sample Analysis of Two-Timescale Stochastic Approximation with Applications to Reinforcement Learning. In *COLT*.
- Dalal, G.; Szorenyi, B.; Thoppe, G.; et al. 2020. A Tale of Two-Timescale Reinforcement Learning with the Tightest Finite-Time Bound. *AAAI*.
- Dann, C.; Neumann, G.; Peters, J.; et al. 2014. Policy evaluation with temporal differences: A survey and comparison. *JMLR* 15(1): 809–883.
- Du, S. S.; Chen, J.; Li, L.; Xiao, L.; and Zhou, D. 2017. Stochastic variance reduction methods for policy evaluation. *ICML*.
- Du, S. S.; and Hu, W. 2019. Linear convergence of the primal-dual gradient method for convex-concave saddle point problems without strong convexity. *AISTATS*.
- Geist, M.; and Scherrer, B. 2014. Off-policy learning with eligibility traces: A survey. *JMLR*.
- Ghosh, D.; and Bellemare, M. G. 2020. Representations for Stable Off-Policy Reinforcement Learning. *ICML*.
- Gupta, H.; Srikant, R.; Ying, L.; et al. 2019. Finite-time performance bounds and adaptive learning rate selection for two time-scale reinforcement learning. In *NeurIPS*, 4704–4713.
- Kakade, S.; Shalev-Shwartz, S.; Tewari, A.; et al. 2009. On the duality of strong convexity and strong smoothness: Learning applications and matrix regularization.
- Kaledin, M.; Moulines, E.; Naumov, A.; Tadic, V.; and Wai, H.-T. 2020. Finite time analysis of linear two-timescale stochastic approximation with Markovian noise. *COLT*.
- Khalil, H. K.; and Grizzle, J. W. 2002. *Nonlinear systems*, volume 3. Prentice hall Upper Saddle River, NJ.
- Lagoudakis, M. G.; and Parr, R. 2003. Least-squares policy iteration. *JMLR* 4: 1107–1149.
- Lakshminarayanan, C.; Szepesvari, Csaba; et al. 2018. Linear Stochastic Approximation: How Far Does Constant Step-Size and Iterate Averaging Go? In *AISTATS*.
- Liu, B.; Liu, J.; Ghavamzadeh, M.; Mahadevan, S.; and Petrik, M. 2015. Finite-Sample Analysis of Proximal Gradient TD Algorithms. In *UAI*.
- Maei, H. R. 2011. *Gradient temporal-difference learning algorithms*. Ph.D. thesis, University of Alberta Edmonton, Alberta.
- Maei, H. R.; and Sutton, R. S. 2010. GQ(λ): A general gradient algorithm for temporal-difference prediction learning with eligibility traces. In *Proceedings of the third conference on artificial general intelligence*, volume 1, 91–96.
- Mahmood, A. R.; Yu, H.; Sutton, R. S.; et al. 2017b. Multi-step off-policy learning without importance sampling ratios.
- Nathaniel, K.; Prashanth, L.; Prashanth, L.; et al. 2015. On TD(0) with function approximation: Concentration bounds and a centered variant with exponential convergence. In *ICML*.
- Nemirovski, A.; Juditsky, A.; Lan, G.; and Shapiro, A. 2009. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization* 19(4): 1574–1609.
- Precup, D.; Sutton, R. S.; and Singh, S. P. 2000. Eligibility Traces for Off-Policy Policy Evaluation. In *International Conference on Machine Learning*, 759–766. Citeseer.
- Rockafellar, R. T. 1970. *Convex analysis*. 28. Princeton university press.
- Rubinstein, R. Y.; and Kroese, D. P. 2016. *Simulation and the Monte Carlo method*, volume 10. John Wiley & Sons.
- Srikant; and Ying, L. 2019. Finite-Time Error Bounds For Linear Stochastic Approximation and TD Learning. In *COLT*.
- Sutton, R. S. 1988. Learning to predict by the methods of temporal differences. *Machine learning* 3(1): 9–44.
- Sutton, R. S.; and Barto, A. G. 1998. *Reinforcement learning: An introduction*. MIT press Cambridge.
- Sutton, R. S.; and Barto, A. G. 2018. *Reinforcement learning: An introduction*. MIT press.
- Sutton, R. S.; Maei, H. R.; Precup, D.; Bhatnagar, S.; Silver, D.; Szepesvári, C.; and Wiewiora, E. 2009a. Fast gradient-descent methods for temporal-difference learning with linear function approximation. In *ICML*.
- Sutton, R. S.; Maei, H. R.; Szepesvári, C.; et al. 2009b. A Convergent $O(n)$ Temporal-difference Algorithm for Off-policy Learning with Linear Function Approximation. In *NeurIPS*.
- Sutton, R. S.; Mahmood, A. R.; White, M.; et al. 2016. An emphatic approach to the problem of off-policy temporal-difference learning. *JMLR* 17(1): 2603–2631.
- Thomas, P. S. 2015. *Safe reinforcement learning*. Ph.D. thesis, University of Massachusetts Libraries.
- Touati, A.; Bacon, P.-L.; Precup, D.; and Vincent, P. 2018. Convergent Tree-Backup and Retrace with Function Approximation. In *ICML*.

Voloshin, C.; Le, H. M.; Jiang, N.; and Yue, Y. 2019. Empirical Study of Off-Policy Policy Evaluation for Reinforcement Learning. *arXiv preprint arXiv:1911.06854* .

Wang, G.; Li, B.; Giannakis, G. B.; et al. 2019. A multistep Lyapunov approach for finite-time analysis of biased stochastic approximation. *arXiv preprint arXiv:1909.04299* .

Wang, Y.; Chen, W.; Liu, Y.; Ma, Z.-M.; and Liu, T.-Y. 2017. Finite Sample Analysis of the GTD Policy Evaluation Algorithms in Markov Setting. In *Advances in Neural Information Processing Systems(NeurIPS)*.

Xu, T.; Zou, S.; Liang, Y.; et al. 2019. Two time-scale off-policy TD learning: Non-asymptotic analysis over Markovian samples. In *NeurIPS*.

Zou, S.; Xu, T.; Liang, Y.; et al. 2019. Finite-sample analysis for SARSA with linear function approximation. In *NeurIPS*.

Appendix A: Importance Sampling and λ -Operator

For the discussion of off-policy learning, we need the background of importance sampling. Thus, the basic common conclusion about importance sampling (IS) and pre-decision importance sampling (PDIS) (Precup, Sutton, and Singh 2000) is necessary.

A.1. Off-Policy Learning via Importance Sampling

Usually, we require that every action taken by π is also taken by μ , which is often called *coverage* (Sutton and Barto 2018) in reinforcement learning.

Assumption 4 (Coverage). $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$, we require that $\pi(a|s) > 0 \Rightarrow \mu(a|s) > 0$.

