

---

# Is Minimizing Errors the Only Option for Value-based Reinforcement Learning?

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       The existing research on value-based reinforcement learning also minimizes the  
2       error. However, is error minimization really the only option for value-based  
3       reinforcement learning? We can easily observe that the policy on action choosing  
4       probabilities is often related to the relative values, and has nothing to do with  
5       their absolute values. Based on this observation, we propose the objective of  
6       variance minimization instead of error minimization, derive many new variance  
7       minimization algorithms, both including a traditional parameter  $\omega$ , and conduct an  
8       analysis of the convergence rate and experiments. The experimental results show  
9       that our proposed variance minimization algorithms converge much faster.

## 10   1 Introduction

11   Reinforcement learning can be mainly divided into two categories: value-based reinforcement  
12   learning and policy gradient-based reinforcement learning. This paper focuses on temporal difference  
13   learning based on linear approximated valued functions. Its research is usually divided into two steps:  
14   the first step is to establish the convergence of the algorithm, and the second step is to accelerate the  
15   algorithm.

16   In terms of stability, Sutton [1988] established the convergence of on-policy TD(0), and Tsitsiklis  
17   and Van Roy [1997] established the convergence of on-policy TD( $\lambda$ ). However, “The deadly triad”  
18   consisting of off-policy learning, bootstrapping, and function approximation makes the stability a  
19   difficult problem [Sutton and Barto, 2018]. To solve this problem, convergent off-policy temporal  
20   difference learning algorithms are proposed, e.g., BR Baird and others [1995], GTD Sutton *et al.*  
21   [2008], GTD2 and TDC Sutton *et al.* [2009], ETD Sutton *et al.* [2016], and MRetrace Chen *et al.*  
22   [2023].

23   In terms of acceleration, Hackman [2012] proposed Hybrid TD algorithm with on-policy matrix. Liu  
24   *et al.* [2015, 2016, 2018] proposed true stochastic algorithms, i.e., GTD-MP and GTD2-MP, from a  
25   convex-concave saddle-point formulation. Second-order methods are used to accelerate TD learning,  
26   e.g., Quasi Newton TD Givchi and Palhang [2015] and accelerated TD (ATD) [Pan *et al.*, 2017].  
27   Hallak *et al.* [2016] introduced an new parameter to reduce variance for ETD. Zhang and Whiteson  
28   [2022] proposed truncated ETD with a lower variance. Variance Reduced TD with direct variance  
29   reduction technique [Johnson and Zhang, 2013] is proposed by Korda and La [2015] and analysed by  
30   Xu *et al.* [2019]. How to further improve the convergence rates of reinforcement learning algorithms  
31   is currently still an open problem.

32   Algorithm stability is prominently reflected in the changes to the objective function, transitioning  
33   from mean squared errors (MSE) [Sutton and Barto, 2018] to mean squared bellman errors (MSBE)  
34   Baird and others [1995], then to norm of the expected TD update Sutton *et al.* [2009], and further to  
35   mean squared projected Bellman errors (MSPBE) Sutton *et al.* [2009]. On the other hand, algorithm

acceleration is more centered around optimizing the iterative update formula of the algorithm itself without altering the objective function, thereby speeding up the convergence rate of the algorithm. The emergence of new optimization objective functions often leads to the development of novel algorithms. The introduction of new algorithms, in turn, tends to inspire researchers to explore methods for accelerating algorithms, leading to the iterative creation of increasingly superior algorithms.

The kernel loss function can be optimized using standard gradient-based methods, addressing the issue of double sampling in residual gradient algorithm Feng *et al.* [2019]. It ensures convergence in both on-policy and off-policy scenarios. The logistic bellman error is convex and smooth in the action-value function parameters, with bounded gradients Bas-Serrano *et al.* [2021]. In contrast, the squared Bellman error is not convex in the action-value function parameters, and RL algorithms based on recursive optimization using it are known to be unstable.

It is necessary to propose a new objective function, but the mentioned objective functions above are all some form of error. Is minimizing error the only option for value-based reinforcement learning?

For policy evaluation experiments, differences in objective functions may result in inconsistent fixed points. This inconsistency makes it difficult to uniformly compare the superiority of algorithms derived from different objective functions. However, for control experiments, since the choice of actions depends on the relative values of the Q values rather than their absolute values, the presence of solution bias is acceptable.

Based on this observation, we propose alternate objective functions instead of minimizing errors. We minimize Variance of Bellman Error (VBE), Variance of Projected Bellman Error (VPBE), and Variance of the norm of the expected TD update (VNEU) and derive Variance Minimization (VM) algorithms. These algorithms preserve the invariance of the optimal policy in the control environments, but significantly reduce the variance of gradient estimation, and thus hastening convergence.

The contributions of this paper are as follows: (1) Introduction of novel objective functions based on the invariance of the optimal policy. (2) Derived mang variance minimization algorithms, including on-policy and one off-policy. (3) Proof of their convergence. (4) Analysis of the convergence rate of on-policy algorithm. (5) Experiments demonstrating the faster convergence speed of the proposed algorithms.

## 2 Background

Reinforcement learning agent interacts with environment, observes state, takes sequential decision makings to influence environment, and obtains rewards. Consider an infinite-horizon discounted Markov Decision Process (MDP), defined by a tuple  $\langle S, A, R, P, \gamma \rangle$ , where  $S = \{1, 2, \dots, N\}$  is a finite set of states of the environment;  $A$  is a finite set of actions of the agent;  $R : S \times A \times S \rightarrow \mathbb{R}$  is a bounded deterministic reward function;  $P : S \times A \times S \rightarrow [0, 1]$  is the transition probability distribution; and  $\gamma \in (0, 1)$  is the discount factor Sutton and Barto [2018]. Due to the requirements of online learning, value iteration based on sampling is considered in this paper. In each sampling, an experience (or transition)  $\langle s, a, s', r \rangle$  is obtained.

A policy is a mapping  $\pi : S \times A \rightarrow [0, 1]$ . The goal of the agent is to find an optimal policy  $\pi^*$  to maximize the expectation of a discounted cumulative rewards in a long period. State value function  $V^\pi(s)$  for a stationary policy  $\pi$  is defined as:

$$V^\pi(s) = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_k | s_0 = s \right].$$

Linear value function for state  $s \in S$  is defined as:

$$V_\theta(s) := \theta^\top \phi(s) = \sum_{i=1}^m \theta_i \phi_i(s), \quad (1)$$

where  $\theta := (\theta_1, \theta_2, \dots, \theta_m)^\top \in \mathbb{R}^m$  is a parameter vector,  $\phi := (\phi_1, \phi_2, \dots, \phi_m)^\top \in \mathbb{R}^m$  is a feature function defined on state space  $S$ , and  $m$  is the feature size.

Tabular temporal difference (TD) learning Sutton and Barto [2018] has been successfully applied to small-scale problems. To deal with the well-known curse of dimensionality of large scale MDPs,

Table 1: Classification accuracies for naive Bayes and flexible Bayes on various data sets.

ACTION	$Q$ VALUE	$Q$ VALUE WITH BIAS
$Q(s, a_0)$	1	5
$Q(s, a_1)$	2	6
$Q(s, a_2)$	3	7
$Q(s, a_3)$	4	8
$\arg \min_a Q(s, a)$	$a_3$	$a_3$

value function is usually approximated by a linear model, kernel methods, decision trees, or neural networks, etc. This paper focuses on the linear model, where features are usually hand coded by domain experts.

TD learning can also be used to find optimal strategies. The problem of finding an optimal policy is often called the control problem. Two popular TD methods are Sarsa and Q-learning. The former is an on-policy TD control, while the latter is an off-policy control.

It is well known that TDC algorithm Sutton *et al.* [2009] guarantees convergence under off-policy conditions while the off-policy TD algorithm may diverge. The objective function of TDC is MSPBE. TDC is essentially an adjustment or correction of the TD update so that it follows the gradient of the MSPBE objective function. In the context of the TDC algorithm, the control algorithm is known as Greedy-GQ( $\lambda$ ) Sutton *et al.* [2009]. When  $\lambda$  is set to 0, it is denoted as GQ(0).

### 3 Variance Minimization Algorithms

#### 3.1 Motivation

As shown in Table 1, although there is a bias between the true value and the predicted value, action  $a_3$  is still chosen under the greedy-policy. On the contrary, supervised learning is usually used to predict temperature, humidity, morbidity, etc. If the bias is too large, the consequences could be serious.

In addition, reward shaping can significantly speed up the learning by adding a shaping reward  $F(s, s')$  to the original reward  $r$ , where  $F(s, s')$  is the general form of any state-based shaping reward. Static potential-based reward shaping (Static PBRS) maintains the policy invariance if the shaping reward follows from  $F(s, s') = \gamma f(s') - f(s)$  Ng *et al.* [1999].

This means that we can make changes to the TD error  $\delta = r + \gamma \theta^\top \phi' - \theta^\top \phi$  while still ensuring the invariance of the optimal policy,

$$\delta - \omega = r + \gamma \theta^\top \phi' - \theta^\top \phi - \omega,$$

where  $\omega$  is a constant, acting as a static PBRS. This also means that algorithms with the optimization goal of minimizing errors, after introducing reward shaping, may result in larger or smaller bias. Fortunately, as discussed above, bias is acceptable in reinforcement learning. However, the problem is that selecting an appropriate  $\omega$  requires expert knowledge. This forces us to learn  $\omega$  dynamically, i.e.,  $\omega = \omega_t$  and dynamic PBRS can also maintain the policy invariance if the shaping reward is  $F(s, t, s', t') = \gamma f(s', t') - f(s, t)$ , where  $t$  is the time-step the agent reaches in state  $s$  Devlin and Kudenko [2012]. However, this result requires the convergence guarantee of the dynamic potential function  $f(s, t)$ . If  $f(s, t)$  does not converge as the time-step  $t \rightarrow \infty$ , the Q-values of dynamic PBRS are not guaranteed to converge.

Let  $f_{\omega_t}(s) = \frac{\omega_t}{\gamma-1}$ . Thus,  $F_{\omega_t}(s, s') = \gamma f_{\omega_t}(s') - f_{\omega_t}(s) = \omega_t$  is a dynamic PBRS. And if  $\omega$  converges finally, the dynamic potential function  $f(s, t)$  will converge. Bias is the expected difference between the predicted value and the true value. Therefore, under the premise of bootstrapping, we first think of letting  $\omega \doteq \mathbb{E}[\delta|s] = \mathbb{E}[\delta]$ .

As we all know, the optimization process of linear TD(0) (semi-gradient) and linear TDC are as follows, respectively:

$$\theta^* = \arg \min_{\theta} \mathbb{E}[(\mathbb{E}[\delta|s])^2],$$

---

**Algorithm 1** VMTD algorithm with linear function approximation in the on-policy setting

---

**Input:**  $\theta_0, \omega_0, \gamma$ , learning rate  $\alpha_t$  and  $\beta_t$

**repeat**

For any episode, initialize  $\theta_0$  arbitrarily,  $\omega_0$  to 0,  $\gamma \in (0, 1]$ , and  $\alpha_t$  and  $\beta_t$  are constant.

