# Eigensubspace of Temporal-Difference Dynamics and How It Improves Value Approximation in Reinforcement Learning

Qiang He[1](✉), Tianyi Zhou[2], Meng Fang[3], and Setareh Maghsudi[1]

[1] University of Tübingen, Tübingen, Germany
{qiang.he, setareh.maghsudi}@uni-tuebingen.de
[2] University of Maryland, College Park, USA
tianyi@umd.edu
[3] University of Liverpool, Liverpool, UK
Meng.Fang@liverpool.ac.uk

**Abstract.** We propose a novel value approximation method, namely "Eigensubspace Regularized Critic (ERC)" for deep reinforcement learning (RL). ERC is motivated by an analysis of the dynamics of Q-value approximation error in the Temporal-Difference (TD) method, which follows a path defined by the 1-eigensubspace of the transition kernel associated with the Markov Decision Process (MDP). It reveals a fundamental property of TD learning that has remained unused in previous deep RL approaches. In ERC, we propose a regularizer that guides the approximation error tending towards the 1-eigensubspace, resulting in a more efficient and stable path of value approximation. Moreover, we theoretically prove the convergence of the ERC method. Besides, theoretical analysis and experiments demonstrate that ERC effectively reduces the variance of value functions. Among 26 tasks in the DMControl benchmark, ERC outperforms state-of-the-art methods for 20. Besides, it shows significant advantages in Q-value approximation and variance reduction. Our code is available at https://sites.google.com/view/erc-ecml23/.

## 1 Introduction

In recent years, deep reinforcement learning (RL), which is built upon the basis of the Markov decision process (MDP), has achieved remarkable success in a wide range of sequential decision-making tasks [29], including board games [27], video games [23, 35], and robotics manipulation [11]. Leveraging the rich structural information in MDP, such as the Markov assumption [29], the stochasticity of transition kernel [1], and low-rank MDP [2, 26, 33, 39], supports designing efficient RL algorithms.

Motivated by the potential benefits of utilizing structural information to design efficient RL algorithms [15, 19–21, 37], we investigate the dynamics of Q value approximation induced by the Temporal-Difference (TD) method. The Q value, commonly used to design DRL algorithms, is analyzed in the context of the matrix form of the Bellman equation. We study continuous learning dynamics. We
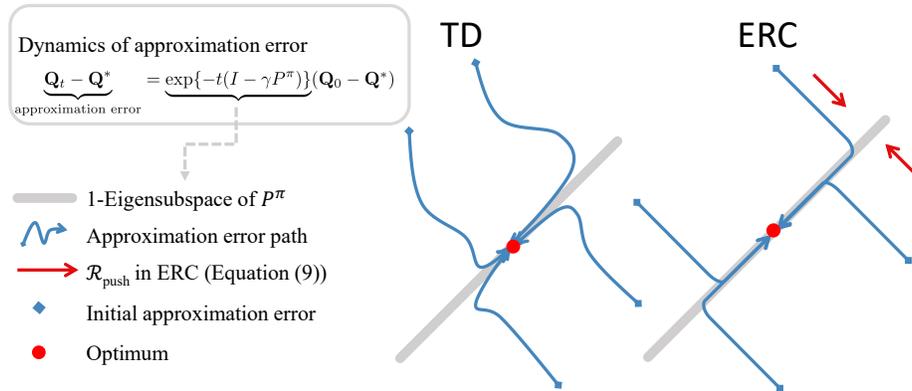
Fig. 1: Value approximation error $(Q - Q^*)$ path for TD and our proposed ERC algorithms. The approximation error of the TD method gradually approaches the 1-eigensubspace (as defined in Remark 3) before ultimately converging to the optimum. This path is referred to as the inherent path. ERC leverages this inherent path by directly pushing the approximation error towards the 1-eigensubspace through $\mathcal{R}_{\text{push}}$ (defined in Equation (9)), resulting in a more efficient and stable learning process.

also examine the crucial role of the transition kernel of MDP in the dynamics of the Bellman equation. We perform an eigenvalue decomposition of that dynamic process, which reveals that for any MDP, the TD method induces an inherent learning path of approximation error in the value approximation process.

