

AAAI Press Anonymous Submission

Instructions for Authors Using L^AT_EX

Anonymous submission

Abstract

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

Introduction

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

In terms of stability, Sutton (1988) established the convergence of on-policy TD(0), and Tsitsiklis and Van Roy (1997) established the convergence of on-policy TD(λ). However, “The deadly triad” consisting of off-policy learning, bootstrapping, and function approximation makes the stability a difficult problem (Sutton and Barto 2018). To solve this problem, convergent off-policy temporal difference learning algorithms are proposed, e.g., BR (Baird et al. 1995), GTD (Sutton, Maei, and Szepesvári 2008), GTD2 and TDC (Sutton et al. 2009), ETD (Sutton, Mahmood, and White 2016), and MReTrace (Chen et al. 2023).

In terms of acceleration, Hackman (2012) proposed Hybrid TD algorithm with on-policy matrix. Liu et al. (2015, 2016, 2018) proposed true stochastic algorithms, i.e., GTD-MP and GTD2-MP, from a convex-concave saddle-point formulation. Second-order methods are used to accelerate TD learning, e.g., Quasi Newton TD (Givchi and Palhang 2015) and accelerated TD (ATD) (Pan, White, and White 2017). Hallak et al. (2016) introduced an new parameter to reduce variance for ETD. Zhang and Whiteson (2022) proposed truncated ETD with a lower variance. Variance Reduced TD with direct variance reduction technique (Johnson and

Zhang 2013) is proposed by (Korda and La 2015) and analysed by (Xu et al. 2019). How to further improve the convergence rates of reinforcement learning algorithms is currently still an open problem.

Algorithm stability is prominently reflected in the changes to the objective function, transitioning from mean squared errors (MSE) (Sutton and Barto 2018) to mean squared bellman errors (MSBE) (Baird et al. 1995), then to norm of the expected TD update (Sutton et al. 2009), and further to 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, Li, and Liu 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) and Variance of Projected Bellman Error (VPBE) 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 three variance minimization algorithms, including on-policy and 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.

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 π^* to maximize the expectation of a discounted cumulative rewards in a long period. State value function $V^\pi(s)$ for a stationary policy π is defined as:

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

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, 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.

Variance Minimization Algorithms

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

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

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, Harada, and Russell 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 ω 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 ω requires expert knowledge. This forces us to learn ω 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 ω 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 $\omega \doteq \mathbb{E}[\mathbb{E}[\delta|s]] = \mathbb{E}[\delta]$.

Variance Minimization TD Learning: VMTD

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

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

Clearly, it is no longer to minimize Bellman errors.

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

Input: $\theta_0, \omega_0, \gamma$, learning rate α_t and β_t
repeat
 For any episode, initialize θ_0 arbitrarily, ω_0 to 0, $\gamma \in (0, 1]$, and α_t and β_t are constant.
for $t = 0$ **to** $T - 1$ **do**
 Take A_t from S_t according to policy π , 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

First, the parameter ω is derived directly based on stochastic gradient descent:

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

where δ_k is the TD error as follows:

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

Then, based on stochastic semi-gradient descent, the update of the parameter θ is as follows:

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

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

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

and

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

where δ_k delta in VMSarsa is:

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

and δ_k delta in VMQ is:

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

Variance Minimization TDC Learning: VMTDC

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

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

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

The derivation process of the VMTDC algorithm is the same as that of the TDC algorithm, the only difference is

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

Input: $\theta_0, \mathbf{u}_0, \omega_0, \gamma$, learning rate α_t, ζ_t and β_t , behavior policy μ and target policy π
repeat
 For any episode, initialize θ_0 arbitrarily, \mathbf{u}_0 and ω_0 to 0, $\gamma \in (0, 1]$, and α_t, ζ_t and β_t are constant.
for $t = 0$ **to** $T - 1$ **do**
 Take A_t from S_t according to μ , 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)
 $\delta_t = R_{t+1} + \gamma \theta_t^\top \phi_{t+1} - \theta_t^\top \phi_t$
 $\rho_t \leftarrow \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)}$
 $\theta_{t+1} \leftarrow \theta_t + \alpha_t \rho_t [(\delta_t - \omega_t) \phi_t - \gamma \phi_{t+1} (\phi_t^\top \mathbf{u}_t)]$
 $\mathbf{u}_{t+1} \leftarrow \mathbf{u}_t + \zeta_t [\rho_t (\delta_t - \omega_t) - \phi_t^\top \mathbf{u}_t] \phi_t$
 $\omega_{t+1} \leftarrow \omega_t + \beta_t \rho_t (\delta_t - \omega_t)$
 $S_t = S_{t+1}$
end for
until terminal episode

that the original δ is replaced by $\delta - \omega$. Therefore, we can easily get the updated formula of VMTDC, as follows:

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

$$\mathbf{u}_{k+1} \leftarrow \mathbf{u}_k + \zeta_k [\delta_k - \omega_k - \phi^\top(s_k) \mathbf{u}_k] \phi(s_k), \quad (12)$$

and

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

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

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

$$\mathbf{u}_{k+1} \leftarrow \mathbf{u}_k + \zeta_k [(\delta_k - \mathbf{u}_k) - \phi^\top(s_k, a_k) \mathbf{u}_k] \phi(s_k, a_k), \quad (15)$$

and

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

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

Variance Minimization ETD Learning: VMETD

Based on the off-policy TD algorithm, a scalar, F , is introduced to obtain the ETD algorithm, which ensures convergence under off-policy conditions. This paper further introduces a scalar, ω , based on the ETD algorithm to obtain VMETD. VMETD by the following update:

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

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

$$\theta_{k+1} \leftarrow \theta_k + \alpha_k (F_k \rho_k \delta_k - \omega_k) \phi_k, \quad (19)$$

$$\omega_{k+1} \leftarrow \omega_k + \beta_k (F_k \rho_k \delta_k - \omega_k), \quad (20)$$

where μ is behavior policy and π is target policy, F_t is a scalar variable, $F_0 = 1$, ω is used to estimate $\mathbb{E}[\delta]$, i.e.,

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

Input: $\theta_0, F_0, \omega_0, \gamma$, learning rate α_t, ζ_t and β_t , behavior policy μ and target policy π

repeat

For any episode, initialize θ_0 arbitrarily, F_0 to 1 and ω_0 to 0, $\gamma \in (0, 1]$, and α_t, ζ_t and β_t are constant.

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

Take A_t from S_t according to μ , 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

$\omega \doteq \mathbb{E}[\delta]$, and \mathbf{F} is a diagonal matrix with diagonal elements $f(s) \doteq d_\mu(s) \lim_{t \rightarrow \infty} \mathbb{E}_\mu[F_t | S_t = s]$, which we assume exists. The vector $\mathbf{f} \in \mathbb{R}^N$ with components $[\mathbf{f}]_s \doteq f(s)$ can be written as

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

Theoretical Analysis

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

Theorem 1. (Convergence of VMTD). *In the case of on-policy learning, consider the iterations (3) and (5) with (4) of VMTD. Let the step-size sequences α_k and β_k , $k \geq 0$ satisfy in this 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 $\alpha_k = o(\beta_k)$. Assume that (ϕ_k, r_k, ϕ'_k) is an i.i.d. sequence with uniformly bounded second moments, where ϕ_k and ϕ'_k are sampled from the same Markov chain. Let $\mathbf{A} = \text{Cov}(\phi, \phi - \gamma \phi')$, $\mathbf{b} = \text{Cov}(r, \phi)$. Assume that matrix θ is non-singular. Then the parameter vector θ_k converges with probability one to $\mathbf{A}^{-1} \mathbf{b}$.*

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:

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

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

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:

$$\dot{\theta}(t) = 0, \quad (22)$$

$$\dot{\omega}(t) = \mathbb{E}[\delta_t | \theta(t)] - \omega(t). \quad (23)$$

Based on the ODE (22), $\theta(t) \equiv \theta$ 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 θ that depends on the initial condition θ_0 of recursion (5). Thus, the ODE pair (22)-(23) can be written as

$$\dot{\omega}(t) = \mathbb{E}[\delta_t | \theta] - \omega(t). \quad (24)$$

Consider the function $h(\omega) = \mathbb{E}[\delta | \theta] - \omega$, i.e., the driving vector field of the ODE (24). 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 (24), $\omega^* = \mathbb{E}[\delta | \theta]$ is the unique globally asymptotically stable equilibrium. For the ODE

$$\dot{\omega}(t) = h_\infty(\omega(t)) = -\omega(t), \quad (25)$$

apply $\vec{V}(\omega) = (-\omega)^\top (-\omega)/2$ as its associated strict Liapunov function. Then, the origin of (25) 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 ϕ_k, r_k , and ϕ'_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 α_k, β_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 $Z_{k+1} = Y_k - \mathbb{E}[Y_k | \mathcal{G}(k)]$, where