**for**  $t = 0$  **to**  $T - 1$  **do**

Take  $A_t$  from  $S_t$  according to policy  $\mu$ , and arrive at  $S_{t+1}$

Observe sample  $(S_t, R_{t+1}, S_{t+1})$  at time step  $t$  (with their corresponding state feature vectors)

$$\delta_t = R_{t+1} + \gamma \theta_t^\top \phi'_t - \theta_t^\top \phi_t$$

$$\theta_{t+1} \leftarrow \theta_t + \alpha_t (\delta_t - \omega_t) \phi_t$$

$$\omega_{t+1} \leftarrow \omega_t + \beta_t (\delta_t - \omega_t)$$

$$S_t = S_{t+1}$$

**end for**

**until** terminal episode

---

118 and

$$\theta^* = \arg \min_{\theta} \mathbb{E}[\delta \phi]^\top \mathbb{E}[\phi \phi^\top]^{-1} \mathbb{E}[\delta \phi].$$

119 As a result, two novel objective functions and their corresponding algorithms are proposed, where  $\omega$   
120 is subsequently proven to converge, meaning that these two algorithms can maintain the invariance of  
121 the optimal strategy.

### 122 3.2 Variance Minimization TD Learning: VMTD

123 For on-policy learning, a novel objective function, Variance of Bellman Error (VBE), is proposed as  
124 follows:

$$\begin{aligned} \arg \min_{\theta} \text{VBE}(\theta) &= \arg \min_{\theta} \mathbb{E}[(\mathbb{E}[\delta|s] - \mathbb{E}[\mathbb{E}[\delta|s]])^2] \\ &= \arg \min_{\theta, \omega} \mathbb{E}[(\mathbb{E}[\delta|s] - \omega)^2]. \end{aligned} \quad (2)$$

125 Clearly, it is no longer to minimize Bellman errors.

126 First, the parameter  $\omega$  is derived directly based on stochastic gradient descent:

$$\omega_{k+1} \leftarrow \omega_k + \beta_k (\delta_k - \omega_k), \quad (3)$$

127 where  $\delta_k$  is the TD error as follows:

$$\delta_k = r_{k+1} + \gamma \theta_k^\top \phi'_k - \theta_k^\top \phi_k. \quad (4)$$

128 Then, based on stochastic semi-gradient descent, the update of the parameter  $\theta$  is as follows:

$$\theta_{k+1} \leftarrow \theta_k + \alpha_k (\delta_k - \omega_k) \phi_k. \quad (5)$$

129 The pseudocode of the VMTD algorithm is shown in Algorithm 1.

130 For control tasks, two extensions of VMTD are named VMSarsa and VMQ respectively, and the  
131 update formulas are shown below:

$$\theta_{k+1} \leftarrow \theta_k + \alpha_k (\delta_k - \omega_k) \phi(s_k, a_k). \quad (6)$$

132 and

$$\omega_{k+1} \leftarrow \omega_k + \beta_k (\delta_k - \omega_k), \quad (7)$$

133 where  $\delta_k$  delta in VMSarsa is:

$$\delta_k = r_{k+1} + \gamma \theta_k^\top \phi(s_{k+1}, a_{k+1}) - \theta_k^\top \phi(s_k, a_k), \quad (8)$$

134 and  $\delta_k$  delta in VMQ is:

$$\delta_k = r_{k+1} + \gamma \max_{a \in A} \theta_k^\top \phi(s_{k+1}, a) - \theta_k^\top \phi(s_k, a_k). \quad (9)$$

### 135 3.3 Variance Minimization TDC Learning: VMTDC

136 For off-policy learning, we employ a projection operator. The objective function is called Variance of  
137 Projected Bellman error (VPBE), and the corresponding algorithm is called VMTDC.

$$\begin{aligned} \text{VPBE}(\theta) &= \mathbb{E}[(\delta - \mathbb{E}[\delta])\phi]^\top \mathbb{E}[\phi\phi^\top]^{-1} \mathbb{E}[(\delta - \mathbb{E}[\delta])\phi] \\ &= \mathbb{E}[(\delta - \omega)\phi]^\top \mathbb{E}[\phi\phi^\top]^{-1} \mathbb{E}[(\delta - \omega)\phi], \end{aligned} \quad (10)$$

138 where  $\omega$  is used to estimate  $\mathbb{E}[\delta]$ , i.e.,  $\omega \doteq \mathbb{E}[\delta]$ .

139 The derivation process of the VMTDC algorithm is the same as that of the TDC algorithm, the only  
140 difference is that the original  $\delta$  is replaced by  $\delta - \omega$ . Therefore, we can easily get the updated formula  
141 of VMTDC, as follows:

$$\theta_{k+1} \leftarrow \theta_k + \alpha_k [(\delta_k - \omega_k)\phi(s_k) - \gamma\phi(s_{k+1})(\phi^\top(s_k)u_k)], \quad (11)$$

$$u_{k+1} \leftarrow u_k + \zeta_k [\delta_k - \omega_k - \phi^\top(s_k)u_k]\phi(s_k), \quad (12)$$

143 and

$$\omega_{k+1} \leftarrow \omega_k + \beta_k (\delta_k - \omega_k), \quad (13)$$

144 The pseudocode of the VMTDC algorithm for importance-sampling scenario is shown in Algorithm  
145 2 of Appendix A.3.

146 Now, we will introduce the improved version of the GQ(0) algorithm, named VMGQ(0):

$$\begin{aligned} \theta_{k+1} \leftarrow \theta_k &+ \alpha_k [(\delta_k - \omega_k)\phi(s_k, a_k) \\ &- \gamma\phi(s_{k+1}, A_{k+1}^*)(\phi^\top(s_k, a_k)u_k)], \end{aligned} \quad (14)$$

147

$$u_{k+1} \leftarrow u_k + \zeta_k [(\delta_k - u_k) - \phi^\top(s_k, a_k)u_k]\phi(s_k, a_k), \quad (15)$$

148 and

$$\omega_{k+1} \leftarrow \omega_k + \beta_k (\delta_k - \omega_k), \quad (16)$$

149 where  $\delta_k$  is (9) and  $A_{k+1}^* = \arg \max_a (\theta_k^\top \phi(s_{k+1}, a))$ .

### 150 3.4 Variance Minimization ETD Learning: VMETD

151 VMETD by the following update:

$$\rho_k \leftarrow \frac{\pi(A_k|S_k)}{\mu(A_k|S_k)} \quad (17)$$

152

$$F_k \leftarrow \gamma\rho_{k-1}F_{k-1} + 1, \quad (18)$$

153

$$\omega_{k+1} \leftarrow \omega_k + \beta_k (F_k \rho_k \delta_k - \omega_k), \quad (19)$$

154

$$\theta_{k+1} \leftarrow \theta_k + \alpha_k F_k \rho_k (R_{k+1} + \gamma\theta_k^\top \phi_{k+1} - \theta_k^\top \phi_k)\phi_k - \alpha_k \omega_{k+1} \phi_k, \quad (20)$$

155 where  $\mu$  is behavior policy and  $\pi$  is target policy,  $F_t$  is a scalar variable,  $F_0 = 1$ ,  $\omega$  is  
156 used to estimate  $\mathbb{E}[\delta]$ , i.e.,  $\omega \doteq \mathbb{E}[\delta]$ , and  $\mathbf{F}$  is a diagonal matrix with diagonal elements  
157  $f(s) \doteq d_\mu(s) \lim_{t \rightarrow \infty} \mathbb{E}_\mu[F_t | S_k = s]$ , which we assume exists. The vector  $\mathbf{f} \in \mathbb{R}^N$  with components  
158  $[\mathbf{f}]_s \doteq f(s)$  can be written as

$$\begin{aligned} \mathbf{f} &= \mathbf{d}_\mu + \gamma \mathbf{P}_\pi^\top \mathbf{d}_\mu + (\gamma \mathbf{P}_\pi^\top)^2 \mathbf{d}_\mu + \dots \\ &= (\mathbf{I} - \gamma \mathbf{P}_\pi^\top)^{-1} \mathbf{d}_\mu. \end{aligned} \quad (21)$$

## 159 4 Theoretical Analysis

160 The purpose of this section is to establish the stabilities of the VMTD algorithm and the VMTDC  
161 algorithm, and also presents a corollary on the convergence rate of VMTD.

162 **Theorem 4.1.** (Convergence of VMTD). In the case of on-policy learning, consider the iterations (3)  
163 and (5) with (4) of VMTD. Let the step-size sequences  $\alpha_k$  and  $\beta_k$ ,  $k \geq 0$  satisfy in this case  $\alpha_k, \beta_k > 0$ ,  
164 for all  $k$ ,  $\sum_{k=0}^{\infty} \alpha_k = \sum_{k=0}^{\infty} \beta_k = \infty$ ,  $\sum_{k=0}^{\infty} \alpha_k^2 < \infty$ ,  $\sum_{k=0}^{\infty} \beta_k^2 < \infty$ , and  $\alpha_k = o(\beta_k)$ . Assume  
165 that  $(\phi_k, r_k, \phi'_k)$  is an i.i.d. sequence with uniformly bounded second moments, where  $\phi_k$  and  $\phi'_k$   
166 are sampled from the same Markov chain. Let  $A = \text{Cov}(\phi, \phi - \gamma\phi')$ ,  $b = \text{Cov}(r, \phi)$ . Assume that  
167 matrix  $A$  is non-singular. Then the parameter vector  $\theta_k$  converges with probability one to  $A^{-1}b$ .

168 Please refer to the appendix A.1 for detailed proof process.

169 Theorem 3 in Dalal *et al.* [2020] provides a general conclusion on the convergence speed of all linear  
170 two-timescale algorithms. VMTD satisfies the assumptions of this theorem, leading to the following  
171 corollary.

172 **Corollary 4.2.** *Consider the Sparsely Projected variant of VMTD. Then, for  $\alpha_k = 1/(k+1)^\alpha$ ,  
173  $\beta_k = 1/(k+1)^\beta$ ,  $0 < \beta < \alpha < 1$ ,  $p > 1$ , with probability larger than  $1 - \tau$ , for all  $k \geq N_3$ , we have*

$$\|\theta'_k - \theta^*\| \leq C_{3,\theta} \frac{\sqrt{\ln(4d_1^2(k+1)^p/\tau)}}{(k+1)^{\alpha/2}} \quad (22)$$

174

$$\|\omega'_n - \omega^*\| \leq C_{3,\omega} \frac{\sqrt{\ln(4d_2^2(k+1)^p/\tau)}}{(k+1)^{\omega/2}}, \quad (23)$$

175 where  $d_1$  and  $d_2$  represent the dimensions of  $\theta$  and  $\omega$ , respectively. For VMTD,  $d_2 = 1$ . The  
176 meanings of  $N_3, C_{3,\theta}$  and  $C_{3,\omega}$  are explained in Dalal *et al.* [2020]. The formulas for  $\theta'_k$  and  $\omega'_n$  can  
177 be found in (35) and (36).