Furthermore, as Figure 1 shows, the approximation error of the TD method is optimized towards the 1-eigensubspace before ultimately converging to zero. It is in direct contrast with the optimum that the Monte Carlo (MC) method achieves [29]. Thus, such an inherent path is non-trivial. Such a path, which corresponds to prior knowledge of MDP, has been neglected in designing DRL algorithms [19–21, 37], despite its great potential. Our main idea is utilizing the prior knowledge to improve the efficiency of value approximation by directly guiding the approximation error towards the 1-eigensubspace, leading to a more efficient and stable path. To that end, we propose a novel value approximation method, Eigensubspace Regularized Critic (ERC), as shown in Figure 1. We also establish the convergence of our proposal. Theoretical analysis and experiments demonstrate that ERC effectively reduces the variance of value functions.

To evaluate the effectiveness of our proposed algorithm, ERC, we conduct extensive experiments on the continuous control suite DMControl [32]. Our empirical results demonstrate that ERC performs well in terms of approximation error. Moreover, by examining the ERC variance reduction, we verify its superior performance compared to the algorithms specifically designed for variance control (e.g., TQC [17], REDQ [5]), as consistent with our theoretical analysis. All in all, comparisons show that ERC outperforms the majority of the state-of-the-art

methods (**20** out of **26** tasks) while achieving a similar performance otherwise. Specifically, on average, the ERC algorithm surpasses TQC, REDQ, SAC [11], and TD3 [9] by 13%, 25.6%, 16.6%, and 27.9%, respectively.

The main contributions of this work include i) identifying the existence of an inherent path of the approximation error for TD learning, ii) leveraging the inherent path to introduce a more efficient and stable method, ERC, with a convergence guarantee, and iii) demonstrating, through comparison with state-of-the-art deep RL methods, the superiority of ERC in terms of performance, variance, and approximation error.

## 2  Preliminaries

To formalize RL, one uses an MDP framework consisting of 6-tuple $(\mathcal{S}, \mathcal{A}, R, P, \gamma, \rho_0)$, where $\mathcal{S}$ denotes a state space, $\mathcal{A}$ an action space, $R : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ a reward function, $P : \mathcal{S} \times \mathcal{A} \to p(s)$ a transition kernel, $\gamma \in [0, 1)$ a discount factor, and $\rho_0$ an initial state distribution.

Deep RL focuses on optimizing the policy through return, defined as $R_t = \sum_{i=t}^{T} \gamma^{i-t} r(s_i, a_i)$. The action value (Q) function, $Q^\pi(s, a)$, represents the quality of a specific action, $a$, in a state, $s$, for a given policy $\pi$. Formally, the Q function is defined as

$$Q^\pi(s, a) = \mathbb{E}_{\tau \sim \pi, p}[R_\tau | s_0 = s, a_0 = a], \tag{1}$$

where $\tau$ is a state-action sequence $(s_0, a_0, s_1, a_1, s_2, a_2 \cdots)$ induced by a policy $\pi$ and $P$. The state value (V) function is $V^\pi(s) = \mathbb{E}_{\tau \sim \pi, p}[R_\tau | s_0 = s]$. A four-tuple $(s_t, a_t, r_t, s_{t+1})$ is referred to as a transition. The $Q$ value can be recursively computed by Bellman equation [29]

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s', a'}[Q^\pi(s', a')], \tag{2}$$

where $s' \sim p(\cdot|s, a)$ and $a \sim \pi(\cdot|s)$. The process of using a function approximator (e.g. neural networks) to estimate Q or V values is referred to as value approximation.