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

Consequently,

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

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 ϕ_k, r_k , and ϕ'_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):

$$\begin{aligned} \dot{\theta}(t) &= \text{Cov}(\delta | \theta(t), \phi) \\ &= \text{Cov}(r + (\gamma\phi' - \phi)^\top \theta(t), \phi) \\ &= \text{Cov}(r, \phi) - \text{Cov}(\theta(t)^\top (\phi - \gamma\phi'), \phi) \\ &= \text{Cov}(r, \phi) - \theta(t)^\top \text{Cov}(\phi - \gamma\phi', \phi) \\ &= \text{Cov}(r, \phi) - \text{Cov}(\phi - \gamma\phi', \phi)^\top \theta(t) \\ &= \text{Cov}(r, \phi) - \text{Cov}(\phi, \phi - \gamma\phi')\theta(t) \\ &= -\mathbf{A}\theta(t) + \mathbf{b}. \end{aligned} \quad (26)$$

Let $\vec{h}(\theta(t))$ be the driving vector field of the ODE (26).

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

Consider the cross-covariance matrix,

$$\begin{aligned} \mathbf{A} &= \text{Cov}(\phi, \phi - \gamma\phi') \\ &= \frac{\text{Cov}(\phi, \phi) + \text{Cov}(\phi - \gamma\phi', \phi - \gamma\phi') - \text{Cov}(\gamma\phi', \gamma\phi')}{2} \\ &= \frac{\text{Cov}(\phi, \phi) + \text{Cov}(\phi - \gamma\phi', \phi - \gamma\phi') - \gamma^2 \text{Cov}(\phi', \phi')}{2} \\ &= \frac{(1 - \gamma^2) \text{Cov}(\phi, \phi) + \text{Cov}(\phi - \gamma\phi', \phi - \gamma\phi')}{2}, \end{aligned} \quad (27)$$

where we eventually used $\text{Cov}(\phi', \phi') = \text{Cov}(\phi, \phi)$ ¹. 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 \mathbf{A} is semi-positive definite because \mathbf{A} is linearly combined by two positive-weighted semi-positive definite matrices (27). Furthermore, \mathbf{A} is nonsingular due to the assumption. Hence, the cross-covariance matrix \mathbf{A} is positive definite.

Therefore, $\theta^* = \mathbf{A}^{-1}\mathbf{b}$ can be seen to be the unique globally asymptotically stable equilibrium for ODE (26). Let $\vec{h}_\infty(\theta) = \lim_{r \rightarrow \infty} \frac{\vec{h}(r\theta)}{r}$. Then $\vec{h}_\infty(\theta) = -\mathbf{A}\theta$ is well-defined. Consider now the ODE

$$\dot{\theta}(t) = -\mathbf{A}\theta(t). \quad (28)$$

The ODE (28) has the origin as its unique globally asymptotically stable equilibrium. Thus, the assumption (A1) and (A2) are verified. \square

Theorem 2. (Convergence of VMTDC). *In the case of off-policy learning, consider the iterations (13), (12) and (11) of VMTDC. Let the step-size sequences α_k, ζ_k and $\beta_k, k \geq 0$ satisfy in this case $\alpha_k, \zeta_k, \beta_k > 0$, for all k , $\sum_{k=0}^{\infty} \alpha_k = \sum_{k=0}^{\infty} \beta_k = \sum_{k=0}^{\infty} \zeta_k = \infty$, $\sum_{k=0}^{\infty} \alpha_k^2 < \infty$, $\sum_{k=0}^{\infty} \zeta_k^2 < \infty$, $\sum_{k=0}^{\infty} \beta_k^2 < \infty$, and $\alpha_k = o(\zeta_k), \zeta_k = o(\beta_k)$. Assume that (ϕ_k, r_k, ϕ'_k) is an i.i.d. sequence with uniformly*

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

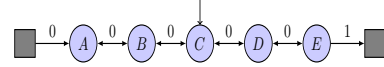


Figure 1: Random walk.

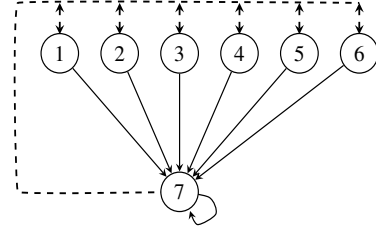


Figure 2: 7-state version of Baird's off-policy counterexample.

bounded second moments. Let $\mathbf{A} = \text{Cov}(\phi, \phi - \gamma\phi')$, $\mathbf{b} = \text{Cov}(r, \phi)$, and $\mathbf{C} = \mathbb{E}[\phi\phi^\top]$. Assume that \mathbf{A} and \mathbf{C} are non-singular matrices. Then the parameter vector θ_k converges with probability one to $\mathbf{A}^{-1}\mathbf{b}$.