178 Please refer to the appendix A.2 for detailed proof process.

179 **Theorem 4.3.** *(Convergence of VMTDC). In the case of off-policy learning, consider the iterations  
180 (13), (12) and (11) of VMTDC. Let the step-size sequences  $\alpha_k, \zeta_k$  and  $\beta_k, k \geq 0$  satisfy in this case  
181  $\alpha_k, \zeta_k, \beta_k > 0$ , for all  $k$ ,  $\sum_{k=0}^{\infty} \alpha_k = \sum_{k=0}^{\infty} \beta_k = \sum_{k=0}^{\infty} \zeta_k = \infty$ ,  $\sum_{k=0}^{\infty} \alpha_k^2 < \infty$ ,  $\sum_{k=0}^{\infty} \zeta_k^2 < \infty$ ,  
182  $\sum_{k=0}^{\infty} \beta_k^2 < \infty$ , and  $\alpha_k = o(\zeta_k)$ ,  $\zeta_k = o(\beta_k)$ . Assume that  $(\phi_k, r_k, \phi'_k)$  is an i.i.d. sequence with  
183 uniformly bounded second moments. Let  $A = \text{Cov}(\phi, \phi - \gamma\phi')$ ,  $b = \text{Cov}(r, \phi)$ , and  $C = \mathbb{E}[\phi\phi^\top]$ .  
184 Assume that  $A$  and  $C$  are non-singular matrices. Then the parameter vector  $\theta_k$  converges with  
185 probability one to  $A^{-1}b$ .*

186 Please refer to the appendix A.3 for detailed proof process.

187 **Theorem 4.4.** *(Convergence of VMETD). In the case of off-policy learning, consider the iterations  
188 (19) and (20) with (4) of VMETD. Let the step-size sequences  $\alpha_k$  and  $\beta_k, k \geq 0$  satisfy in this  
189 case  $\alpha_k, \beta_k > 0$ , for all  $k$ ,  $\sum_{k=0}^{\infty} \alpha_k = \sum_{k=0}^{\infty} \beta_k = \infty$ ,  $\sum_{k=0}^{\infty} \alpha_k^2 < \infty$ ,  $\sum_{k=0}^{\infty} \beta_k^2 < \infty$ , and  
190  $\alpha_k = o(\beta_k)$ . Assume that  $(\phi_k, r_k, \phi'_k)$  is an i.i.d. sequence with uniformly bounded second moments,  
191 where  $\phi_k$  and  $\phi'_k$  are sampled from the same Markov chain. Let  $A = \text{Cov}(\phi, \phi - \gamma\phi')$ ,  $b = \text{Cov}(r, \phi)$ .  
192 Assume that matrix  $A$  is non-singular. Then the parameter vector  $\theta_k$  converges with probability one  
193 to  $A^{-1}b$ .*

194 Please refer to the appendix A.4 for detailed proof process.

## 195 5 Experimental Studies

196 This section assesses algorithm performance through experiments, which are divided into policy  
197 evaluation experiments and control experiments.

### 198 5.1 Testing Tasks

199 **Random-walk:** as shown in Figure 1, all episodes start in the center state,  $C$ , and proceed to left  
200 or right by one state on each step, equiprobably. Episodes terminate either on the extreme left or  
201 the extreme right, and get a reward of +1 if terminate on the right, or 0 in the other case. In this  
202 task, the true value for each state is the probability of starting from that state and terminating on  
203 the right Sutton and Barto [2018]. Thus, the true values of states from  $A$  to  $E$  are  $\frac{1}{6}, \frac{2}{6}, \frac{3}{6}, \frac{4}{6}, \frac{5}{6}$ ,  
204 respectively. The discount factor  $\gamma = 1.0$ . There are three standard kinds of features for random-walk  
205 problems: tabular feature, inverted feature and dependent feature Sutton *et al.* [2009]. The feature  
206 matrices corresponding to three random walks are shown in Appendix B. Conduct experiments using  
207 an on-policy approach in the Random-walk environment.

208 **Baird's off-policy counterexample:** This task is well known as a counterexample, in which TD  
209 diverges Baird and others [1995]; Sutton *et al.* [2009]. As shown in Figure 2, reward for each  
210 transition is zero. Thus the true values are zeros for all states and for any given policy. The behaviour  
211 policy chooses actions represented by solid lines with a probability of  $\frac{1}{7}$  and actions represented by

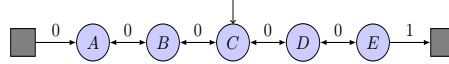


Figure 1: Random walk.

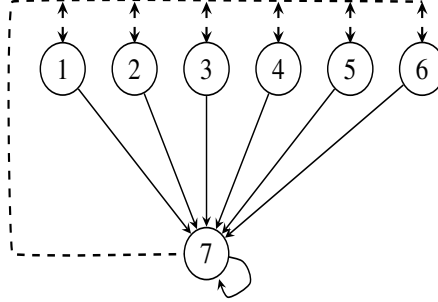
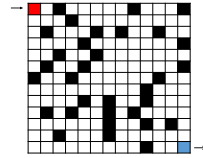


Figure 2: 7-state version of Baird's off-policy counterexample.

212 dotted lines with a probability of  $\frac{6}{7}$ . The target policy is expected to choose the solid line with more  
 213 probability than  $\frac{1}{7}$ , and it chooses the solid line with probability of 1 in this paper. The discount  
 214 factor  $\gamma = 0.99$ , and the feature matrix is defined in Appendix B Baird and others [1995]; Sutton *et*  
 215 *al.* [2009]; Maei [2011].

216 **Maze:** The learning agent should find a shortest path from the upper left corner to the lower  
 217 right corner. In each state, there are four alternative actions: *up*, *down*, *left*, and *right*,  
 218 which takes the agent deterministically to the corresponding neighbour state, except when  
 219 a movement is blocked by an obstacle or the edge of the maze. Rewards are  
 220  $-1$  in all transitions until the agent reaches the goal state. The discount factor  
 221  $\gamma = 0.99$ , and states  $s$  are represented by tabular features. The maximum  
 222 number of moves in the game is set to 1000.

223 **The other three control environments:** Cliff Walking, Mountain Car, and  
 224 Acrobot are selected from the gym official website and correspond to the  
 225 following versions: “CliffWalking-v0”, “MountainCar-v0” and “Acrobot-v1”.  
 226 For specific details, please refer to the gym official website. The maximum  
 227 number of steps for the Mountain Car environment is set to 1000, while the default settings are used  
 228 for the other two environments. In Mountain car and Acrobot, features are generated by tile coding.  
 229 Please, refer to the Appendix B for the selection of learning rates for all experiments.



## 230 5.2 Experimental Results and Analysis

231 For policy evaluation experiments, compare the performance of the VMTD, VMTDC, TD, and TDC  
 232 algorithms. The vertical axis is unified as RVBE.

233 For policy evaluation experiments, the criteria for evaluating algorithms vary. The objective function  
 234 minimized by our proposed new algorithm differs from that of other algorithms. Therefore, to ensure  
 235 fairness in comparisons, this study only contrasts algorithm experiments in controlled settings.

236 This study will compare the performance of Sarsa, Q-learning, GQ(0), AC, VMSarsa, VMQ, and  
 237 VMGQ(0) in four control environments.

238 The learning curves of the algorithms corresponding to policy evaluation experiments and control  
 239 experiments are shown in Figures 3 and 4, respectively. The shaded area in Figure 3, 4 represents the  
 240 standard deviation (std).

241 In the random-walk tasks, VMTD and VMTDC exhibit excellent performance, outperforming TD  
 242 and TDC in the case of dependent random-walk.

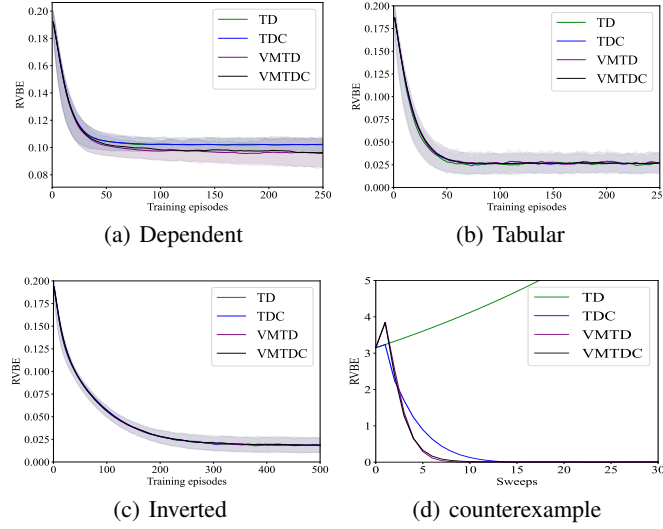


Figure 3: Learning curves of four evaluation environments.

In the 7-state example counter task, TD diverges, while VMTDC converges and performs better than TDC. From the update formula, it can be observed that the VMTD algorithm, like TDC, is also an adjustment or correction of the TD update. What is more surprising is that VMTD also maintains convergence and demonstrates the best performance.

In Maze, Mountain Car, and Acrobot, the convergence speed of VMSarsa, VMQ, and VMGQ(0) has been significantly improved compared to Sarsa, Q-learning, and GQ(0), respectively. The performance of the AC algorithm is at an intermediate level. The performances of VMSarsa, VMQ, and VMGQ(0) in these three experimental environments have no significant differences.

In Cliff Walking, Sarsa and VMSarsa converge to slightly worse solutions compared to other algorithms. The convergence speed of VMSarsa is significantly better than that of Sarsa. The convergence speed of VMGQ(0) and VMQ is better than other algorithms, and the performance of VMGQ(0) is slightly better than that of VMQ.

In summary, the performance of VMSarsa, VMQ, and VMGQ(0) is better than that of other algorithms. In the Cliff Walking environment, the performance of VMGQ(0) is slightly better than that of VMSarsa and VMQ. In the other three experimental environments, the performances of VMSarsa, VMQ, and VMGQ(0) are close.