**Bellman equation in matrix form.** Let $\mathbf{Q}^\pi$ denote the vector of all Q value with length $|\mathcal{S}| \cdot |\mathcal{A}|$, and $\mathbf{r}$ as vectors of the same length. We overload notation and let P refer to a matrix of dimension $(|\mathcal{S}| \cdot |\mathcal{A}|) \times |\mathcal{S}|$, with entry $P_{s,a,s'}$ equal to $P(s'|s, a)$. We define $P^\pi$ to be the transition matrix on state-action pairs induced by a stationary policy $\pi$

$$P^\pi_{s,a,s',a'} := P(s'|s, a)\pi(a'|s'). \tag{3}$$

Then, it is straightforward to verify

$$\mathbf{Q}^\pi = \mathbf{r} + \gamma P^\pi \mathbf{Q}^\pi, \tag{4}$$

where $P^\pi \in \mathbb{R}^{|\mathcal{S}| \cdot |\mathcal{A}| \times |\mathcal{S}| \cdot |\mathcal{A}|}$. The following famous eigenpair result holds.

*Remark 1 (Eigenpair for stochastic matrix $P^\pi$ [22]).* The spectral radius of $P^\pi$ is 1. The eigenvector corresponding to 1 is $\mathbf{e}$, where $\mathbf{e}$ is a column of all 1's.

## 3    Method

In this section, we start with the dynamics of Q value approximation induced by the TD method. This analysis reveals that in the value approximation, the approximation error has an inherent path. We leverage that path to design our proposed practical algorithm, ERC. We also investigate the convergence property of our method. All proofs are available in the appendix.

### 3.1    An Inherent Path of Value Approximation

Motivated by the potential benefits of using structural information to design novel RL algorithms [14, 19–21, 37], we examine the dynamics of Q value induced by the Bellman equation. Given Equation (4) in matrix form, the true Q function of policy $\pi$ can be directly solved.

*Remark 2 (True Q value of a policy [1]).*  Given a transition matrix induced by a policy $\pi$, as defined in Equation (3), the true Q function of policy $\pi$ can be directly solved by

$$\mathbf{Q}^* = (I - \gamma P^\pi)^{-1}\, \mathbf{r}, \tag{5}$$

where $I$ is the identity matrix and $\mathbf{Q}^*$ is the true Q function of policy $\pi$.

To simplify the notation, let $X = (s, a)$. The one-step temporal difference (TD) continuous learning dynamics follows as

$$\partial_t Q_t(x) = \mathbb{E}_\pi[r_t + \gamma Q_t(x_{t+1})|x_t] - Q_t(x). \tag{6}$$

According to Equation (4), we have the matrix form

$$\partial_t \mathbf{Q}_t = -(I - \gamma P^\pi)\mathbf{Q}_t + \mathbf{r}. \tag{7}$$

Equation (7) is a differential equation that is directly solvable using Remark 2.

**Lemma 1 (Dynamics of approximation error).**  *Consider a continuous sequence $\{Q_t | t \geq 0\}$, satisfy Equation (7) with initial condition $\mathbf{Q}_0$ at time step $t = 0$, then*

$$\mathbf{Q}_t - \mathbf{Q}^* = \exp\{-t(I - \gamma P^\pi)\}(\mathbf{Q}_0 - \mathbf{Q}^*). \tag{8}$$

From Lemma 1, $\mathbf{Q}_t$ converges to $\mathbf{Q}^*$, as $t \to \infty$. The approximation error, $\mathbf{Q}_t - \mathbf{Q}^*$, appears in Equation (8), which reveals the dynamics of approximation error is related to the structure of transition kernel and policy $\pi$. Moreover, there is rich structural information in $P^\pi$, which inspires us to consider Equation (8) in a more fine-grained way. To better understand the approximation error, following Ghosh and Bellemare [10], Lyle et al. [20, 21], we make the following assumption.

**Assumption 1** $P^\pi$ *is a real-diagonalizable matrix with a strictly decreasing eigenvalue sequence* $|\lambda_1|, |\lambda_2|, \cdots, |\lambda_{|\mathcal{S}| \cdot |\mathcal{A}|}|$, *and the corresponding eigenvector* $H_1, H_2, \cdots, H_{|\mathcal{S}| \cdot |\mathcal{A}|}$.