Theorem 3. (Convergence of VMETD). *In the case of off-policy learning, consider the iterations (18), (20) and (19) with (4) of VMETD. Let the step-size sequences α_k and $\beta_k, k \geq 0$ satisfy in this 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 $\alpha_k = o(\beta_k)$. Assume that (ϕ_k, r_k, ϕ'_k) is an i.i.d. sequence with uniformly bounded second moments, where ϕ_k and ϕ'_k are sampled from the same Markov chain. Let $\mathbf{A}_{\text{VMETD}} = \Phi^\top (\mathbf{F}(\mathbf{I} - \gamma\mathbf{P}_\pi) - \mathbf{d}_\mu \mathbf{d}_\mu^\top) \Phi$, $\mathbf{b}_{\text{VMETD}} = \Phi^\top (\mathbf{F} - \mathbf{d}_\mu \mathbf{d}_\mu^\top) \mathbf{r}_\pi$. Assume that matrix \mathbf{A} is non-singular. Then the parameter vector θ_k converges with probability one to $\mathbf{A}_{\text{VMETD}}^{-1} \mathbf{b}_{\text{VMETD}}$.*

Experimental Studies

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

Testing Tasks

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

Baird's off-policy counterexample: This task is well known as a counterexample, in which TD diverges (Baird

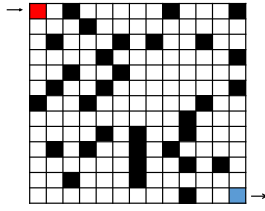


Figure 3: Maze.

et al. 1995; Sutton et al. 2009). As shown in Figure 2, reward for each transition is zero. Thus the true values are zeros for all states and for any given policy. The behaviour policy chooses actions represented by solid lines with a probability of $\frac{1}{7}$ and actions represented by dotted lines with a probability of $\frac{6}{7}$. The target policy is expected to choose the solid line with more probability than $\frac{1}{7}$, and it chooses the solid line with probability of 1 in this paper. The discount factor $\gamma = 0.99$, and the feature matrix is defined in Appendix.

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

The other three control environments: Cliff Walking, Mountain Car, and Acrobot are selected from the gym official website and correspond to the following versions: “CliffWalking-v0”, “MountainCar-v0” and “Acrobot-v1”. For specific details, please refer to the gym official website. The maximum number of steps for the Mountain Car environment is set to 1000, while the default settings are used for the other two environments. In Mountain car and Acrobot, features are generated by tile coding.

Please, refer to the Appendix for the selection of learning rates for all experiments.

Experimental Results and Analysis

The experiment needs further elaboration.

Related Work

Difference between VMQ and R-learning

Tabular VMQ’s update formula bears some resemblance to R-learning’s update formula. As shown in Table ??, 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 ω in the tabular VMQ update formula eventually 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.

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.

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.

Conclusion and Future Work

Value-based reinforcement learning typically aims to minimize error as an optimization objective. As an alternative, this study proposes new objective functions: VBE and VPBE, and derives many variance minimization algorithms, including VMTD, VMTDC 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 and VMETD. (2) extensions of VBE and VPBE to multi-step returns. (3) extensions to nonlinear approximations, such as neural networks.