## 6 Related Work

### 6.1 Difference between VMQ and R-learning

Table 2: Difference between R-learning and tabular VMQ.

algorithms	update formula
R-learning	$Q_{k+1}(s, a) \leftarrow Q_k(s, a) + \alpha_k(r_{k+1} - m_k + \max_{b \in A} Q_k(s, b) - Q_k(s, a))$ $m_{k+1} \leftarrow m_k + \beta_k(r_{k+1} + \max_{b \in A} Q_k(s, b) - Q_k(s, a) - m_k)$
tabular VMQ	$Q_{k+1}(s, a) \leftarrow Q_k(s, a) + \alpha_k(r_{k+1} + \gamma \max_{b \in A} Q_k(s, b) - Q_k(s, a) - \omega_k)$ $\omega_{k+1} \leftarrow \omega_k + \beta_k(r_{k+1} + \gamma \max_{b \in A} Q_k(s, b) - Q_k(s, a) - \omega_k)$

Tabular VMQ’s update formula bears some resemblance to R-learning’s update formula. As shown in Table 2, the update formulas of the two algorithms have the following differences:

(1) The goal of the R-learning algorithm Schwartz [1993] is to maximize the average reward, rather than the cumulative reward, by learning an estimate of the average reward. This estimate  $m$  is then used to update the  $Q$ -values. On the contrary, the  $\omega$  in the tabular VMQ update formula eventually

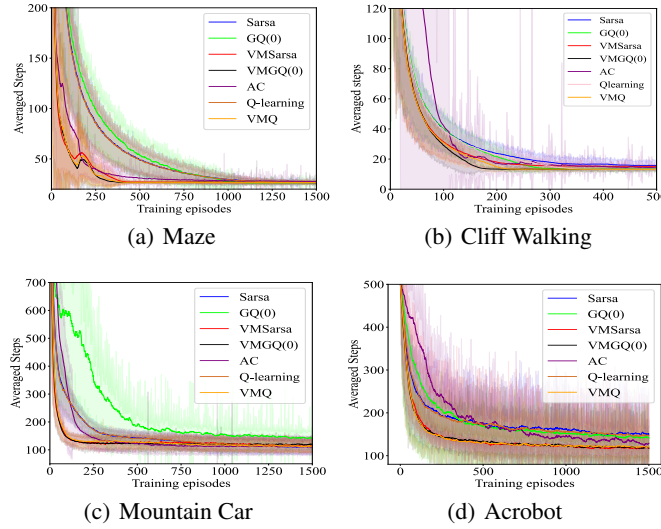


Figure 4: Learning curves of four contral environments.

converges to  $\mathbb{E}[\delta]$ .

(2) When  $\gamma = 1$  in the tabular VMQ update formula, the R-learning update formula is formally the same as the tabular VMQ update formula. Therefore, R-learning algorithm can be considered as a special case of VMQ algorithm in form.

## 6.2 Variance Reduction for TD Learning

The TD with centering algorithm (CTD) Korda and La [2015] was proposed, which directly applies variance reduction techniques to the TD algorithm. The CTD algorithm updates its parameters using the average gradient of a batch of Markovian samples and a projection operator. Unfortunately, the authors' analysis of the CTD algorithm contains technical errors. The VRTD algorithm Xu *et al.* [2020] is also a variance-reduced algorithm that updates its parameters using the average gradient of a batch of i.i.d. samples. The authors of VRTD provide a technically sound analysis to demonstrate the advantages of variance reduction.

## 6.3 Variance Reduction for Policy Gradient Algorithms

Policy gradient algorithms are a class of reinforcement learning algorithms that directly optimize cumulative rewards. REINFORCE is a Monte Carlo algorithm that estimates gradients through sampling, but may have a high variance. Baselines are introduced to reduce variance and to accelerate learning Sutton and Barto [2018]. In Actor-Critic, value function as a baseline and bootstrapping are used to reduce variance, also accelerating convergence Sutton and Barto [2018]. TRPO Schulman *et al.* [2015] and PPO Schulman *et al.* [2017] use generalized advantage estimation, which combines multi-step bootstrapping and Monte Carlo estimation to reduce variance, making gradient estimation more stable and accelerating convergence.

In Variance Minimization, the incorporation of  $\omega \doteq \mathbb{E}[\delta]$  bears a striking resemblance to the use of a baseline in policy gradient methods. The introduction of a baseline in policy gradient techniques does not alter the expected value of the update; rather, it significantly impacts the variance of gradient estimation. The addition of  $\omega \doteq \mathbb{E}[\delta]$  in Variance Minimization preserves the invariance of the optimal policy while stabilizing gradient estimation, reducing the variance of gradient estimation, and hastening convergence.

## 7 Conclusion and Future Work

Value-based reinforcement learning typically aims to minimize error as an optimization objective. As an alternation, this study proposes new objective functions: VBE, VPBE and VNEU, and derives many variance minimization algorithms, including VMTD, VMTDC, VMGTD, VMGTD2 and VMETD. All algorithms demonstrated superior performance in policy evaluation and control experiments. Future work may include, but are not limited to, (1) analysis of the convergence rate of VMTDC. (2) extensions of VBE and VPBE to multi-step returns. (3) extensions to nonlinear approximations, such as neural networks.

### A Relevant proofs

#### A.1 Proof of Theorem 4.1

*Proof.* The proof is based on Borkar's Theorem for general stochastic approximation recursions with two time scales Borkar [1997].

A new one-step linear TD solution is defined as:

$$0 = \mathbb{E}[(\delta - \mathbb{E}[\delta])\phi] = -A\theta + b.$$

Thus, the VMTD's solution is  $\theta_{\text{VMTD}} = A^{-1}b$ .

First, note that recursion (5) can be rewritten as

$$\theta_{k+1} \leftarrow \theta_k + \beta_k \xi(k),$$

where

$$\xi(k) = \frac{\alpha_k}{\beta_k}(\delta_k - \omega_k)\phi_k$$

Due to the settings of step-size schedule  $\alpha_k = o(\beta_k)$ ,  $\xi(k) \rightarrow 0$  almost surely as  $k \rightarrow \infty$ . That is the increments in iteration (3) are uniformly larger than those in (5), thus (3) is the faster recursion. Along the faster time scale, iterations of (3) and (5) are associated to ODEs system as follows:

$$\dot{\theta}(t) = 0, \tag{24}$$

$$\dot{\omega}(t) = \mathbb{E}[\delta_t|\theta(t)] - \omega(t). \tag{25}$$

Based on the ODE (24),  $\theta(t) \equiv \theta$  when viewed from the faster timescale. By the Hirsch lemma Hirsch [1989], it follows that  $\|\theta_k - \theta\| \rightarrow 0$  a.s. as  $k \rightarrow \infty$  for some  $\theta$  that depends on the initial condition  $\theta_0$  of recursion (5). Thus, the ODE pair (24)-(25) can be written as

$$\dot{\omega}(t) = \mathbb{E}[\delta_t|\theta] - \omega(t). \tag{26}$$

Consider the function  $h(\omega) = \mathbb{E}[\delta|\theta] - \omega$ , i.e., the driving vector field of the ODE (26). It is easy to find that the function  $h$  is Lipschitz with coefficient  $-1$ . Let  $h_\infty(\cdot)$  be the function defined by  $h_\infty(\omega) = \lim_{x \rightarrow \infty} \frac{h(x\omega)}{x}$ . Then  $h_\infty(\omega) = -\omega$ , is well-defined. For (26),  $\omega^* = \mathbb{E}[\delta|\theta]$  is the unique globally asymptotically stable equilibrium. For the ODE

$$\dot{\omega}(t) = h_\infty(\omega(t)) = -\omega(t), \tag{27}$$

apply  $\vec{V}(\omega) = (-\omega)^\top(-\omega)/2$  as its associated strict Liapunov function. Then, the origin of (27) is a globally asymptotically stable equilibrium.

Consider now the recursion (3). Let  $M_{k+1} = (\delta_k - \omega_k) - \mathbb{E}[(\delta_k - \omega_k)|\mathcal{F}(k)]$ , where  $\mathcal{F}(k) = \sigma(\omega_l, \theta_l, l \leq k; \phi_s, \phi'_s, r_s, s < k)$ ,  $k \geq 1$  are the sigma fields generated by  $\omega_0, \theta_0, \omega_{l+1}, \theta_{l+1}, \phi_l, \phi'_l$ ,  $0 \leq l < k$ . It is easy to verify that  $M_{k+1}, k \geq 0$  are integrable random variables that satisfy  $\mathbb{E}[M_{k+1}|\mathcal{F}(k)] = 0, \forall k \geq 0$ . Because  $\phi_k, r_k$ , and  $\phi'_k$  have uniformly bounded second moments, it can be seen that for some constant  $c_1 > 0, \forall k \geq 0$ ,

$$\mathbb{E}[\|M_{k+1}\|^2|\mathcal{F}(k)] \leq c_1(1 + \|\omega_k\|^2 + \|\theta_k\|^2).$$

Now Assumptions (A1) and (A2) of Borkar and Meyn [2000] are verified. Furthermore, Assumptions (TS) of Borkar and Meyn [2000] is satisfied by our conditions on the step-size sequences  $\alpha_k, \beta_k$ .

Thus, by Theorem 2.2 of Borkar and Meyn [2000] we obtain that  $\|\omega_k - \omega^*\| \rightarrow 0$  almost surely as  $k \rightarrow \infty$ .

Consider now the slower time scale recursion (5). Based on the above analysis, (5) can be rewritten as

$$\theta_{k+1} \leftarrow \theta_k + \alpha_k(\delta_k - \mathbb{E}[\delta_k|\theta_k])\phi_k.$$

Let  $\mathcal{G}(k) = \sigma(\theta_l, l \leq k; \phi_s, \phi'_s, r_s, s < k)$ ,  $k \geq 1$  be the sigma fields generated by  $\theta_0, \theta_{l+1}, \phi_l, \phi'_l$ ,  $0 \leq l < k$ . Let  $Z_{k+1} = Y_t - \mathbb{E}[Y_t|\mathcal{G}(k)]$ , where

$$Y_k = (\delta_k - \mathbb{E}[\delta_k|\theta_k])\phi_k.$$

Consequently,

$$\begin{aligned} \mathbb{E}[Y_t|\mathcal{G}(k)] &= \mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\theta_k])\phi_k|\mathcal{G}(k)] \\ &= \mathbb{E}[\delta_k\phi_k|\theta_k] - \mathbb{E}[\mathbb{E}[\delta_k|\theta_k]\phi_k] \\ &= \mathbb{E}[\delta_k\phi_k|\theta_k] - \mathbb{E}[\delta_k|\theta_k]\mathbb{E}[\phi_k] \\ &= \text{Cov}(\delta_k|\theta_k, \phi_k), \end{aligned}$$

where  $\text{Cov}(\cdot, \cdot)$  is a covariance operator.