Remark 1 shows that $\lambda_1 = 1$ because the $P^\pi$ is a stochastic matrix, and the eigenvector corresponding to 1 is $\mathbf{e}$. That inspires us to perform an eigenvalue decomposition for Equation (8).

**Theorem 1.** *If Assumption 1 holds, we have* $\mathbf{Q}_t - \mathbf{Q}^* = \alpha_1 \exp\{t(\gamma\lambda_1 - 1)\}H_1 + \sum_{i=2}^{|\mathcal{S}|\cdot|\mathcal{A}|} \alpha_i \exp\{t(\gamma\lambda_i - 1)\}H_i = \alpha_1 \exp\{t(\gamma\lambda_1 - 1)\}H_1 + o\Big(\alpha_1 \exp\{t(\gamma\lambda_1 - 1)\}\Big),$ *where* $\alpha_i$ *is a constant.*

Theorem 1 states that the approximation error of the TD method can be decomposed into two components. One of the components is primarily influenced by the 1-eigensubspace. Thus, the approximation error of the value function in the TD method follows a path induced by the 1-eigensubspace, which is a result of the stochasticity inherent in MDP.

*Remark 3 (An inherent path of TD method).* Given an MDP consisting of 6-tuple $(\mathcal{S}, \mathcal{A}, R, P, \gamma, \rho_0)$, policy $\pi$, and a Banach space $(\mathcal{Q}, \|\cdot\|)$, there exists an inherent path that the Bellman approximation error, i.e., $\mathbf{Q}_t - \mathbf{Q}^*$, starts at the initial point and then approaches the 1-eigensubspace before ultimately converging to zero. The 1-eigensubspace, which is induced by $P^\pi$, is defined as $\{c\,\mathbf{e}\}$, where $c \in \mathbb{R}$ and $\mathbf{e}$ is a column of all ones.

**What does a theoretically inherent path really look like in practice?** The aforementioned content discusses the inherent path of approximation error. Does this occur in the practical scene? To empirically show the path of the approximation error, we visualize the path of approximation error given a fixed policy. The results are given in Figure 2, where we perform experiments on FrozenLake-v1 environment since the true Q value $Q^*$ of this environment can be evaluated by the Monte Carlo Method. The approximation error of the TD method is optimized toward the 1-eigensubspace before ultimately converging to zero rather than directly toward the optimum. The inherent path in Figure 2
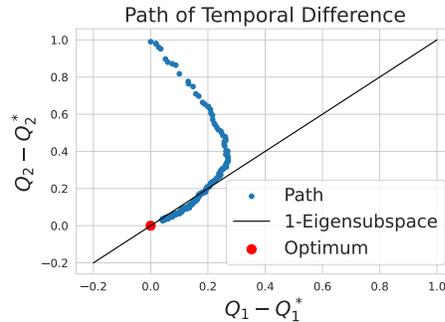


Fig. 2: Path of TD method. There exists an inherent path that approximation error approaches 1-eigensubspace before converging to zero. The empirical fact is consistent with our theoretical analysis in Theorem 1.
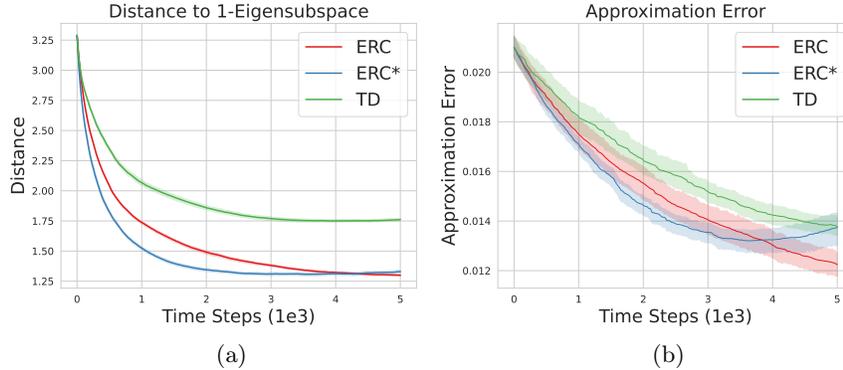
is consistent with Remark 3. Thus, such a path is non-trivial and motivates us to improve the value approximation by guiding the approximation error to 1-eigensubspace, resulting in an efficient and robust path.