The difficulty of off-policy roots in the discrepancy between target policy π and behavior policy μ —we want to learn the target policy while we only get the data generated by behavior policy. One technique to hand this discrepancy is *importance sampling* (IS) (Rubinstein and Kroese 2016). Let $\tau_t^h = \{S_t, A_t, R_{t+1}\}_{t \geq 0}^h$ be a trajectory with finite horizon $h < \infty$. Let $\rho_{t:k} = \prod_{i=t}^k \rho_i$ denote the *cumulated importance sampling ratio*, where $\rho_i = \frac{\pi(A_i|S_i)}{\mu(A_i|S_i)}$ and $k \leq h$. Let $G_t^h = \sum_{k=0}^{h-t-1} \gamma^k R_{t+k+1}$, under Assumption 4 the IS estimator $G_t^{\text{IS}} = \rho_{t:h-1} G_t^h$ is a unbiased estimation of q^π . However, it is known that IS estimator suffers from large variance of the product $\rho_{t:h-1}$ (Sutton and Barto 1998). Pre-decision importance sampling (PDIS) (Precup, Sutton, and Singh 2000) $G_t^{\text{PDIS}} = \sum_{k=0}^{h-t-1} \gamma^k \rho_{t:t+k} R_{t+k+1}$ is a practical variance reduction method without introducing bias, i.e. $\mathbb{E}_\mu[G_t^{\text{PDIS}} | S_t = s, A_t = a] = q^\pi(s, a)$.

$$\begin{aligned} \mathbb{E}_\mu[\rho_{t:h-1} G_t^h] &= \mathbb{E}_\mu[\underbrace{\rho_{t:h-1} R_{t+1} + \rho_{t:h-1} \gamma R_{t+2} + \cdots + \rho_{t:h-1} \gamma^{h-t-1} R_h}_{\stackrel{\text{def}}{=} G_t^{\text{IS}} \text{ IS-return}}] \\ &= \mathbb{E}_\mu[\underbrace{\rho_t R_{t+1} + \rho_{t:t+1} \gamma R_{t+2} + \cdots + \rho_{t:h-1} \gamma^{h-t-1} R_h}_{\stackrel{\text{def}}{=} G_t^{\text{PDIS}} \text{ PDIS-return}}] = \mathbb{E}_\mu[\sum_{k=0}^{h-t-1} \gamma^k \rho_{t:t+k} R_{t+k+1}]. \end{aligned}$$

For the equation $\mathbb{E}_\mu[G_t^{\text{IS}}] = \mathbb{E}_\mu[G_t^{\text{PDIS}}]$, please see (Precup, Sutton, and Singh 2000) or section 5.9 in (Sutton and Barto 2018).

Lemma 2 (Section 3.10, (Thomas 2015); Section 5.9, (Sutton and Barto 2018)). *Let $\tau_t^h = \{S_k, A_k, R_{k+1}\}_{k=t}^h$ be the trajectory generated by behavior policy μ , for a given policy π and under Assumption 4, the following holds,*

$$\mathbb{E}_\mu[\rho_{t:h-1} R_{t+k}] = \mathbb{E}_\mu[\rho_{t:t+k-1} R_{t+k}]. \quad (36)$$

Lemma 2 implies that for any time $t+k$ ($k \geq 0$), the importance sampling factors after $t+k$ have no effect in the expectation, thus the following holds: for all $k \geq 0$,

$$\mathbb{E}_\mu[\rho_{t:h-1} R_{t+k}] = \mathbb{E}_\mu[\rho_{t:t+k-1} R_{t+k}] = \mathbb{E}_\pi[R_{t+k}]. \quad (37)$$

A.2. Derivation of Eq.(7)

Proof.

$$\begin{aligned} q + \mathbb{E}_\mu[\sum_{k=t}^{\infty} (\lambda\gamma)^{k-t} \delta_k^{\text{ES}} \rho_{t+1:k}] &\stackrel{(37)}{=} q + \mathbb{E}_\pi[\sum_{k=t}^{\infty} (\lambda\gamma)^{k-t} \delta_k^{\text{ES}}] \\ &= q + (I - \lambda\gamma P^\pi)^{-1} (\mathcal{B}^\pi q - q), \end{aligned} \quad (38)$$

Eq. (38) is a common result in RL, for the details of

$$\mathbb{E}_\pi[\sum_{k=t}^{\infty} (\lambda\gamma)^{k-t} \delta_k^{\text{ES}}] = (I - \lambda\gamma P^\pi)^{-1} (\mathcal{B}^\pi q - q),$$

please refer to (Geist and Scherrer 2014) or Section 6.3.9 in (Bertsekas 2012). \square

Appendix B: Proof of Theorem 1

Theorem 1 (Stability Criteria) *Under Assumption 1, the off-line update (9) is stable if and only if the eigenvalues of the matrix A (10) have negative real components, i.e.,*

$$\text{Spec}(A) \subset \mathbb{C}_-. \quad (39)$$

Before we present the details of its proof, we need some notations of the matrix. Recall $\text{Spec}(A)$ are the eigenvalues of the matrix $A \in \mathbb{C}^{p \times p}$, we use $\rho(A)$ to denote its spectral radius of the matrix A , i.e.,

$$\rho(A) = \sup_{\lambda \in \text{Spec}(A)} \{|\lambda|\}.$$

Proof. Recall $A = \Phi^\top \Xi (I - \gamma \lambda P^\pi)^{-1} (\gamma P^\pi - I) \Phi$, and θ^* satisfies

$$A\theta^* + b = 0,$$

which implies

$$\theta_{t+1} - \theta^* = \theta_t - \alpha(A\theta_t + b) - \theta^* = (I + \alpha A)(\theta_t - \theta^*).$$

Applying above result of recurrently, we have

$$\theta_t - \theta^* = (I + \alpha A)^t (\theta_0 - \theta^*), \quad (40)$$

which implies if the iteration (40) converges, if and only if

$$\rho(I + \alpha A) < 1.$$

Furthermore, if the iteration (40) is stable, if and only if $\rho(I + \alpha A) < 1$. If there is a step-size $\alpha > 0$ such that $\rho(I + \alpha A) < 1$, then we have $\text{Spec}(A) \subset \mathbb{C}_-$. Conversely, if $\text{Spec}(A) \subset \mathbb{C}_-$, and let

$$\alpha = \inf_{\lambda \in \text{Spec}(A)} \left\{ \frac{\text{Re}(\lambda)}{|\lambda|} \right\},$$

then $\rho(I + \alpha A) < 1$. □

Appendix C

C.1: Proof of Eq.(17)

For a given policy π , $Q_\theta = \Phi\theta$, then by the definition of MSPBE objection function, we have,

$$\begin{aligned} \text{MSPBE}(\theta, \lambda) &= \|Q_\theta - \Pi \mathcal{B}_\lambda^\pi Q_\theta\|_\Xi^2 \\ &= \|\Pi Q_\theta - \Pi \mathcal{B}_\lambda^\pi Q_\theta\|_\Xi^2 \\ &= \|\Phi^T \Xi (Q_\theta - \mathcal{B}_\lambda^\pi Q_\theta)\|_{(\Phi^T \Xi \Phi)^{-1}}^2 \\ &= \|\Phi^T \Xi (I - \lambda \gamma P^\pi)^{-1} (\Phi\theta - \gamma P^\pi \Phi\theta - R^\pi)\|_{(\Phi^T \Xi \Phi)^{-1}}^2 \\ &= \|\Phi^T \Xi (I - \lambda \gamma P^\pi)^{-1} ((I - \gamma P^\pi) \Phi\theta - R^\pi)\|_{(\Phi^T \Xi \Phi)^{-1}}^2 \\ &= \|b + A\theta\|_{(\Phi^T \Xi \Phi)^{-1}}^2, \end{aligned} \quad (41)$$

where $A = \Phi^T \Xi (I - \lambda \gamma P^\pi)^{-1} (\gamma P^\pi - I) \Phi$, $b = \Phi^T \Xi (I - \lambda \gamma P^\pi)^{-1} r$.