Thus,

$$Z_{k+1} = (\delta_k - \mathbb{E}[\delta_k|\theta_k])\phi_k - \text{Cov}(\delta_k|\theta_k, \phi_k).$$

It is easy to verify that  $Z_{k+1}, k \geq 0$  are integrable random variables that satisfy  $\mathbb{E}[Z_{k+1}|\mathcal{G}(k)] = 0$ ,  $\forall k \geq 0$ . Also, because  $\phi_k, r_k$ , and  $\phi'_k$  have uniformly bounded second moments, it can be seen that for some constant  $c_2 > 0$ ,  $\forall k \geq 0$ ,

$$\mathbb{E}[\|Z_{k+1}\|^2|\mathcal{G}(k)] \leq c_2(1 + \|\theta_k\|^2).$$

Consider now the following ODE associated with (5):

$$\begin{aligned} \dot{\theta}(t) &= \text{Cov}(\delta|\theta(t), \phi) \\ &= \text{Cov}(r + (\gamma\phi' - \phi)^\top \theta(t), \phi) \\ &= \text{Cov}(r, \phi) - \text{Cov}(\theta(t)^\top (\phi - \gamma\phi'), \phi) \\ &= \text{Cov}(r, \phi) - \theta(t)^\top \text{Cov}(\phi - \gamma\phi', \phi) \\ &= \text{Cov}(r, \phi) - \text{Cov}(\phi - \gamma\phi', \phi)^\top \theta(t) \\ &= \text{Cov}(r, \phi) - \text{Cov}(\phi, \phi - \gamma\phi')\theta(t) \\ &= -A\theta(t) + b. \end{aligned} \tag{28}$$

Let  $\vec{h}(\theta(t))$  be the driving vector field of the ODE (28).

$$\vec{h}(\theta(t)) = -A\theta(t) + b.$$

Consider the cross-covariance matrix,

$$\begin{aligned} A &= \text{Cov}(\phi, \phi - \gamma\phi') \\ &= \frac{\text{Cov}(\phi, \phi) + \text{Cov}(\phi - \gamma\phi', \phi - \gamma\phi') - \text{Cov}(\gamma\phi', \gamma\phi')}{2} \\ &= \frac{\text{Cov}(\phi, \phi) + \text{Cov}(\phi - \gamma\phi', \phi - \gamma\phi') - \gamma^2 \text{Cov}(\phi', \phi')}{2} \\ &= \frac{(1 - \gamma^2) \text{Cov}(\phi, \phi) + \text{Cov}(\phi - \gamma\phi', \phi - \gamma\phi')}{2}, \end{aligned} \tag{29}$$

where we eventually used  $\text{Cov}(\phi', \phi') = \text{Cov}(\phi, \phi)$ <sup>1</sup>. Note that the covariance matrix  $\text{Cov}(\phi, \phi)$  and  $\text{Cov}(\phi - \gamma\phi', \phi - \gamma\phi')$  are semi-positive definite. Then, the matrix  $A$  is semi-positive definite because  $A$  is linearly combined by two positive-weighted semi-positive definite matrices (29). Furthermore,  $A$  is nonsingular due to the assumption. Hence, the cross-covariance matrix  $A$  is positive definite.

Therefore,  $\theta^* = A^{-1}b$  can be seen to be the unique globally asymptotically stable equilibrium for ODE (28). Let  $\vec{h}_\infty(\theta) = \lim_{r \rightarrow \infty} \frac{\vec{h}(r\theta)}{r}$ . Then  $\vec{h}_\infty(\theta) = -A\theta$  is well-defined. Consider now the ODE

$$\dot{\theta}(t) = -A\theta(t). \tag{30}$$

The ODE (30) has the origin as its unique globally asymptotically stable equilibrium. Thus, the assumption (A1) and (A2) are verified.  $\square$

<sup>1</sup>The covariance matrix  $\text{Cov}(\phi', \phi')$  is equal to the covariance matrix  $\text{Cov}(\phi, \phi)$  if the initial state is reachable or initialized randomly in a Markov chain for on-policy update.

## A.2 Proof of Corollary 4.2

The update formulas in linear two-timescale algorithms are as follows:

$$\theta_{k+1} = \theta_k + \alpha_k [h_1(\theta_k, \omega_k) + M_{k+1}^{(1)}], \quad (31)$$

$$\omega_{k+1} = \omega_k + \alpha_k [h_2(\theta_k, \omega_k) + M_{k+1}^{(2)}]. \quad (32)$$

where  $\alpha_k, \beta_k \in \mathbb{R}$  are stepsizes and  $M^{(1)} \in \mathbb{R}^{d_1}, M^{(2)} \in \mathbb{R}^{d_2}$  denote noise.  $h_1 : \mathbb{R}^{d_1} \times \mathbb{R}^{d_2} \rightarrow \mathbb{R}^{d_1}$  and  $h_2 : \mathbb{R}^{d_1} \times \mathbb{R}^{d_2} \rightarrow \mathbb{R}^{d_2}$  have the form, respectively,

$$h_1(\theta, \omega) = v_1 - \Gamma_1 \theta - W_1 \omega, \quad (33)$$

$$h_2(\theta, \omega) = v_2 - \Gamma_2 \theta - W_2 \omega, \quad (34)$$

where  $v_1 \in \mathbb{R}^{d_1}, v_2 \in \mathbb{R}^{d_2}, \Gamma_1 \in \mathbb{R}^{d_1 \times d_1}, \Gamma_2 \in \mathbb{R}^{d_2 \times d_1}, W_1 \in \mathbb{R}^{d_1 \times d_2}$  and  $W_2 \in \mathbb{R}^{d_2 \times d_2}$ .  $d_1$  and  $d_2$  are the dimensions of vectors  $\theta$  and  $\omega$ , respectively.

For Theorem 3 in Dalal *et al.* [2020], the theorem still holds even when  $d^{-1}$  is not equal to  $d_2$ . For the VMTD algorithm,  $d_2$  is equal to 1. Dalal *et al.* [2020] presents the matrix assumption, step size assumption, and defines sparse projection.

**Assumption A.1.** (Matrix Assumption).  $W_2$  and  $X_1 = \Gamma_1 - W_1 W_2^{-1} \Gamma_2$  are positive definite(not necessarily symmetric).

**Assumption A.2.** (Step Size Assumption).  $\alpha_k = (k+1)^{-\alpha}$  and  $\beta_k = (k+1)^{-\beta}$ , where  $1 > \alpha > \beta > 0$ .

**Definition A.3.** (Sparse Projection). For  $R > 0$ , let  $\Pi_R(x) = \min\{1, R/\|x\|\}$ .  $x$  be the projection into the ball with radius  $R$  around the origin. The sparse projection operator

$$\Pi_{n,R} = \begin{cases} \Pi_R, & \text{if } k = n^n - 1 \text{ for some } n \in \mathbb{Z}_{>0}, \\ I, & \text{otherwise.} \end{cases}$$

We call it sparse as it projects only on specific indices that are exponentially far apart.

Pick an arbitrary  $p > 1$ . Fix some constant  $R_{\text{proj}}^\theta > 0$  and  $R_{\text{proj}}^\omega > 0$  for the radius of the projection ball. Further, let

$$\theta^* = X_1^{-1} b_1, \omega^* = W_2^{-1} (v_2 - \Gamma_2 \theta^*)$$

with  $b_1 = v_1 - W_1 W_2^{-1} v_2$ . The formula for the sparse projection update in linear two-timescale algorithms is as follows:

$$\theta'_{k+1} = \Pi_{k+1, R_{\text{proj}}^\theta} (\theta'_k + \alpha_k [h_1(\theta'_k, \omega'_k) + M_{k+1}^{(1')}]), \quad (35)$$

$$\omega'_{k+1} = \Pi_{k+1, R_{\text{proj}}^\omega} (\omega'_k + \beta_k [h_2(\theta'_k, \omega'_k) + M_{k+1}^{(2')}]). \quad (36)$$

*Proof.* As long as the VMTD algorithm satisfies Assumption A.1, the convergence speed of the VMTD algorithm can be obtained.

VMTD's update rule is

$$\theta_{k+1} = \theta_k + \alpha_k (\delta_k - \omega_k) \phi_k.$$

$$\omega_{k+1} = \omega_k + \beta_k (\delta_k - \omega_k).$$

Thus,  $h_1(\theta, \omega) = \text{Cov}(r, \phi) - \text{Cov}(\phi, \phi - \gamma \phi') \theta$ ,  $h_2(\theta, \omega) = \mathbb{E}[r] + \mathbb{E}[\gamma \phi'^\top - \phi^\top] \theta - \omega$ ,  $\Gamma_1 = \text{Cov}(\phi, \phi - \gamma \phi')$ ,  $W_1 = 0$  and  $\Gamma_2 = -\mathbb{E}[\gamma \phi'^\top - \phi^\top]$ ,  $W_2 = 1$ ,  $v_2 = \mathbb{E}[r]$ . Additionally,  $X_1 = \Gamma_1 - W_1 W_2^{-1} \Gamma_2 = \text{Cov}(\phi, \phi - \gamma \phi')$ . It can be deduced from the proof A.1 that  $X_1$  is a positive definite matrix. The VMTD algorithm satisfies the Assumption A.1. By the proof A.1, Definition 1 in Dalal *et al.* [2020] is satisfied. We can apply the Theorem 3 in Dalal *et al.* [2020] to get the Corollary 4.2.

□

### 386 A.3 Proof of Theorem 4.3

387 *Proof.* The proof is similar to that given by Sutton *et al.* [2009] for TDC, but it is based on multi-  
 388 time-scale stochastic approximation.

389 For the VMTDC algorithm, a new one-step linear TD solution is defined as:

$$0 = \mathbb{E}[(\phi - \gamma\phi' - \mathbb{E}[\phi - \gamma\phi'])\phi^\top] \mathbb{E}[\phi\phi^\top]^{-1} \mathbb{E}[(\delta - \mathbb{E}[\delta])\phi] = A^\top C^{-1}(-A\theta + b).$$

390 The matrix  $A^\top C^{-1}A$  is positive definite. Thus, the VMTD's solution is  $\theta_{\text{VMTDC}} = \theta_{\text{VMTD}} = A^{-1}b$ .