References

- Baird, L.; et al. 1995. Residual algorithms: Reinforcement learning with function approximation. In *Proc. 12th Int. Conf. Mach. Learn.*, 30–37.
- Bas-Serrano, J.; Curi, S.; Krause, A.; and Neu, G. 2021. Logistic Q-Learning. In *International Conference on Artificial Intelligence and Statistics*, 3610–3618.
- Borkar, V. S. 1997. Stochastic approximation with two time scales. *Syst. & Control Letters*, 29(5): 291–294.
- Borkar, V. S.; and Meyn, S. P. 2000. The ODE method for convergence of stochastic approximation and reinforcement learning. *SIAM J. Control Optim.*, 38(2): 447–469.
- Chen, X.; Ma, X.; Li, Y.; Yang, G.; Yang, S.; and Gao, Y. 2023. Modified Retrace for Off-Policy Temporal Difference Learning. In *Uncertainty in Artificial Intelligence*, 303–312. PMLR.
- Devlin, S.; and Kudenko, D. 2012. Dynamic potential-based reward shaping. In *Proc. 11th Int. Conf. Autonomous Agents and Multiagent Systems*, 433–440.
- Feng, Y.; Li, L.; and Liu, Q. 2019. A kernel loss for solving the Bellman equation. In *Advances in Neural Information Processing Systems*, 15430–15441.
- Givchi, A.; and Palhang, M. 2015. Quasi newton temporal difference learning. In *Asian Conference on Machine Learning*, 159–172.
- Hackman, L. 2012. *Faster Gradient-TD Algorithms*. Ph.D. thesis, University of Alberta.
- Hallak, A.; Tamar, A.; Munos, R.; and Mannor, S. 2016. Generalized emphatic temporal difference learning: bias-variance analysis. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence*, 1631–1637.
- Hirsch, M. W. 1989. Convergent activation dynamics in continuous time networks. *Neural Netw.*, 2(5): 331–349.
- Johnson, R.; and Zhang, T. 2013. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems*, 315–323.
- Korda, N.; and La, P. 2015. On TD (0) with function approximation: Concentration bounds and a centered variant with exponential convergence. In *International conference on machine learning*, 626–634. PMLR.
- Liu, B.; Gemp, I.; Ghavamzadeh, M.; Liu, J.; Mahadevan, S.; and Petrik, M. 2018. Proximal gradient temporal difference learning: Stable reinforcement learning with polynomial sample complexity. *Journal of Artificial Intelligence Research*, 63: 461–494.
- Liu, B.; Liu, J.; Ghavamzadeh, M.; Mahadevan, S.; and Petrik, M. 2015. Finite-sample analysis of proximal gradient TD algorithms. In *Proceedings of the 21st Conference on Uncertainty in Artificial Intelligence*, 504–513.
- Liu, B.; Liu, J.; Ghavamzadeh, M.; Mahadevan, S.; and Petrik, M. 2016. Proximal Gradient Temporal Difference Learning Algorithms. In *Proceedings of the International Joint Conference on Artificial Intelligence*, 4195–4199.
- Ng, A. Y.; Harada, D.; and Russell, S. 1999. Policy invariance under reward transformations: Theory and application to reward shaping. In *Proc. 16th Int. Conf. Mach. Learn.*, 278–287.
- Pan, Y.; White, A.; and White, M. 2017. Accelerated gradient temporal difference learning. In *Proceedings of the 21st AAAI Conference on Artificial Intelligence*, 2464–2470.
- Schulman, J.; Levine, S.; Abbeel, P.; Jordan, M.; and Moritz, P. 2015. Trust region policy optimization. In *International Conference on Machine Learning*, 1889–1897.
- Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; and Klimov, O. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Schwartz, A. 1993. A reinforcement learning method for maximizing undiscounted rewards. In *Proc. 10th Int. Conf. Mach. Learn.*, volume 298, 298–305.
- Sutton, R.; Maei, H.; Precup, D.; Bhatnagar, S.; Silver, D.; Szepesvári, C.; and Wiewiora, E. 2009. Fast gradient-descent methods for temporal-difference learning with linear function approximation. In *Proc. 26th Int. Conf. Mach. Learn.*, 993–1000.
- 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. 2018. *Reinforcement Learning: An Introduction*. The MIT Press, second edition.
- Sutton, R. S.; Maei, H. R.; and Szepesvári, C. 2008. A Convergent $O(n)$ Temporal-difference Algorithm for Off-policy Learning with Linear Function Approximation. In *Advances in Neural Information Processing Systems*, 1609–1616. Cambridge, MA: MIT Press.
- Sutton, R. S.; Mahmood, A. R.; and White, M. 2016. An emphatic approach to the problem of off-policy temporal-difference learning. *The Journal of Machine Learning Research*, 17(1): 2603–2631.
- Tsitsiklis, J. N.; and Van Roy, B. 1997. Analysis of temporal-difference learning with function approximation. In *Advances in Neural Information Processing Systems*, 1075–1081.
- Xu, T.; Wang, Z.; Zhou, Y.; and Liang, Y. 2019. Reanalysis of Variance Reduced Temporal Difference Learning. In *International Conference on Learning Representations*.
- Xu, T.; Wang, Z.; Zhou, Y.; and Liang, Y. 2020. Reanalysis of variance reduced temporal difference learning. *arXiv preprint arXiv:2001.01898*.
- Zhang, S.; and Whiteson, S. 2022. Truncated emphatic temporal difference methods for prediction and control. *The Journal of Machine Learning Research*, 23(1): 6859–6917.