C.2: Proof of Proposition 1

Proposition 1 If (θ_*, ω_*) is the solution of the problem (19), then θ_* is the solution of original problem (17), i.e.,

$$\theta_* = \arg \min_{\theta} \text{MSPBE}(\theta, \lambda).$$

Proof. If $\omega_* = \arg \max_{\omega} (A\theta + b)^\top \omega - \frac{1}{2} \|\omega\|_M^2$, then $\omega_* = M^{-1}(A\theta + b)$. Taking ω_* into (19), then (19) is reduced to $\min_{\theta} \frac{1}{2} \|A\theta + b\|_{M^{-1}}^2$. Now, let θ_* be the solution of (19), then we have

$$\theta_* = \min_{\theta} \frac{1}{2} \|A\theta + b\|_{M^{-1}}^2 = \arg \min_{\theta} \text{MSPBE}(\theta, \lambda).$$

□

C.3: Proof of Theorem 2

Theorem 2 Let $\{(\theta_t, \omega_t)\}_{t=0}^T$ be generated by Algorithm1. Recall $g(\omega)$ has been defined in (18) and θ^* is the optimal solution of (17). Let $\nu = \frac{2\kappa^2(A)\kappa(M)\sigma_{\max}(A)}{\sigma_{\min}(M)}$, $\Delta_{\theta_t} = \|\theta_t - \theta^*\|_2^2$, $\Delta_{\omega_t} = \|\omega_t - \nabla g^*(A\theta_t)\|_2^2$, and $D_t = \nu \Delta_{\theta_t} + \Delta_{\omega_t}$. Under Assumption 1, and we assume $\text{rank}(\Phi) = p$. Let

$$\alpha = \frac{\sigma_{\min}(M)}{(\sigma_{\max}(M) + \sigma_{\min}(M))(\frac{\sigma_{\max}^2(A)}{\sigma_{\min}(M)} + \nu \sigma_{\max}(A))}, \beta = \frac{2}{\sigma_{\max}(M) + \sigma_{\min}(M)},$$

under Assumption 1, we have

$$\mathbb{E}[D_{t+1}] \leq \left(1 - \frac{1}{12 \kappa^3(M) \kappa^4(A)}\right) \mathbb{E}[D_t]. \quad (42)$$

Furthermore, we have

$$\mathbb{E}[\|\theta_t - \theta^*\|_2^2] \leq \frac{1}{\nu} \left(1 - \frac{1}{12 \kappa^3(M) \kappa^4(A)}\right)^t \mathbb{E}[D_0].$$

Proof. The proof is inspired by a general analysis that appears in (Du and Hu 2019); we refer the reader to that reference for further technical details.

Firstly, we apply Theorem 3.1 of (Du and Hu 2019) to achieve (42), which requires us to check $g(\omega)$ is a Lipschitz-smooth and strongly convex function. Recall $g(\omega) = \frac{1}{2} \|\omega\|_M^2 - b^\top \omega$, for any ω_1, ω_2 , we have

$$\|\nabla g(\omega_1) - \nabla g(\omega_2)\|_2 = \|M(\omega_1 - \omega_2)\|_2 \leq \|M\|_{\text{op}} \|\omega_1 - \omega_2\|_2 \stackrel{\text{(a)}}{=} \sigma_{\max}(M) \|\omega_1 - \omega_2\|_2, \quad (43)$$

where the last equation (a) of (43) holds since $M = \mathbb{E}_\mu[\phi_t \phi_t^\top] = \Phi^\top \Xi \Phi$ is a positive symmetric matrix if $\text{rank}(\Phi) = p$, then

$$\|M\|_{\text{op}} = \max\{\lambda_1, \lambda_2, \dots, \lambda_p\} = \sigma_{\max}(M).$$

Eq.(43) implies $g(\omega)$ is $\sigma_{\max}(M)$ -smooth function. Under Assumption 1, $\Xi \succ 0$; recall $\text{rank}(\Phi) = p$, so $M = \Phi^\top \Xi \Phi \succ 0$, then $\nabla^2 g(\omega) = M$ implies

$$\nabla^2 g(\omega) - \sigma_{\min}(M)I \succ 0,$$

thus $g(\omega)$ is a $\sigma_{\min}(M)$ -strongly convex function. Let step-size satisfy the following condition

$$\alpha = \frac{\sigma_{\min}(M)}{(\sigma_{\max}(M) + \sigma_{\min}(M))(\frac{\sigma_{\max}^2(A)}{\sigma_{\min}(M)} + \nu\sigma_{\max}(A))}, \beta = \frac{2}{\sigma_{\max}(M) + \sigma_{\min}(M)},$$

then according to Theorem 3.1 of (Du and Hu 2019), we have

$$\mathbb{E}[D_{t+1}] \leq \left(1 - \frac{1}{12} \frac{1}{\kappa^3(M)\kappa^4(A)}\right) \mathbb{E}[D_t].$$

Furthermore, the fact $D_t = \nu\Delta_{\theta_t} + \Delta_{\omega_t} \geq \nu\Delta_{\theta_t}$ implies

$$\mathbb{E}[\|\theta_t - \theta_\star\|_2^2] \leq \frac{\mathbb{E}[D_t]}{\nu} \stackrel{(42)}{\leq} \frac{1}{\nu} \left(1 - \frac{1}{12} \frac{1}{\kappa^3(M)\kappa^4(A)}\right)^t \mathbb{E}[D_0].$$

Therefore the proof is completed. \square

Corollary 1. Consider all the conditions and notations in Theorem 2, the output of Algorithm 1 closes to $(\theta^\star, \omega^\star)$ as follows,

$$\mathbb{E}[\|\theta_t - \theta^\star\|_2^2] \leq \delta^2, \quad \mathbb{E}[\|\omega_t - \omega^\star\|_2^2] \leq \delta^2,$$

if after a computational cost of

$$\mathcal{O}\left(\max\left\{1, \frac{\lambda_{\max}(A)}{\lambda_{\max}(M)\nu}\right\} \left(1 - \frac{1}{12\kappa^3(M)\kappa^4(A)}\right) \log\left(\frac{1}{\delta}\right)\right).$$

Proof. Furthermore, recall $\omega^\star = (\nabla g)^{-1}(A\theta^\star) = \nabla g^\star(A\theta^\star)$, then we have

$$\begin{aligned} \|\omega_t - \omega^\star\|_2 &\leq \|\omega_t - \nabla g^\star(A\theta_t)\|_2 + \|\nabla g^\star(A\theta_t) - \omega^\star\|_2 \\ &\leq \Delta_{\omega_t} + \|\nabla g^\star(A\theta_t) - g^\star(A\theta^\star)\|_2 \\ &\leq \Delta_{\omega_t} + \frac{\lambda_{\max}(A)}{\lambda_{\max}(M)} \Delta_{\theta_t} \leq \max\left\{1, \frac{\lambda_{\max}(A)}{\lambda_{\max}(M)\nu}\right\} D_t. \end{aligned}$$