391 First, note that recursion (11) and (12) can be rewritten as, respectively,

$$\theta_{k+1} \leftarrow \theta_k + \zeta_k x(k),$$

$$u_{k+1} \leftarrow u_k + \beta_k y(k),$$

393 where

$$x(k) = \frac{\alpha_k}{\zeta_k} [(\delta_k - \omega_k)\phi_k - \gamma\phi'_k(\phi_k^\top u_k)],$$

$$y(k) = \frac{\zeta_k}{\beta_k} [\delta_k - \omega_k - \phi_k^\top u_k]\phi_k.$$

395 Recursion (11) can also be rewritten as

$$\theta_{k+1} \leftarrow \theta_k + \beta_k z(k),$$

396 where

$$z(k) = \frac{\alpha_k}{\beta_k} [(\delta_k - \omega_k)\phi_k - \gamma\phi'_k(\phi_k^\top u_k)],$$

397 Due to the settings of step-size schedule  $\alpha_k = o(\zeta_k)$ ,  $\zeta_k = o(\beta_k)$ ,  $x(k) \rightarrow 0$ ,  $y(k) \rightarrow 0$ ,  $z(k) \rightarrow 0$   
 398 almost surely as  $k \rightarrow \infty$ . That is that the increments in iteration (13) are uniformly larger than those  
 399 in (12) and the increments in iteration (12) are uniformly larger than those in (11), thus (13) is the  
 400 fastest recursion, (12) is the second fast recursion and (11) is the slower recursion. Along the fastest  
 401 time scale, iterations of (11), (12) and (13) are associated to ODEs system as follows:

$$\dot{\theta}(t) = 0, \tag{37}$$

$$\dot{u}(t) = 0, \tag{38}$$

$$\dot{\omega}(t) = \mathbb{E}[\delta_t | u(t), \theta(t)] - \omega(t). \tag{39}$$

404 Based on the ODE (37) and (38), both  $\theta(t) \equiv \theta$  and  $u(t) \equiv u$  when viewed from the fastest  
 405 timescale. By the Hirsch lemma Hirsch [1989], it follows that  $\|\theta_k - \theta\| \rightarrow 0$  a.s. as  $k \rightarrow \infty$   
 406 for some  $\theta$  that depends on the initial condition  $\theta_0$  of recursion (11) and  $\|u_k - u\| \rightarrow 0$  a.s. as  
 407  $k \rightarrow \infty$  for some  $u$  that depends on the initial condition  $u_0$  of recursion (12). Thus, the ODE pair  
 408 (37)-(39) can be written as

$$\dot{\omega}(t) = \mathbb{E}[\delta_t | u, \theta] - \omega(t). \tag{40}$$

409 Consider the function  $h(\omega) = \mathbb{E}[\delta | \theta, u] - \omega$ , i.e., the driving vector field of the ODE (40). It is  
 410 easy to find that the function  $h$  is Lipschitz with coefficient  $-1$ . Let  $h_\infty(\cdot)$  be the function defined  
 411 by  $h_\infty(\omega) = \lim_{r \rightarrow \infty} \frac{h(r\omega)}{r}$ . Then  $h_\infty(\omega) = -\omega$ , is well-defined. For (40),  $\omega^* = \mathbb{E}[\delta | \theta, u]$  is the  
 412 unique globally asymptotically stable equilibrium. For the ODE

$$\dot{\omega}(t) = h_\infty(\omega(t)) = -\omega(t), \tag{41}$$

413 apply  $\vec{V}(\omega) = (-\omega)^\top(-\omega)/2$  as its associated strict Liapunov function. Then, the origin of (41) is a  
 414 globally asymptotically stable equilibrium.

415 Consider now the recursion (13). Let  $M_{k+1} = (\delta_k - \omega_k) - \mathbb{E}[(\delta_k - \omega_k) | \mathcal{F}(k)]$ , where  
 416  $\mathcal{F}(k) = \sigma(\omega_l, u_l, \theta_l, l \leq k; \phi_s, \phi'_s, r_s, s < k)$ ,  $k \geq 1$  are the sigma fields generated by  
 417  $\omega_0, u_0, \theta_0, \omega_{l+1}, u_{l+1}, \theta_{l+1}, \phi_l, \phi'_l$ ,  $0 \leq l < k$ . It is easy to verify that  $M_{k+1}, k \geq 0$  are inte-  
 418 grable random variables that satisfy  $\mathbb{E}[M_{k+1} | \mathcal{F}(k)] = 0, \forall k \geq 0$ . Because  $\phi_k, r_k$ , and  $\phi'_k$  have  
 419 uniformly bounded second moments, it can be seen that for some constant  $c_1 > 0, \forall k \geq 0$ ,

$$\mathbb{E}[\|M_{k+1}\|^2 | \mathcal{F}(k)] \leq c_1(1 + \|\omega_k\|^2 + \|u_k\|^2 + \|\theta_k\|^2).$$

Now Assumptions (A1) and (A2) of Borkar and Meyn [2000] are verified. Furthermore, Assumptions (TS) of Borkar and Meyn [2000] is satisfied by our conditions on the step-size sequences  $\alpha_k, \zeta_k, \beta_k$ . Thus, by Theorem 2.2 of Borkar and Meyn [2000] we obtain that  $\|\omega_k - \omega^*\| \rightarrow 0$  almost surely as  $k \rightarrow \infty$ .

Consider now the second time scale recursion (12). Based on the above analysis, (12) can be rewritten as

$$\dot{\theta}(t) = 0, \quad (42)$$

$$\dot{u}(t) = \mathbb{E}[(\delta_t - \mathbb{E}[\delta_t|u(t), \theta(t)])\phi_t|\theta(t)] - Cu(t). \quad (43)$$

The ODE (42) suggests that  $\theta(t) \equiv \theta$  (i.e., a time invariant parameter) when viewed from the second fast timescale. By the Hirsch lemma Hirsch [1989], it follows that  $\|\theta_k - \theta\| \rightarrow 0$  a.s. as  $k \rightarrow \infty$  for some  $\theta$  that depends on the initial condition  $\theta_0$  of recursion (11).

Consider now the recursion (12). Let  $N_{k+1} = ((\delta_k - \mathbb{E}[\delta_k]) - \phi_k \phi_k^\top u_k) - \mathbb{E}[(\delta_k - \mathbb{E}[\delta_k]) - \phi_k \phi_k^\top u_k] | \mathcal{I}(k)$ , where  $\mathcal{I}(k) = \sigma(u_l, \theta_l, l \leq k; \phi_s, \phi'_s, r_s, s < k)$ ,  $k \geq 1$  are the sigma fields generated by  $u_0, \theta_0, u_{l+1}, \theta_{l+1}, \phi_l, \phi'_l, 0 \leq l < k$ . It is easy to verify that  $N_{k+1}, k \geq 0$  are integrable random variables that satisfy  $\mathbb{E}[N_{k+1} | \mathcal{I}(k)] = 0, \forall k \geq 0$ . Because  $\phi_k, r_k$ , and  $\phi'_k$  have uniformly bounded second moments, it can be seen that for some constant  $c_2 > 0, \forall k \geq 0$ ,

$$\mathbb{E}[\|N_{k+1}\|^2 | \mathcal{I}(k)] \leq c_2(1 + \|u_k\|^2 + \|\theta_k\|^2).$$

Because  $\theta(t) \equiv \theta$  from (42), the ODE pair (42)-(43) can be written as

$$\dot{u}(t) = \mathbb{E}[(\delta_t - \mathbb{E}[\delta_t|\theta])\phi_t|\theta] - Cu(t). \quad (44)$$

Now consider the function  $h(u) = \mathbb{E}[\delta_t - \mathbb{E}[\delta_t|\theta]|\theta] - Cu$ , i.e., the driving vector field of the ODE (44). For (44),  $u^* = C^{-1}\mathbb{E}[(\delta - \mathbb{E}[\delta|\theta])\phi|\theta]$  is the unique globally asymptotically stable equilibrium. Let  $h_\infty(u) = -Cu$ . For the ODE

$$\dot{u}(t) = h_\infty(u(t)) = -Cu(t), \quad (45)$$

the origin of (45) is a globally asymptotically stable equilibrium because  $C$  is a positive definite matrix (because it is nonnegative definite and nonsingular). Now Assumptions (A1) and (A2) of Borkar and Meyn [2000] are verified. Furthermore, Assumptions (TS) of Borkar and Meyn [2000] is satisfied by our conditions on the step-size sequences  $\alpha_k, \zeta_k, \beta_k$ . Thus, by Theorem 2.2 of Borkar and Meyn [2000] we obtain that  $\|u_k - u^*\| \rightarrow 0$  almost surely as  $k \rightarrow \infty$ .

Consider now the slower timescale recursion (11). In the light of the above, (11) can be rewritten as

$$\theta_{k+1} \leftarrow \theta_k + \alpha_k(\delta_k - \mathbb{E}[\delta_k|\theta_k])\phi_k - \alpha_k\gamma\phi'_k(\phi_k^\top C^{-1}\mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\theta_k])\phi|\theta_k]). \quad (46)$$

Let  $\mathcal{G}(k) = \sigma(\theta_l, l \leq k; \phi_s, \phi'_s, r_s, s < k)$ ,  $k \geq 1$  be the sigma fields generated by  $\theta_0, \theta_{l+1}, \phi_l, \phi'_l, 0 \leq l < k$ . Let

$$\begin{aligned} Z_{k+1} &= ((\delta_k - \mathbb{E}[\delta_k|\theta_k])\phi_k - \gamma\phi'_k\phi_k^\top C^{-1}\mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\theta_k])\phi|\theta_k]) \\ &\quad - \mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\theta_k])\phi_k - \gamma\phi'_k\phi_k^\top C^{-1}\mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\theta_k])\phi|\theta_k]) | \mathcal{G}(k)] \\ &= ((\delta_k - \mathbb{E}[\delta_k|\theta_k])\phi_k - \gamma\phi'_k\phi_k^\top C^{-1}\mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\theta_k])\phi|\theta_k]) \\ &\quad - \mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\theta_k])\phi_k|\theta_k] - \gamma\mathbb{E}[\phi'_k\phi_k^\top]C^{-1}\mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\theta_k])\phi_k|\theta_k]. \end{aligned}$$

It is easy to see that  $Z_k, k \geq 0$  are integrable random variables and  $\mathbb{E}[Z_{k+1} | \mathcal{G}(k)] = 0, \forall k \geq 0$ . Further,

$$\mathbb{E}[\|Z_{k+1}\|^2 | \mathcal{G}(k)] \leq c_3(1 + \|\theta_k\|^2), k \geq 0$$

for some constant  $c_3 \geq 0$ , again because  $\phi_k, r_k$ , and  $\phi'_k$  have uniformly bounded second moments, it can be seen that for some constant.

Consider now the following ODE associated with (11):

$$\dot{\theta}(t) = (I - \mathbb{E}[\gamma\phi'\phi^\top]C^{-1})\mathbb{E}[(\delta - \mathbb{E}[\delta|\theta(t)])\phi|\theta(t)]. \quad (47)$$

Let

$$\begin{aligned} \vec{h}(\theta(t)) &= (I - \mathbb{E}[\gamma\phi'\phi^\top]C^{-1})\mathbb{E}[(\delta - \mathbb{E}[\delta|\theta(t)])\phi|\theta(t)] \\ &= (C - \mathbb{E}[\gamma\phi'\phi^\top])C^{-1}\mathbb{E}[(\delta - \mathbb{E}[\delta|\theta(t)])\phi|\theta(t)] \\ &= (\mathbb{E}[\phi\phi^\top] - \mathbb{E}[\gamma\phi'\phi^\top])C^{-1}\mathbb{E}[(\delta - \mathbb{E}[\delta|\theta(t)])\phi|\theta(t)] \\ &= A^\top C^{-1}(-A\theta(t) + b), \end{aligned}$$