Fig. 3: Value function approximation process for various methods on FrozenLake-v1 environment. (a) illustrates the distance between the approximation error and the 1-eigensubspace for various methods, where ERC* denotes ERC utilizing an oracle true Q value, $Q^*$, to push the approximation error towards the 1-eigensubspace. The results demonstrate that the ERC method is closer to the 1-eigensubspace at the same time compared to the TD method. (b) represents the absolute approximation error for various algorithms. The result illustrates that the ERC method has a smaller approximation error at the same time than the TD method. For both metrics, ERC is observed to outperform or be at least as good as ERC* in later stages. The shaded area represents a standard deviation over ten trials.

### 3.2   Using Eigensubspace Regularization to Improve Value Approximation

The inherent path, which can be viewed as prior knowledge of MDP, has not been utilized to design DRL algorithms in previous work [19–21, 37]. We improve the value approximation by directly guiding the approximation error towards the 1-eigensubspace, leading to a more efficient and stable path.

The true Q value, $Q^*$, is always unknown when designing DRL algorithms. However, we can use a target Q as an approximate substitute for $Q^*$ for two reasons. Firstly, from the perspective of value function optimization, the objective of the Bellman equation optimization is to make the learned Q-value as close as possible to $Q^*$. That is achievable by minimizing the distance between the learned Q and the target Q. Instead of using $Q^*$ directly in learning, the target Q is used to approximate $Q^*$ through a bootstrap approach. Similarly, we use target Q in our ERC algorithm design. Secondly, in our experiments, we find that using target Q to replace $Q^*$ produces an effect that approximates the effect produced by using $Q^*$, as illustrated in Figure 3. We calculate the distance of the approximation error to the 1-eigensubspace and approximation error during the optimization of the value function, and we can see that: i) the ERC algorithm using target Q instead of $Q^*$ allows the approximation error to reach the 1-eigensubspace

faster than the TD algorithm and the effect can be approximated to that of the ERC using $Q^*$ (ERC* algorithm). And ii) for the approximation error, the ERC algorithm obtains a relatively small approximation error compared to the TD method. The approximation error of ERC is smaller than that of ERC* in the later optimization stages. Therefore, it is reasonable to replace $Q^*$ with the target Q in the practical design of the algorithm. Using target Q-value is a feasible solution that allows us to achieve results similar or better to using $Q^*$, and it has the advantage of being more easily implemented in practice. Thus, the Bellman error is pushed to 1-eigensubspace in the ERC algorithm.

To push the Bellman error toward the 1-eigensubspace, it is essential to project the error onto that subspace. Therefore, one must determine the projected point in the 1-eigensubspace to the Bellman error.

**Lemma 2.** *Consider a Banach space $(\mathfrak{B}, \|\cdot\|)$ of dimension N, and let the N-dimensional Bellman error at timestep t, represented by $\mathbf{B}^t$, have coordinates $(B_1, B_2, \cdots, B_N)$ in $(\mathfrak{B}, \|\cdot\|)$. Within this Banach space, the projected point in the 1-eigensubspace, which is closest to $B^t$, is $Z^t$ whose coordinates are $(z^t, z^t, \cdots, z^t)$, where $z^t = \frac{1}{N}\sum_{j=1}^{N} B_i^t$.*