Since $\mathbb{E}[D_{t+1}] \leq \left(1 - \frac{1}{12} \frac{1}{\kappa^3(M)\kappa^4(A)}\right) \mathbb{E}[D_t]$, then we have

$$\mathbb{E}[\|\omega_t - \omega^\star\|_2] \leq \max\left\{1, \frac{\lambda_{\max}(A)}{\lambda_{\max}(M)\nu}\right\} \mathbb{E}[D_{t+1}] \leq \max\left\{1, \frac{\lambda_{\max}(A)}{\lambda_{\max}(M)\nu}\right\} \left(1 - \frac{1}{12} \frac{1}{\kappa^3(M)\kappa^4(A)}\right)^t \mathbb{E}[D_0].$$

If, $\mathbb{E}[\|\omega_t - \omega^\star\|_2] \leq \delta$, then the times of update at least

$$\mathcal{O}\left(\max\left\{1, \frac{\lambda_{\max}(A)}{\lambda_{\max}(M)\nu}\right\} \left(1 - \frac{1}{12\kappa^3(M)\kappa^4(A)}\right) \log\left(\frac{1}{\delta}\right)\right).$$

This conclude the result of Remark 1. \square

Appendix D: Proof of Theorem 3

D.1: Some Preliminaries

Assumption 2 (Boundedness of Feature Map, Reward) *The features $\{\phi_t\}_{t \geq 0}$ is uniformly bounded by ϕ_{\max} . The reward function is uniformly bounded by R_{\max} . The importance sampling ρ_t is uniformly bounded by ρ_{\max} .*

Remark 5. *Assumption 2 implies the boundedness of \hat{A}_t , \hat{M}_t and \hat{b}_t . In fact,*

$$\|\hat{M}_t\|_{op}^2 = \|\phi_t \phi_t^\top\|_{op}^2 \leq (p\sqrt{p}\phi_{\max}^2)^2 =: C_M^2,$$

where we use a basic fact for matrix: for each $A \in \mathbb{R}^{p \times p}$, we have $\|A\|_{op} \leq p\sqrt{p}\|A\|_\infty$. The boundedness of \hat{M}_t also imply the boundedness of M , i.e.,

$$\|M\|_{op} \leq C_M.$$

Furthermore, since $M = \Phi^\top \Xi \Phi$ and $\text{rank}(\Phi) = p$, thus, both M and M^{-1} are all positive definite matrix, then, there exists a positive scalar $C_{M^{-1}}$ such that

$$\|M^{-1}\|_{op} \leq C_{M^{-1}}.$$

Recall $e_t = \lambda\gamma\rho_t e_{t-1} + \phi_t$, then

$$\|e_t\|_2^2 = \left\| \sum_{k=0}^t (\gamma\lambda)^{t-k} \rho_{k+1:t} \phi_k \right\|_2^2 \leq \left(\frac{\phi_{\max}}{1 - \gamma\lambda\rho_{\max}} \right)^2 =: C_e^2.$$

Since then,

$$\begin{aligned} \|\hat{b}_t\|_2^2 &= \|R_{t+1}e_t\|_2^2 \leq R_{\max}^2 C_e^2 =: C_b^2, \\ \|\hat{A}_t\|_2^2 &= \|e_t(\gamma\mathbb{E}_\pi[\phi(S_{t+1}, \cdot)] - \phi_t)^\top\|_2^2 \leq (\gamma + 1)^2 C_e^2 \phi_{\max}^2 =: C_A^2. \end{aligned}$$

Above results also imply the boundedness of A and b , i.e.,

$$\|A\|_{op} \leq C_A, \quad \text{and} \quad \|b\|_2 \leq C_b.$$

Lemma 3 ((Gupta et al. 2019)). *Consider the matrix A as follows*

$$A(x) = \frac{1}{a_1 + a_2} \begin{pmatrix} a_2 & -a_1 a_2 \\ -a_1 a_2 & \frac{1}{x} a_1 - x y a_1 \end{pmatrix},$$

where $a_1, a_2, c, y > 0$. Then we have

$$\lambda_{\min}(A(x)) \geq \kappa_1 - \kappa_2 x,$$

where $\kappa_1 = \frac{a_2}{a_1 + a_2}$ and κ_2 is a constant that only depends on a_1, a_2 and x .

D.2: Proof of Lemma 1

Lemma 1 *Under Assumption 1-3, there exists a positive scalar τ such that: $t \geq \tau$,*

$$\mathbb{E}[L(z_{t+1})] - \mathbb{E}[L(z_t)] \leq -\alpha \left(\frac{1}{2} \varkappa_1 - \frac{\alpha}{\beta} \varkappa_2 \right) \mathbb{E}[L(z_t)] + 2\alpha^2 \zeta^2 \lambda_{\max}(Q) \tilde{c}_b^2,$$

where $\varkappa_1, \varkappa_2, \zeta, \tilde{c}_b$ are constants that we will concretize them later.

We define

$$\begin{aligned} \hat{H}_t &= -\hat{A}_t^\top M^{-1} A, & H &= -A^\top M^{-1} A, \\ \hat{L}_t &= \hat{A}_t + \hat{M}_t M^{-1} A, & L &= 2A. \end{aligned}$$

By the boundedness of the previous term, we have

$$\begin{aligned} \|\hat{H}_t\|_{op} &\leq \|\hat{A}_t^\top M^{-1} A\|_{op} \leq C_A^2 C_{M^{-1}}, \\ \|\hat{L}_t\|_{op} &\leq \|\hat{A}_t\|_{op} + \|\hat{M}_t M^{-1} A\|_{op} \leq C_A + C_M C_{M^{-1}} C_A, \quad \|L\|_{op} \leq 2C_A. \end{aligned}$$

Furthermore, let Q_1 and Q_2 be the solutions of the following equations

$$\begin{cases} H^\top Q_1 + Q_1 H &= -I \\ M^\top Q_2 + Q_2 M &= I, \end{cases} \quad (44)$$

since both M and L are Hurwitz matrix, then solution of (44) always exists. Furthermore, we define a matrix P as follows,

$$Q = \begin{pmatrix} \frac{\|Q_1 A^\top\|_{op}}{\|Q_1 A^\top\|_{op} + \|Q_2 M^{-1} A L\|_{op}} Q_1 & 0 \\ 0 & \frac{\|Q_2 M^{-1} A L\|_{op}}{\|Q_1 A^\top\|_{op} + \|Q_2 M^{-1} A L\|_{op}} Q_2 \end{pmatrix}. \quad (45)$$