---

**Algorithm 2** VMTDC algorithm with linear function approximation in the off-policy setting

---

**Input:**  $\theta_0, u_0, \omega_0, \gamma$ , learning rate  $\alpha_t, \zeta_t$  and  $\beta_t$ , behavior policy  $\mu$  and target policy  $\pi$   
**repeat**  
  For any episode, initialize  $\theta_0$  arbitrarily,  $u_t$  and  $\omega_0$  to 0,  $\gamma \in (0, 1]$ , and  $\alpha_t, \zeta_t$  and  $\beta_t$  are constant.  
  **Output:**  $\theta^*$ .  
  **for**  $t = 0$  **to**  $T - 1$  **do**  
    Take  $A_t$  from  $S_t$  according to  $\mu$ , and arrive at  $S_{t+1}$   
    Observe sample  $(S_t, R_{t+1}, S_{t+1})$  at time step  $t$  (with their corresponding state feature vectors)  
     $\delta_t = R_{t+1} + \gamma \theta_t^\top \phi_{t+1} - \theta_t^\top \phi_t$   
     $\rho_t \leftarrow \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)}$   
     $\theta_{t+1} \leftarrow \theta_t + \alpha_t \rho_t [(\delta_t - \omega_t) \phi_t - \gamma \phi_{t+1} (\phi_t^\top u_t)]$   
     $u_{t+1} \leftarrow u_t + \zeta_t [\rho_t (\delta_t - \omega_t) - \phi_t^\top u_t] \phi_t$   
     $\omega_{t+1} \leftarrow \omega_t + \beta_t \rho_t (\delta_t - \omega_t)$   
     $S_t = S_{t+1}$   
  **end for**  
**until** terminal episode

---

---

**Algorithm 3** VMETD algorithm with linear function approximation in the off-policy setting

---

**Input:**  $\theta_0, u_0, \omega_0, \gamma$ , learning rate  $\alpha_t, \zeta_t$  and  $\beta_t$ , behavior policy  $\mu$  and target policy  $\pi$   
**repeat**  
  For any episode, initialize  $\theta_0$  arbitrarily,  $u_t$  to 1 and  $\omega_0$  to 0,  $\gamma \in (0, 1]$ , and  $\alpha_t, \zeta_t$  and  $\beta_t$  are constant.  
  **Output:**  $\theta^*$ .  
  **for**  $t = 0$  **to**  $T - 1$  **do**  
    Take  $A_t$  from  $S_t$  according to  $\mu$ , and arrive at  $S_{t+1}$   
    Observe sample  $(S_t, R_{t+1}, S_{t+1})$  at time step  $t$  (with their corresponding state feature vectors)  
     $\delta_t = R_{t+1} + \gamma \theta_t^\top \phi_{t+1} - \theta_t^\top \phi_t$   
     $\rho_t \leftarrow \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)}$   
     $F_t \leftarrow \gamma \rho_t F_{t-1} + 1$   
     $\theta_{t+1} \leftarrow \theta_t + \alpha_t (F_t \rho_t \delta_t - \omega_t) \phi_t$   
     $\omega_{t+1} \leftarrow \omega_t + \beta_t (F_t \rho_t \delta_t - \omega_t)$   
     $S_t = S_{t+1}$   
  **end for**  
**until** terminal episode

---

453 because  $\mathbb{E}[(\delta - \mathbb{E}[\delta|\theta(t)])\phi|\theta(t)] = -A\theta(t) + b$ , where  $A = \text{Cov}(\phi, \phi - \gamma\phi')$ ,  $b = \text{Cov}(r, \phi)$ , and  
454  $C = \mathbb{E}[\phi\phi^\top]$

455 Therefore,  $\theta^* = A^{-1}b$  can be seen to be the unique globally asymptotically stable equilibrium for  
456 ODE (47). Let  $\vec{h}_\infty(\theta) = \lim_{r \rightarrow \infty} \frac{\vec{h}(r\theta)}{r}$ . Then  $\vec{h}_\infty(\theta) = -A^\top C^{-1}A\theta$  is well-defined. Consider  
457 now the ODE

$$\dot{\theta}(t) = -A^\top C^{-1}A\theta(t). \quad (48)$$

458 Because  $C^{-1}$  is positive definite and  $A$  has full rank (as it is nonsingular by assumption), the matrix  
459  $A^\top C^{-1}A$  is also positive definite. The ODE (48) has the origin as its unique globally asymptotically  
460 stable equilibrium. Thus, the assumption (A1) and (A2) are verified.

461 The proof is given above. In the fastest time scale, the parameter  $w$  converges to  $\mathbb{E}[\delta|u_k, \theta_k]$ . In the  
462 second fast time scale, the parameter  $u$  converges to  $C^{-1}\mathbb{E}[(\delta - \mathbb{E}[\delta|\theta_k])\phi|\theta_k]$ . In the slower time  
463 scale, the parameter  $\theta$  converges to  $A^{-1}b$ .  $\square$

464 **A.4 Proof of VMETD convergence**

465 VMETD's  $\theta$  by the following update:

$$\begin{aligned}
 \theta_{k+1} &\leftarrow \theta_k + \alpha_k F_k \rho_k (R_{k+1} + \gamma \theta_k^\top \phi_{k+1} - \theta_k^\top \phi_k) \phi_k - \alpha_k \omega_{k+1} \phi_k \\
 &= \theta_k + \alpha_k F_k \rho_k (R_{k+1} + \gamma \theta_k^\top \phi_{k+1} - \theta_k^\top \phi_k) \phi_k - \alpha_k \mathbb{E}_\mu[F_k \rho_k \delta_k] \phi_k \\
 &= \theta_k + \alpha_k \underbrace{\{ (F_k \rho_k R_{k+1} - \mathbb{E}_\mu[F_k \rho_k R_{k+1}]) \phi_k \}}_{\mathbf{b}_{\text{VMETD},k}} - \underbrace{\{ (F_k \rho_k \phi_k (\phi_k - \gamma \phi_{k+1})^\top - \phi_k \mathbb{E}_\mu[F_k \rho_k (\phi_k - \gamma \phi_{k+1})]^\top) \theta_k \}}_{\mathbf{A}_{\text{VMETD},k}}
 \end{aligned} \tag{49}$$

466

$$\begin{aligned}
 \mathbf{A}_{\text{VMETD}} &= \lim_{k \rightarrow \infty} \mathbb{E}[\mathbf{A}_{\text{VMETD},k}] \\
 &= \lim_{k \rightarrow \infty} \mathbb{E}_\mu[F_k \rho_k \phi_k (\phi_k - \gamma \phi_{k+1})^\top] - \lim_{k \rightarrow \infty} \mathbb{E}_\mu[\phi_k] \mathbb{E}_\mu[F_k \rho_k (\phi_k - \gamma \phi_{k+1})]^\top \\
 &= \lim_{k \rightarrow \infty} \mathbb{E}_\mu \left[ \underbrace{\phi_k}_X \underbrace{F_k \rho_k (\phi_k - \gamma \phi_{k+1})^\top}_Y \right] - \lim_{k \rightarrow \infty} \mathbb{E}_\mu[\phi_k] \mathbb{E}_\mu[F_k \rho_k (\phi_k - \gamma \phi_{k+1})]^\top \\
 &= \lim_{k \rightarrow \infty} \mathbb{E}_\mu[F_k \rho_k \phi_k (\phi_k - \gamma \phi_{k+1})^\top] - \lim_{k \rightarrow \infty} \mathbb{E}_\mu[\phi_k] \lim_{k \rightarrow \infty} \mathbb{E}_\mu[F_k \rho_k (\phi_k - \gamma \phi_{k+1})]^\top \\
 &= \sum_s f(s) \phi(s) (\phi(s) - \gamma \sum_{s'} [\mathbf{P}_\pi]_{ss'} \phi(s'))^\top - \sum_s d_\mu(s) \phi(s) * \sum_s f(s) (\phi(s) - \gamma \sum_{s'} [\mathbf{P}_\pi]_{ss'} \phi(s'))^\top \\
 &= \Phi^\top \mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) \Phi - \Phi^\top \mathbf{d}_\mu \mathbf{f}^\top (\mathbf{I} - \gamma \mathbf{P}_\pi) \Phi \\
 &= \Phi^\top (\mathbf{F} - \mathbf{d}_\mu \mathbf{f}^\top) (\mathbf{I} - \gamma \mathbf{P}_\pi) \Phi \\
 &= \Phi^\top (\mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) - \mathbf{d}_\mu \mathbf{f}^\top (\mathbf{I} - \gamma \mathbf{P}_\pi)) \Phi \\
 &= \Phi^\top (\mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) - \mathbf{d}_\mu \mathbf{d}_\mu^\top) \Phi
 \end{aligned} \tag{50}$$

467 *Proof.* Any matrix  $\mathbf{M}$  is positive definite if and only if the symmetric matrix  $\mathbf{S} = \mathbf{M} + \mathbf{M}^\top$  is positive  
 468 definite. Any symmetric real matrix  $\mathbf{S}$  is positive definite if the absolute values of its diagonal entries  
 469 are greater than the sum of the absolute values of the corresponding off-diagonal entries.

$$\begin{aligned}
 (\mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) - \mathbf{d}_\mu \mathbf{d}_\mu^\top) \mathbf{1} &= \mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) \mathbf{1} - \mathbf{d}_\mu \mathbf{d}_\mu^\top \mathbf{1} \\
 &= \mathbf{F}(\mathbf{1} - \gamma \mathbf{P}_\pi \mathbf{1}) - \mathbf{d}_\mu \mathbf{d}_\mu^\top \mathbf{1} \\
 &= (1 - \gamma) \mathbf{F} \mathbf{1} - \mathbf{d}_\mu \mathbf{d}_\mu^\top \mathbf{1} \\
 &= (1 - \gamma) \mathbf{f} - \mathbf{d}_\mu \mathbf{d}_\mu^\top \mathbf{1} \\
 &= (1 - \gamma) \mathbf{f} - \mathbf{d}_\mu \\
 &= (1 - \gamma) (\mathbf{I} - \gamma \mathbf{P}_\pi^\top)^{-1} \mathbf{d}_\mu - \mathbf{d}_\mu \\
 &= (1 - \gamma) [(\mathbf{I} - \gamma \mathbf{P}_\pi^\top)^{-1} - \mathbf{I}] \mathbf{d}_\mu \\
 &= (1 - \gamma) \left[ \sum_{t=0}^{\infty} (\gamma \mathbf{P}_\pi^\top)^t - \mathbf{I} \right] \mathbf{d}_\mu \\
 &= (1 - \gamma) \left[ \sum_{t=1}^{\infty} (\gamma \mathbf{P}_\pi^\top)^t \right] \mathbf{d}_\mu > 0
 \end{aligned} \tag{51}$$

470

$$\begin{aligned}
 \mathbf{1}^\top (\mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) - \mathbf{d}_\mu \mathbf{d}_\mu^\top) &= \mathbf{1}^\top \mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) - \mathbf{1}^\top \mathbf{d}_\mu \mathbf{d}_\mu^\top \\
 &= \mathbf{d}_\mu^\top - \mathbf{1}^\top \mathbf{d}_\mu \mathbf{d}_\mu^\top \\
 &= \mathbf{d}_\mu^\top - \mathbf{d}_\mu^\top \\
 &= 0
 \end{aligned} \tag{52}$$

471 (51) and (52) show that the matrix  $\mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) - \mathbf{d}_\mu \mathbf{d}_\mu^\top$  of diagonal entries are positive and its  
 472 off-diagonal entries are negative. So its each row sum plus the corresponding column sum is positive.