The Bellman error can be pushed towards 1-eigensubspace at each timestep with the help of Lemma 2. To accomplish this, we minimize the regularization term

$$\mathcal{R}_{\text{push}}(\theta) = \frac{1}{N}\sum_{i=1}^{N}\|B_i - Z\|_2^2, \qquad (9)$$

where $B_i$ represents the Bellman error at the $i$-dimension and $Z = \frac{1}{N}\sum_{j=1}^{N} B_i(\theta)$. By combining Equation (9) with policy evaluation loss $\mathcal{L}_{\mathbf{PE}}$, the $\underline{E}$igensubspace $\underline{R}$egularized $\underline{C}$ritic (ERC) algorithm is defined as

$$\mathcal{L}_{\text{ERC}}(\theta) = \mathcal{L}_{\text{PE}}(\theta) + \beta\mathcal{R}_{\text{push}}(\theta), \qquad (10)$$

where $\mathcal{L}_{\text{PE}}(\theta)$ is a policy evaluation phase loss such as

$$\mathcal{L}_{\text{PE}}(\theta) = \left[Q(s,a;\theta) - \left(r(s,a) + \gamma\mathbb{E}_{s',a'}\left[Q(s',a';\theta')\right]\right)\right]^2,$$

and $\beta$ is a hyper-parameter that controls the degree to which the Bellman error is pushed toward the 1-eigensubspace. ERC enhances the value approximation by pushing the Bellman error to 1-eigensubspace, leading to a more stable and efficient value approximation path. To evaluate the effectiveness of ERC, a case study is conducted on the FrozenLake environment. The distance between the approximation error and the 1-eigensubspace, as well as the absolute approximation error, are used as metrics to compare the performance of ERC with that of the TD method and that of ERC* (ERC utilizing oracle $Q^*$). The results in Figure 3 demonstrate that ERC is superior to the TD method, and is either superior to or at least as effective as ERC*.

The theoretical benefit of the ERC method can be understood by examining Equation (9). Minimizing $\mathcal{R}_{\text{push}}$ explicitly reduces the variance of the Bellman

error. This reduction in variance leads to two benefits: minimizing the variance of the Q-value and minimizing the variance of the target Q-value. That is observable by rewriting Equation (9) as

$$
\begin{aligned}
\mathcal{R}_{\text{push}}(\theta) =& \mathbb{E}\Big( (Q - \mathcal{B}Q) - \mathbb{E}[Q - \mathcal{B}Q] \Big)^2 \\
=& \underbrace{\mathbb{E}[(Q - \mathbb{E}[Q])^2]}_{\text{variance of } Q} + \underbrace{\mathbb{E}[(\mathcal{B}Q - \mathbb{E}\mathcal{B}Q)^2]}_{\text{variance of } \mathcal{B}Q} - \\
& 2\underbrace{\mathbb{E}(Q - \mathbb{E}[Q])(\mathcal{B}Q - \mathbb{E}\mathcal{B}Q)}_{\text{covariance between } Q \text{ and } \mathcal{B}Q},
\end{aligned}
\tag{11}
$$

where $\mathcal{B}$ is a bellman backup operator and $\mathcal{B}Q$ is a target Q. These facts highlight the benefits of the ERC algorithm, as it leads to a more stable and efficient Q value approximation.

---

**Algorithm 1:** ERC (based on SAC [11])

---

**Initialize** actor network $\pi$, and critic network $Q$ with random parameters;
**Initialize** target networks and replay buffer $\mathcal{D}$;
**Initialize** $\beta$, total steps $T$, and $t = 0$;
Reset the environment and receive the initial state $s$;
**while** $t < T$ **do**
　Select action w.r.t. its policy $\pi$ and receive reward $r$, new state $s'$;
　Store transition tuple $(s, a, r, s')$ to $\mathcal{D}$;
　Sample $N$ transitions $(s, a, r, s')$ from $\mathcal{D}$;
　Compute $\hat{\mathcal{R}}_{\text{push}}$ by Equations (9) and (16);
　Update critic by minimizing Equation (13);
　Update actor by minimizing Equation (14);
　Update $\alpha$ by minimizing Equation (15);
　Update target networks;
　$t \leftarrow t + 1, s \leftarrow s'$;
**end**

---

### 3.3   Theoretical Analysis

We are also interested in examining the convergence property of ERC. To do this, we first obtain the new Q value after one step of updating in tabular form.