We define

$$\varrho_t = \omega_t - M^{-1} A \theta_t, \quad z_t = \begin{pmatrix} \theta_t - \theta^* \\ \varrho_t - \varrho^* \end{pmatrix}. \quad (46)$$

where $\varrho^* = \omega^* - M^{-1} A \theta^*$. Let $L(z_t)$ be the following Lyapunov function,

$$L(z_t) = z_t^\top Q z_t. \quad (47)$$

Proof. Recall the update (23) and (24):

$$\omega_{t+1} = \omega_t + \beta(\hat{A}_t \theta_t + \hat{b}_t - \hat{M}_t \omega_t), \quad \theta_{t+1} = \theta_t - \alpha \hat{A}_t^\top \omega_t,$$

we can rewrite z_t (46) as the following recursion:

$$z_{t+1} = z_t + \alpha(\tilde{G}_t z_t + \tilde{g}_t), \quad (48)$$

where

$$\tilde{G}_t = \begin{pmatrix} \hat{H}_t, & -\hat{A}_t^\top \\ -M^{-1} A \hat{H}_t + \frac{\beta}{\alpha} \hat{L}_t & M^{-1} A \hat{A}_t^\top - \frac{\beta}{\alpha} \hat{M}_t \end{pmatrix}, \quad \tilde{g}_t = \begin{pmatrix} 0 \\ \frac{\beta}{\alpha} \hat{b}_t \end{pmatrix}. \quad (49)$$

Let $\tilde{G}_\infty =: \lim_{t \rightarrow \infty} \tilde{G}_t$ and $\tilde{g}_\infty =: \lim_{t \rightarrow \infty} \tilde{g}_t$. Furthermore, we define

$$G_\infty = \lim_{t \rightarrow \infty} \mathbb{E}[\tilde{G}_t] = \begin{pmatrix} H, & A^\top \\ -M^{-1} A H + \frac{\beta}{\alpha} L & M^{-1} A A^\top - \frac{\beta}{\alpha} M \end{pmatrix}, \quad g_\infty = \begin{pmatrix} 0 \\ -\frac{\beta}{\alpha} b \end{pmatrix}.$$

Boundedness of $\tilde{G}_t, \tilde{G}_\infty, \tilde{g}_t, \tilde{g}_\infty$.

$$\|\tilde{G}_t\|_{op} \leq \|\hat{H}_t\|_{op} + \|\hat{A}_t\|_{op} + \|M^{-1} A \hat{H}_t\|_{op} + \left\| \frac{\beta}{\alpha} \hat{L}_t \right\|_{op} + \|M^{-1} A \hat{A}_t\|_{op} + \left\| \frac{\beta}{\alpha} \hat{M}_t \right\|_{op}.$$

Recall the boundedness of $\hat{H}_t, \hat{A}_t, M^{-1}, A, \hat{L}_t$ and \hat{M}_t , then the following holds

$$\|\tilde{G}_t\|_{op} \leq \underbrace{2C_A^2 C_{M^{-1}} + C_A + C_A^3 C_{M^{-1}}}_{=: C_1} + \frac{\beta}{\alpha} \underbrace{(C_A + C_M C_{M^{-1}} C_A + C_M)}_{=: C_2} = C_1 + \frac{\beta}{\alpha} C_2 =: \zeta. \quad (50)$$

Furthermore, we have

$$\|\tilde{G}_\infty\|_{op} \leq \zeta, \quad \|\tilde{g}_t\|_2 \leq \frac{\beta}{\alpha} c_b = C_2 \frac{\beta}{\alpha} \frac{c_b}{C_2} \leq \frac{c_b}{C_2} \zeta =: \tilde{c}_b \zeta, \quad \|\tilde{g}_\infty\|_2 = 0.$$

Recall (48), we have

$$\|z_{t+1} - z_t\|_2 \leq \alpha \|\tilde{G}_t z_t + \tilde{g}_t\|_2 \leq \alpha \zeta (\|z_t\|_2 + \tilde{c}_b). \quad (51)$$

Furthermore, according to the same analysis of Lemma 3 in (Srikant and Ying 2019), we have

$$\begin{aligned} \|z_\tau - z_0\|_2 &\leq 2\alpha\zeta\tau\|z_0\|_2 + 2\alpha\zeta\tau\tilde{c}_b, \\ \|z_\tau - z_0\|_2 &\leq 4\alpha\zeta\tau\|z_\tau\|_2 + 4\alpha\zeta\tau\tilde{c}_b, \\ \|z_\tau - z_0\|_2^2 &\leq 32\alpha^2\zeta^2\tau^2\|z_\tau\|_2^2 + 32\alpha^2\zeta^2\tau^2\tilde{c}_b^2. \end{aligned}$$

Recall (51), we have

$$|(z_{t+1} - z_t)^\top P(z_{t+1} - z_t)| \leq \lambda_{\max}(Q) \|z_{t+1} - z_t\|_2^2 \leq \lambda_{\max}(Q) (\alpha\zeta(\|z_t\|_2 + \tilde{c}_b))^2 \leq 2\alpha^2\zeta^2\lambda_{\max}(Q) (\|z_t\|_2^2 + \tilde{c}_b^2),$$

where the last inequality holds since $(a + b)^2 \leq 2(a^2 + b^2)$. Let

$$X_t =: \{S_0, A_0, S_1, A_1, \dots, S_t, A_t\},$$

after some careful calculations, we have: for all $t \geq \tau$,

$$\begin{aligned} & \left| \mathbb{E} \left[z_t^\top \left(PG_\infty z_t - \frac{1}{\alpha} (z_{t+1} - z_t) \right) \middle| z_{t-\tau}, X_{t-\tau} \right] \right| \\ &= \left| \mathbb{E} \left[z_t^\top \left(PG_\infty z_t - \tilde{G}_t z_t \right) \middle| z_{t-\tau}, X_{t-\tau} \right] \right| \leq \zeta (1 + \lambda_{\max}(Q)) \mathbb{E}[\|z_t\|_2^2 | z_{t-\tau}, X_{t-\tau}] \end{aligned} \quad (52)$$

For $t \geq \tau$, we have

$$\begin{aligned} & \mathbb{E}[L(z_{t+1}) - L(z_t) | z_{t-\tau}, X_{t-\tau}] \\ &= \mathbb{E}[z_{t+1}^\top P z_{t+1} - z_t^\top P z_t | z_{t-\tau}, X_{t-\tau}] \\ &= \mathbb{E}[(z_{t+1} - z_t)^\top P (z_{t+1} - z_t) + 2z_t^\top P (z_{t+1} - z_t) | z_{t-\tau}, X_{t-\tau}] \\ &= \mathbb{E}[\underbrace{(z_{t+1} - z_t)^\top P (z_{t+1} - z_t)}_{\leq 2\alpha^2 \zeta^2 \lambda_{\max}(Q) (\|z_t\|_2^2 + \tilde{c}_b^2)} + 2z_t^\top P (z_{t+1} - z_t - \alpha G_\infty z_t) + 2\alpha z_t^\top P G_\infty z_t | z_{t-\tau}, X_{t-\tau}] \\ &\leq \mathbb{E}[2\alpha^2 \zeta^2 \lambda_{\max}(Q) (\|z_t\|_2^2 + \tilde{c}_b^2) + 2z_t^\top P (z_{t+1} - z_t - \alpha G_\infty z_t) + 2\alpha z_t^\top P G_\infty z_t | z_{t-\tau}, X_{t-\tau}]. \end{aligned} \quad (53)$$