473 The proof is given above  $\square$

## B Experimental details

The feature matrices corresponding to three random walks are shown below respectively:

$$\Phi_{\text{tabular}} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\Phi_{\text{inverted}} = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix}$$

$$\Phi_{\text{dependent}} = \begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \\ \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} \\ 0 & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ 0 & 0 & 1 \end{bmatrix}$$

Three random walk experiments: the  $\alpha$  values for all algorithms are in the range of  $\{0.008, 0.015, 0.03, 0.06, 0.12, 0.25, 0.5\}$ . For the TDC algorithm, the range of the ratio  $\frac{\zeta}{\alpha}$  is  $\{\frac{1}{512}, \frac{1}{256}, \frac{1}{128}, \frac{1}{64}, \frac{1}{32}, \frac{1}{16}, \frac{1}{8}, \frac{1}{4}, \frac{1}{2}, 1, 2\}$ . For the VMTD algorithm, the range of the ratio  $\frac{\beta}{\alpha}$  is  $\{\frac{1}{512}, \frac{1}{256}, \frac{1}{128}, \frac{1}{64}, \frac{1}{32}, \frac{1}{16}, \frac{1}{8}, \frac{1}{4}, \frac{1}{2}, 1, 2\}$ . It can be observed from the update formula of VMTDC that when  $\zeta$  takes a very small value, the VMTDC update tends to be similar to VMTD update. Similarly, when  $\beta$  takes a very small value, the VMTDC update tends to be similar to TDC update. Through experiments, it was found that setting  $\zeta$  to a small value makes VMTDC updates approach VMTD updates, resulting in better performance. Therefore, for the VMTDC algorithm, the range of  $\frac{\beta}{\alpha}$  ratio is  $\{\frac{1}{512}, \frac{1}{256}, \frac{1}{128}, \frac{1}{64}, \frac{1}{32}, \frac{1}{16}, \frac{1}{8}, \frac{1}{4}, \frac{1}{2}, 1, 2\}$ , and the range of  $\zeta$  is  $\{0.1, 0.01, 0.001, 0.0001, 0.00001\}$ . The learning curves in Figure 3 correspond to the optimal parameters.

The feature matrix of 7-state version of Baird’s off-policy counterexample is defined as follow:

$$\Phi_{\text{Counter}} = \begin{bmatrix} 1 & 2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 2 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 2 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 2 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 2 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 2 & 0 \\ 2 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

7-state version of Baird’s off-policy counterexample: for TD algorithm,  $\alpha$  is set to 0.1. For the TDC algorithm, the range of  $\alpha$  is  $\{0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$ , and the range of  $\zeta$  is  $\{0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0, 1.1, 1.2, 1.3, 1.4, 1.5\}$ . For the VMTD algorithm, the range of  $\alpha$  is  $\{0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$ , and the range of  $\beta$  is  $\{0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0, 1.1, 1.2, 1.3, 1.4, 1.5\}$ . Through experiments, it was found that setting  $\zeta$  to a small value makes VMTDC updates approach VMTD updates, resulting in better performance. Therefore, for the VMTDC algorithm, The range of values for  $\alpha$  and  $\beta$  is the same as that of VMTD and the range of  $\zeta$  is  $\{0.1, 0.01, 0.001, 0.0001, 0.00001\}$ . The learning curves in Figure 4 correspond to the optimal parameters. For all policy evaluation experiments, each experiment is independently run 100 times.

For the four control experiments: The learning rates for each algorithm in all experiments are shown in Table 3. For all control experiments, each experiment is independently run 50 times.

## References

Leemon Baird et al. Residual algorithms: Reinforcement learning with function approximation. In *Proc. 12th Int. Conf. Mach. Learn.*, pages 30–37, 1995.

Table 3: Learning rates ( $lr$ ) of four control experiments.

algorithms( $lr$ ) \ envs	Maze	Cliff walking	Mountain Car	Acrobot
Sarsa( $\alpha$ )	0.1	0.1	0.1	0.1
GQ(0)( $\alpha, \zeta$ )	0.1, 0.003	0.1, 0.004	0.1, 0.01	0.1, 0.01
VMSarsa( $\alpha, \beta$ )	0.1, 0.001	0.1, 1e-4	0.1, 1e-4	0.1, 1e-4
VMGQ(0)( $\alpha, \zeta, \beta$ )	0.1, 0.001, 0.001	0.1, 0.005, 1e-4	0.1, 5e-4, 1e-4	0.1, 5e-4, 1e-4
AC( $lr_{actor}, lr_{critic}$ )	0.01, 0.1	0.01, 0.01	0.01, 0.05	0.01, 0.05
Q-learning( $\alpha$ )	0.1	0.1	0.1	0.1
VMQ( $\alpha, \beta$ )	0.1, 0.001	0.1, 1e-4	0.1, 1e-4	0.1, 1e-4

- 504 Joan Bas-Serrano, Sebastian Curi, Andreas Krause, and Gergely Neu. Logistic q-learning. In  
505 *International Conference on Artificial Intelligence and Statistics*, pages 3610–3618, 2021.
- 506 Vivek S Borkar and Sean P Meyn. The ode method for convergence of stochastic approximation and  
507 reinforcement learning. *SIAM J. Control Optim.*, 38(2):447–469, 2000.
- 508 Vivek S Borkar. Stochastic approximation with two time scales. *Syst. & Control Letters*, 29(5):291–  
509 294, 1997.
- 510 Xingguo Chen, Xingzhou Ma, Yang Li, Guang Yang, Shangdong Yang, and Yang Gao. Modified  
511 retrace for off-policy temporal difference learning. In *Uncertainty in Artificial Intelligence*, pages  
512 303–312. PMLR, 2023.
- 513 Gal Dalal, Balazs Szorenyi, and Gudan Thoppe. A tale of two-timescale reinforcement learning with  
514 the tightest finite-time bound. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
515 volume 34, pages 3701–3708, 2020.
- 516 Sam Devlin and Daniel Kudenko. Dynamic potential-based reward shaping. In *Proc. 11th Int. Conf.*  
517 *Autonomous Agents and Multiagent Systems*, pages 433–440, 2012.
- 518 Yihao Feng, Lihong Li, and Qiang Liu. A kernel loss for solving the bellman equation. In *Advances*  
519 *in Neural Information Processing Systems*, pages 15430–15441, 2019.
- 520 Arash Givchi and Maziar Palhang. Quasi newton temporal difference learning. In *Asian Conference*  
521 *on Machine Learning*, pages 159–172, 2015.
- 522 Leah Hackman. *Faster Gradient-TD Algorithms*. PhD thesis, University of Alberta, 2012.
- 523 Assaf Hallak, Aviv Tamar, Remi Munos, and Shie Mannor. Generalized emphatic temporal difference  
524 learning: bias-variance analysis. In *Proceedings of the 30th AAAI Conference on Artificial*  
525 *Intelligence*, pages 1631–1637, 2016.
- 526 Morris W Hirsch. Convergent activation dynamics in continuous time networks. *Neural Netw.*,  
527 2(5):331–349, 1989.
- 528 R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction.  
529 In *Advances in Neural Information Processing Systems*, pages 315–323, 2013.
- 530 Nathaniel Korda and Prashanth La. On td (0) with function approximation: Concentration bounds  
531 and a centered variant with exponential convergence. In *International conference on machine*  
532 *learning*, pages 626–634. PMLR, 2015.
- 533 Bo Liu, Ji Liu, Mohammad Ghavamzadeh, Sridhar Mahadevan, and Marek Petrik. Finite-sample  
534 analysis of proximal gradient td algorithms. In *Proceedings of the 21st Conference on Uncertainty*  
535 *in Artificial Intelligence*, pages 504–513, 2015.
- 536 Bo Liu, Ji Liu, Mohammad Ghavamzadeh, Sridhar Mahadevan, and Marek Petrik. Proximal gradient  
537 temporal difference learning algorithms. In *Proceedings of the International Joint Conference on*  
538 *Artificial Intelligence*, pages 4195–4199, 2016.

539 Bo Liu, Ian Gemp, Mohammad Ghavamzadeh, Ji Liu, Sridhar Mahadevan, and Marek Petrik.  
540 Proximal gradient temporal difference learning: Stable reinforcement learning with polynomial  
541 sample complexity. *Journal of Artificial Intelligence Research*, 63:461–494, 2018.

542 Hamid Reza Maei. *Gradient temporal-difference learning algorithms*. PhD thesis, University of  
543 Alberta, 2011.

544 Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations:  
545 Theory and application to reward shaping. In *Proc. 16th Int. Conf. Mach. Learn.*, pages 278–287,  
546 1999.

547 Yangchen Pan, Adam White, and Martha White. Accelerated gradient temporal difference learning.  
548 In *Proceedings of the 21st AAAI Conference on Artificial Intelligence*, pages 2464–2470, 2017.

549 J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In  
550 *International Conference on Machine Learning*, pages 1889–1897, 2015.

551 J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization  
552 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

553 Anton Schwartz. A reinforcement learning method for maximizing undiscounted rewards. In *Proc.*  
554 *10th Int. Conf. Mach. Learn.*, volume 298, pages 298–305, 1993.

555 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press,  
556 second edition, 2018.

557 Richard S Sutton, Hamid R Maei, and Csaba Szepesvári. A convergent  $o(n)$  temporal-difference  
558 algorithm for off-policy learning with linear function approximation. In *Advances in Neural*  
559 *Information Processing Systems*, pages 1609–1616. Cambridge, MA: MIT Press, 2008.

560 R.S. Sutton, H.R. Maei, D. Precup, S. Bhatnagar, D. Silver, C. Szepesvári, and E. Wiewiora. Fast  
561 gradient-descent methods for temporal-difference learning with linear function approximation. In  
562 *Proc. 26th Int. Conf. Mach. Learn.*, pages 993–1000, 2009.

563 Richard S Sutton, A Rupam Mahmood, and Martha White. An emphatic approach to the problem of  
564 off-policy temporal-difference learning. *The Journal of Machine Learning Research*, 17(1):2603–  
565 2631, 2016.

566 Richard S Sutton. Learning to predict by the methods of temporal differences. *Machine learning*,  
567 3(1):9–44, 1988.

568 John N Tsitsiklis and Benjamin Van Roy. Analysis of temporal-difference learning with function  
569 approximation. In *Advances in Neural Information Processing Systems*, pages 1075–1081, 1997.

570 Tengyu Xu, Zhe Wang, Yi Zhou, and Yingbin Liang. Reanalysis of variance reduced temporal  
571 difference learning. In *International Conference on Learning Representations*, 2019.

572 T. Xu, Z. Wang, Y. Zhou, and Y. Liang. Reanalysis of variance reduced temporal difference learning.  
573 *arXiv preprint arXiv:2001.01898*, 2020.

574 Shangdong Zhang and Shimon Whiteson. Truncated emphatic temporal difference methods for  
575 prediction and control. *The Journal of Machine Learning Research*, 23(1):6859–6917, 2022.