**Lemma 3.** *Given the ERC update rules in Equation* (10), *ERC updates the value function in tabular form in the following way*

$$
Q_{t+1} = (1 - \alpha_t(1 + \beta))Q_t + \alpha_t(1 + \beta)\mathcal{B}Q_t - \alpha_t\beta C_t,
\tag{12}
$$

*where* $\mathcal{B}Q_t(s_t, a_t) = r_t + \gamma\mathbb{E}_{s_{t+1}, a_{t+1}}Q_t(s_{t+1}, a_{t+1})$, $C_t = 2\mathbb{NG}(\mathbb{E}[\mathcal{B}Q_t - Q_t])$, *and* $\mathbb{NG}$ *means stopping gradient.*

To establish the convergence guarantee of ERC, we use an auxiliary lemma from stochastic approximation [28] and the results in Haarnoja et al. [11]. Theorem 2 states the convergence guarantee formally.

**Theorem 2 (Convergence of 1-Eigensubspace Regularized Value Approximation).** *Consider the Bellman backup operator $\mathcal{B}$ and a mapping $Q : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$, and $Q^k$ is updated with Equation (12). Then the sequence $\{Q^k\}_{k=0}^{\infty}$ will converge to 1-eigensubspace regularized optimal $Q$ value of $\pi$ as $k \to \infty$.*

### 3.4 Practical Algorithm

The proposed ERC algorithm, which utilizes the structural information in MDP to improve value approximation, can be formulated as a practical algorithm. We combine the ERC with Soft Actor Critic (SAC) algorithm [11]. For the value approximation, ERC optimizes

$$J_{\text{ERC}}^Q(\theta) = \mathbb{E}_{(s,a)\sim\mathcal{D}}\left[\frac{1}{2}\Big(Q(s,a;\theta) - \big(r(s,a) + \gamma\mathbb{E}_{s'\sim P}[V(s';\theta')]\big)\Big)^2\right] + \beta\mathcal{R}_{\text{push}}, \tag{13}$$

where $V(s;\theta') = \mathbb{E}_{a\sim\pi(\cdot|s;\phi)}[Q(s,a;\theta') - \alpha\log\pi(a|s;\phi)]$. For policy improvement, ERC optimizes

$$J_{\text{ERC}}^\pi(\phi) = \mathbb{E}_{s\sim\mathcal{D}}[\mathbb{E}_{a\sim\pi(\cdot|s,\phi)}[\alpha\log(\pi(a\mid s;\phi)) - Q(s,a;\theta)]]. \tag{14}$$

Besides, we also use the automated entropy trick. The temperature $\alpha$ is learned by minimizing

$$J_{\text{ERC}}^\alpha(\alpha) = \mathbb{E}_{a\sim\pi^*}[-\alpha\log\pi^*(a|s;\alpha,\phi) - \alpha\mathcal{H}], \tag{15}$$

where $\mathcal{H}$ is a pre-selected target entropy. To help better stabilize the value approximation of ERC, we design a truncation mechanism for $\mathcal{R}_{\text{push}}$

$$\hat{\mathcal{R}}_{\text{push}} = \max\Big\{\min\Big\{\beta\mathcal{R}_{\text{push}}, \mathcal{R}_{\text{max}}\Big\}, \mathcal{R}_{\text{min}}\Big\}. \tag{16}$$

The practical algorithm is summarized in Algorithm 1.