Furthermore, since $\mathbb{E}[z_t^\top P G_\infty z_t | z_{t-\tau}, X_{t-\tau}] \leq -\lambda_{\min}(\Sigma) \mathbb{E}[\|z_t\|_2^2 | z_{t-\tau}, X_{t-\tau}]$, where $\lambda_{\min}(\Sigma)$ is the smallest eigenvalue of the the following matrix Σ :

$$\Sigma = \begin{pmatrix} \frac{\xi_2}{\xi_1 + \xi_2} & -\frac{\xi_1 \xi_2}{\xi_1 + \xi_2} \\ -\frac{\xi_1 \xi_2}{\xi_1 + \xi_2} & \xi_1 \left(\frac{\beta}{\alpha} - 2\|Q_2 A^\top M^{-1} A\|_{op} \right) \end{pmatrix},$$

and $\xi_1 = 2\|Q_1 A^\top\|_{op}$, and $\xi_2 = 2\|Q_2 M^{-1} A L\|_{op}$. Recall the result of Lemma 3, we have

$$\lambda_{\min}(\Sigma) \geq \frac{\xi_2}{\xi_1 + \xi_2} - \frac{\alpha}{\beta} \varkappa_1 =: \varkappa_1 - \frac{\alpha}{\beta} \varkappa_2, \quad (54)$$

where $\varkappa_1 = \frac{\xi_2}{\xi_1 + \xi_2} = \frac{\|Q_1 A^\top\|_{op}}{\|Q_2 M^{-1} A L\|_{op} + \|Q_1 A^\top\|_{op}}$, \varkappa_2 is a constant depends on ξ_1, ξ_2 and $\frac{\alpha}{\beta}$.

From the result of (52) and (53), we have

$$\begin{aligned} & \mathbb{E}[L(z_{t+1}) - L(z_t) | z_{t-\tau}, X_{t-\tau}] \\ &\leq 2\alpha^2 \zeta^2 \lambda_{\max}(Q) \mathbb{E}[\|z_t\|_2^2 | z_{t-\tau}, X_{t-\tau}] + 2\alpha^2 \zeta^2 \lambda_{\max}(Q) \tilde{c}_b^2 + 2\alpha \mathbb{E}[z_t^\top P G_\infty z_t | z_{t-\tau}, X_{t-\tau}] \\ &\quad + 2 \underbrace{\mathbb{E}[z_t^\top P (z_{t+1} - z_t - \alpha G_\infty z_t) | z_{t-\tau}, X_{t-\tau}]}_{= 2\alpha \mathbb{E}[z_t^\top P (\frac{z_{t+1} - z_t}{\alpha} - G_\infty z_t) | z_{t-\tau}, X_{t-\tau}]} \\ &\stackrel{(52)}{\leq} 2\alpha^2 \zeta^2 \lambda_{\max}(Q) \mathbb{E}[\|z_t\|_2^2 | z_{t-\tau}, X_{t-\tau}] + 2\alpha^2 \zeta^2 \lambda_{\max}(Q) \tilde{c}_b^2 + 2\alpha \mathbb{E}[z_t^\top P G_\infty z_t | z_{t-\tau}, X_{t-\tau}] \\ &\quad + 2\alpha \zeta (1 + \lambda_{\max}(Q)) \mathbb{E}[\|z_t\|_2^2 | z_{t-\tau}, X_{t-\tau}] \\ &\stackrel{(54)}{\leq} 2\alpha^2 \zeta^2 \lambda_{\max}(Q) \mathbb{E}[\|z_t\|_2^2 | z_{t-\tau}, X_{t-\tau}] + 2\alpha^2 \zeta^2 \lambda_{\max}(Q) \tilde{c}_b^2 - 2\alpha \lambda_{\min}(\Sigma) \mathbb{E}[\|z_t\|_2^2 | z_{t-\tau}, X_{t-\tau}] \\ &\quad + 2\alpha \zeta (1 + \lambda_{\max}(Q)) \mathbb{E}[\|z_t\|_2^2 | z_{t-\tau}, X_{t-\tau}] \\ &\leq 2\alpha^2 \zeta^2 \lambda_{\max}(Q) \mathbb{E}[\|z_t\|_2^2 | z_{t-\tau}, X_{t-\tau}] + 2\alpha^2 \zeta^2 \lambda_{\max}(Q) \tilde{c}_b^2 - 2\alpha \left(\varkappa_1 - \frac{\alpha}{\beta} \varkappa_2 \right) \mathbb{E}[\|z_t\|_2^2 | z_{t-\tau}, X_{t-\tau}] \\ &\quad + 2\alpha \zeta (1 + \lambda_{\max}(Q)) \mathbb{E}[\|z_t\|_2^2 | z_{t-\tau}, X_{t-\tau}] \\ &= \left(-2\alpha \left(\varkappa_1 - \frac{\alpha}{\beta} \varkappa_2 \right) + 2\alpha \zeta (1 + \lambda_{\max}(Q)) + 2\alpha^2 \zeta^2 \lambda_{\max}(Q) \right) \mathbb{E}[\|z_t\|_2^2 | z_{t-\tau}, X_{t-\tau}] + 2\alpha^2 \zeta^2 \lambda_{\max}(Q) \tilde{c}_b^2 \\ &\leq \left(-\alpha \left(\varkappa_1 - \frac{\alpha}{\beta} \varkappa_2 \right) \right) \mathbb{E}[\|z_t\|_2^2 | z_{t-\tau}, X_{t-\tau}] + 2\alpha^2 \zeta^2 \lambda_{\max}(Q) \tilde{c}_b^2 \end{aligned} \quad (55)$$

where the last (55) holds since we need an additional condition as follows,

$$\left(-\alpha \left(\varkappa_1 - \frac{\alpha}{\beta} \varkappa_2 \right) + \alpha \zeta (1 + \lambda_{\max}(Q)) + \alpha^2 \zeta^2 \lambda_{\max}(Q) \right) \leq 0.$$

□

D.3: Proof of Theorem 3

Proof. Lemma 1 implies

$$\begin{aligned}\mathbb{E}[L(z_{t+1})] - \mathbb{E}[L(z_t)] &\leq -\alpha\left(\frac{1}{2}\varkappa_1 - \frac{\alpha}{\beta}\varkappa_2\right)\mathbb{E}[\|z_t\|_2^2] + 2\alpha^2\zeta^2\lambda_{\max}(Q)\tilde{c}_b^2 \\ &\leq -\alpha\left(\frac{1}{2}\varkappa_1 - \frac{\alpha}{\beta}\varkappa_2\right)\frac{1}{\lambda_{\max}(Q)}\mathbb{E}[L(z_t)] + 2\alpha^2\zeta^2\lambda_{\max}(Q)\tilde{c}_b^2.\end{aligned}$$

Rewrite above equation, we have

$$\mathbb{E}[L(z_{t+1})] \leq \left(1 - \alpha\left(\frac{1}{2}\varkappa_1 - \frac{\alpha}{\beta}\varkappa_2\right)\frac{1}{\lambda_{\max}(Q)}\right)\mathbb{E}[L(z_t)] + 2\alpha^2\zeta^2\lambda_{\max}(Q)\tilde{c}_b^2,$$

i.e., we have

$$\mathbb{E}[L(z_{t+1})] - \frac{2\alpha^2\zeta^2\lambda_{\max}(Q)\tilde{c}_b^2}{\alpha\left(\frac{1}{2}\varkappa_1 - \frac{\alpha}{\beta}\varkappa_2\right)\frac{1}{\lambda_{\max}(Q)}} \leq \left(1 - \alpha\left(\frac{1}{2}\varkappa_1 - \frac{\alpha}{\beta}\varkappa_2\right)\frac{1}{\lambda_{\max}(Q)}\right)\left(\mathbb{E}[L(z_t)] - \frac{2\alpha^2\zeta^2\lambda_{\max}(Q)\tilde{c}_b^2}{\alpha\left(\frac{1}{2}\varkappa_1 - \frac{\alpha}{\beta}\varkappa_2\right)\frac{1}{\lambda_{\max}(Q)}}\right).$$

Rewrite above equation, we have

$$\mathbb{E}[L(z_{t+1})] - \frac{2\alpha^2\zeta^2\lambda_{\max}^2(P)\tilde{c}_b^2}{\alpha\left(\frac{1}{2}\varkappa_1 - \frac{\alpha}{\beta}\varkappa_2\right)} \leq \left(1 - \alpha\left(\frac{1}{2}\varkappa_1 - \frac{\alpha}{\beta}\varkappa_2\right)\frac{1}{\lambda_{\max}(Q)}\right)\left(\mathbb{E}[L(z_t)] - \frac{2\alpha^2\zeta^2\lambda_{\max}^2(P)\tilde{c}_b^2}{\alpha\left(\frac{1}{2}\varkappa_1 - \frac{\alpha}{\beta}\varkappa_2\right)}\right).$$

To simplify notations, we introduce

$$u = 1 - \alpha\left(\frac{1}{2}\varkappa_1 - \frac{\alpha}{\beta}\varkappa_2\right)\frac{1}{\lambda_{\max}(Q)}, v = 2\alpha^2\zeta^2\lambda_{\max}(Q)\tilde{c}_b^2.$$

Furthermore, applying above equation recursively, we have

$$\mathbb{E}[L(z_t)] \leq u^{t-\tau}\mathbb{E}[L(z_\tau)] + v\frac{1-u^{t-\tau}}{1-u} \leq u^{t-\tau}\mathbb{E}[L(z_\tau)] + \frac{v}{1-u}.$$

Then, the following equation holds

$$\mathbb{E}[\|z_t\|_2^2] \leq \frac{1}{\lambda_{\min}(P)}\mathbb{E}[L(z_t)] \leq \frac{1}{\lambda_{\min}(P)}\left(u^{t-\tau}\mathbb{E}[L(z_\tau)] + \frac{v}{1-u}\right). \quad (56)$$

Additionally,

$$\begin{aligned}\mathbb{E}[L(z_\tau)] &\leq \lambda_{\max}(Q)\mathbb{E}[\|z_\tau\|_2^2] \leq \lambda_{\max}(Q)(\mathbb{E}[\|z_\tau - z_0\|_2^2] + \|z_0\|_2^2) \\ &\leq \lambda_{\max}(Q)((2\alpha\zeta\tau\|z_0\|_2 + 2\alpha\zeta\tau\tilde{c}_b)^2 + \|z_0\|_2^2) \\ &= \lambda_{\max}(Q)(4\alpha^2\zeta^2\tau^2(\|z_0\|_2 + \tilde{c}_b)^2 + \|z_0\|_2^2).\end{aligned}$$

Taking it to (56), we have

$$\begin{aligned}\mathbb{E}[\|z_t\|_2^2] &\leq \frac{\lambda_{\max}(Q)}{\lambda_{\min}(P)}u^{t-\tau}(4\alpha^2\zeta^2\tau^2(\|z_0\|_2 + \tilde{c}_b)^2 + \|z_0\|_2^2) + \frac{1}{\lambda_{\min}(P)}\frac{v}{1-u} \\ &= \alpha^2\zeta^2u^{t-\tau}\underbrace{\kappa(P)(4\zeta^2\tau^2(\|z_0\|_2 + \tilde{c}_b)^2 + \|z_0\|_2^2)}_{=: \eta_1} + \alpha\kappa(P)\underbrace{\frac{2\lambda_{\max}(Q)\tilde{c}_b^2}{\left(\frac{1}{2}\varkappa_1 - \frac{\alpha}{\beta}\varkappa_2\right)}}_{=: \eta_2} \\ &= \alpha^2\eta_1\left(1 - \frac{\alpha}{\lambda_{\max}(Q)}\left(\frac{1}{2}\varkappa_1 - \frac{\alpha}{\beta}\varkappa_2\right)\right)^{t-\tau} + \alpha\eta_2.\end{aligned} \quad (57)$$

□

Appendix E: Performance over Primal-Dual Gap Error

According to Nemirovski et al., (2009), we can measure the convergence of problem (19) by *primal-dual gap error*.

Definition 2 (Primal-Dual Gap Error). *Recall Ψ defined in (18), the primal-dual gap error $\epsilon_\Psi(\theta, \omega)$ at each solution (ω, θ) is defined as:*

$$\epsilon_\Psi(\theta, \omega) = \max_{\omega'} \Psi(\theta, \omega') - \min_{\theta'} \Psi(\theta', \omega).$$

Theorem 4 (Convergence of Algorithm 1). *Under Assumption 1-2. Consider the sequence $\{(\theta_t, \omega_t)\}_{t=1}^T$ generated Algorithm*

1. *Let C be a constant defined in (61), step-size $\alpha_t = \beta_t = \frac{2}{C\sqrt{5t}}$ and $\tilde{\theta}_T = \frac{\sum_{t=1}^T \alpha_t \theta_t}{\sum_{t=1}^T \alpha_t}$, $\tilde{\omega}_T = \frac{\sum_{t=1}^T \alpha_t \omega_t}{\sum_{t=1}^T \alpha_t}$. Then primal-dual gap error $\epsilon_\Psi(\tilde{\theta}_T, \tilde{\omega}_T)$ is upper-bounded by*

$$\mathbb{E}[\epsilon_\Psi(\tilde{\theta}_T, \tilde{\omega}_T)] \leq C\sqrt{\frac{5}{T}}. \quad (58)$$

Furthermore, for any $\delta \in (0, \frac{2}{e})$, the following holds with probability at least $1 - \delta$,

$$\epsilon_\Psi(\tilde{\theta}_T, \tilde{\omega}_T) \leq C\sqrt{\frac{5}{T}} \left(8 + 2 \log \frac{2}{\delta}\right). \quad (59)$$

Proof. The proof of Theorem 4 relies on some results in the section 3.1 of (Nemirovski et al. 2009), we refer the reader to that reference for further technical details. Let $\hat{G}(\theta, \omega)$ be the stochastic gradient vector of $\Psi(\theta, \omega)$:

$$\hat{G}(\theta, \omega) = \begin{pmatrix} \hat{g}_\theta(\theta, \omega) \\ \hat{g}_\omega(\theta, \omega) \end{pmatrix} = \begin{pmatrix} \hat{A}_t \omega \\ \hat{A}_t \theta + \hat{b}_t - \hat{M}_t \omega \end{pmatrix}.$$

According to Nemirovski et al., (2009), we need to check: I) $\hat{G}(\theta, \omega)$ is an unbiased estimate of the gradient of $\Psi(\theta, \omega)$; II) $\mathbb{E}[\|\hat{G}(\theta, \omega)\|]$ is uniformly bounded on the region $D_\theta \times D_\omega$.

From (22), $\hat{g}_\theta(\theta, \omega)$ and $\hat{g}_\omega(\theta, \omega)$ are the unbiased estimates of $\partial_\theta \Psi(\theta, \omega)$ and $\partial_\omega \Psi(\theta, \omega)$ correspondingly:

$$\mathbb{E}[\hat{G}(\theta, \omega)] = \begin{pmatrix} \mathbb{E}[\hat{g}_\theta(\theta, \omega)] \\ \mathbb{E}[\hat{g}_\omega(\theta, \omega)] \end{pmatrix} = \begin{pmatrix} \partial_\theta \Psi(\theta, \omega) \\ \partial_\omega \Psi(\theta, \omega) \end{pmatrix}. \quad (60)$$

Furthermore, we should check for each (θ_k, ω_k) , the terms $\mathbb{E}[\hat{g}_\theta(\theta_k, \omega_k)]$ and $\mathbb{E}[\hat{g}_\omega(\theta_k, \omega_k)]$ are uniformly bounded:

$$\begin{aligned} \mathbb{E}[\|\hat{g}_\omega(\theta_k, \omega_k)\|_2^2] &= \mathbb{E}[\|\hat{A}_t \theta_t + \hat{b}_t - \hat{M}_t \omega_t\|_2^2] \leq C_b^2 + C_A^2 \text{diam}^2(D_\theta) + C_M^2 \text{diam}^2(D_\omega) \stackrel{\text{def}}{=} \tilde{C}_1^2, \\ \mathbb{E}[\|\hat{g}_\theta(\theta_k, \omega_k)\|_2^2] &= \mathbb{E}[\|\hat{A}_t \omega\|_2^2] \leq C_A^2 \text{diam}^2(D_\omega) \stackrel{\text{def}}{=} \tilde{C}_2^2, \end{aligned}$$

where ‘‘diam’’ is short for diameter. Let C be a constant:

$$C = 4\text{diam}^2(D_\omega)\tilde{C}_1^2 + \text{diam}^2(D_\theta)\tilde{C}_2^2. \quad (61)$$

Then, according to Eq.(3.15) in (Nemirovski et al. 2009), the result (58) holds. Furthermore, by the Proposition 3.2 of (Nemirovski et al. 2009), for any $\eta > 1$:

$$\mathbb{P}\left[\epsilon_\Psi(\tilde{\theta}_T, \tilde{\omega}_T) > C\frac{8+2\eta}{\sqrt{T}}\right] \leq 2e^{-\eta},$$

which implies for any $\delta \in (0, \frac{2}{e})$, the result (59) holds with probability at least $1 - \delta$. □

Discussion 1 (Comparison with Existing Works over Primal-Dual Gap Error). *Liu et al. (2015) firstly derive GTD via convex-concave saddle-point formulation, and their optimal convergence rate reaches $\mathbb{E}[\epsilon_\Psi(\tilde{\theta}_T, \tilde{\omega}_T)] = \mathcal{O}(1/\sqrt{T})$. Later, Wang et al.(2017) extends the work of Liu et al.(2015), they suppose the data is generated from Markov processes rather than I.I.D*

assumption. Wang et al.(2017) prove the convergence rate $\mathbb{E}[\epsilon_\Psi(\tilde{\theta}_T, \tilde{\omega}_T)] = \mathcal{O}\left(\frac{\sum_{t=1}^T \alpha_t^2}{\sum_{t=1}^T \alpha_t}\right)$, the optimal convergence rate also

reaches $\mathcal{O}(1/\sqrt{T})$, where they require step-size satisfies $\sum_{t=1}^\infty \alpha_t = \infty$, $\frac{\sum_{t=1}^T \alpha_t^2}{\sum_{t=1}^T \alpha_t} \leq \infty$. Our work, i.e., Theorem 4 matches the best-known theoretical results Wang et al.(2017), while we point out a simpler and concrete step-size.

Appendix F: Additional Details of Experiments

For the limitation of space, in this section, we present all the details of experiments.

MountainCar Since the state space of mountaincar domain is continuous, we use the open tile coding software <http://incompleteideas.net/rlai.cs.ualberta.ca/RLAI/RLtoolkit/tilecoding.html> to extract feature of states.

In this experiment, we set the number of tilings to be 4 and there are no white noise features. The performance is an average 5 runs and each run contains 5000 episodes. We set $\lambda = 0.99$, $\gamma = 0.99$. The MSPBE/MSE distribution is computed over the combination of step-size, $(\alpha_t, \frac{\beta_t}{\alpha_t}) \in [0.1 \times 2^j | j = -10, -9, \dots, -1, 0]^2$, and $\lambda = 0.99$. Following suggestions from Section 10.1 in (Sutton and Barto 2018), we set all the initial state-action values to be 0, which is optimistic to cause extensive exploration.

Baird Example The Baird example considers the episodic seven-state, two-action MDP. The dashed action takes the system to one of the six upper states with equal probability, whereas the solid action takes the system to the seventh state. The behavior policy b selects the dashed and solid actions with probabilities $\frac{6}{7}$ and $\frac{1}{7}$, so that the next-state distribution under it is uniform (the same for all nonterminal states), which is also the starting distribution for each episode. The target policy π always takes the solid action, and so the on-policy distribution (for π) is concentrated in the seventh state. The reward is zero on all transitions. The discount rate is $\gamma = 0.99$. The feature $\phi(\cdot, \text{dashed})$ and $\phi(\cdot, \text{solid})$ are defined as follows,

$$\begin{aligned}
 \phi(\mathbf{s}_1, \text{dashed}) &= (2, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, , 0, , 0, 0, 0, 0) \\
 \phi(\mathbf{s}_2, \text{dashed}) &= (0, 2, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, , 0, , 0, 0, 0, 0) \\
 &\dots \\
 \phi(\mathbf{s}_7, \text{dashed}) &= (0, 0, 0, 0, 0, 0, 2, 1, 0, 0, 0, 0, , 0, , 0, 0, 0, 0), \tag{62}
 \end{aligned}$$

$$\begin{aligned}
 \phi(\mathbf{s}_1, \text{solid}) &= (0, 0, 0, 0, 0, 0, 0, 0, 2, 0, 0, 0, , 0, , 0, 0, 0, 1) \\
 \phi(\mathbf{s}_2, \text{solid}) &= (0, 0, 0, 0, 0, 0, 0, 0, 2, 0, 0, , 0, , 0, 0, 0, 1) \\
 &\dots \\
 \phi(\mathbf{s}_7, \text{solid}) &= (0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, , 0, , 0, 0, 2, 1). \tag{63}
 \end{aligned}$$