## 4    Experiments

In this section, we thoroughly evaluate the performance of ERC by comparing it to a few baseline methods. Furthermore, we examine the value approximation error and the variance of the value function to gain a deeper understanding of ERC. Additionally, we analyze the individual contributions of each component of ERC to gain insight into its effectiveness.
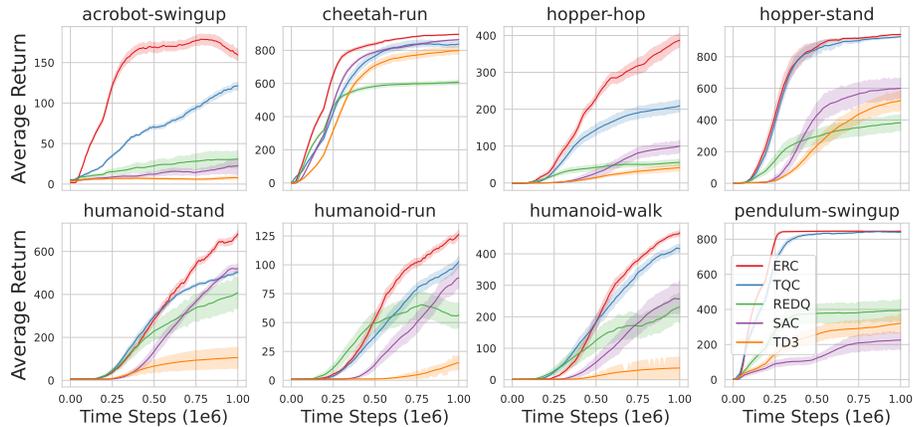
Fig. 4: Performance curves for OpenAI gym continuous control tasks on Deep-Mind Control suite. The proposed algorithm, ERC, is observed to significantly outperform the other tested algorithms. The shaded region represents half of the standard deviation of the average evaluation over 10 seeds. The curves are smoothed with a moving average window of size ten.

### 4.1   Evaluation Setting

**Baselines.** We conduct a comparative study of the proposed ERC algorithm with several well-established baselines in the literature. Specifically, we select TD3 [9] and SAC as our primary baselines as they are commonly used and perform well in various tasks. Additionally, we compare our method with REDQ [5] and TQC [17], which employ different techniques to improve value approximation and reduce the variance of the value function. To ensure a fair comparison, we use the authors' implementation of TD3 and REDQ available on Github, and the public implementation of SAC provided in PyTorch [38], which is also the basis of our ERC implementation. Besides, we use the implementation of TQC available in the stable baselines 3 library and use the default hyper-parameters as suggested by the authors.

   **Environments.** The experimental suite is the state-based DMControl suite [32], which is for physics-based simulation, utilizing the MuJoCo physics engine [31]. We chose the DMControl suite as it offers a diverse range of environments to benchmark the capabilities of RL algorithms. We facilitate the interactions between the algorithm and environment using Gym [4]. We evaluate each algorithm over one million timesteps and obtain the average return of the algorithm every 10k timesteps over ten episodes.

   **Setup.** The magnitude of the regularization effectiveness of ERC is controlled by a hyper-parameter, $\beta$, which is $5e-3$. Additionally, $\mathcal{R}_{\max}$ and $\mathcal{R}_{\min}$ are $1e-2$ and 0, respectively, for all experiments. The remaining hyper-parameters are consistent with the suggestions provided by Haarnoja et al. [11]. To ensure the validity and reproducibility of the experiments, unless otherwise specified, we eval-

Table 1: Average Return after 1M timesteps of training on DMC. ERC demonstrates state-of-the-art performance on the majority (**20** out of **26**) tasks. If not, ERC is still observed to be comparable in performance. Additionally, ERC outperforms its backbone algorithm, SAC, on **all** tasks by a large margin. Specifically, the ERC algorithm outperforms TQC, REDQ, SAC, and TD3 by 13%, 25.6%, 16.6%, and 27.9%, respectively. The best score is marked with  colorbox. $\pm$ corresponds to a standard deviation over ten trials.

| Domain | Task | ERC | TQC | REDQ | SAC | TD3 |
|---|---|---|---|---|---|---|
| Acrobot | Swingup | $151.0 \pm 36.8$ | $136.3 \pm 51.9$ | $31.1 \pm 41.1$ | $26.9 \pm 47.7$ | $5.3 \pm 4.7$ |
| BallInCup | Catch | $979.7 \pm 1.3$ | $981.6 \pm 2.3$ | $978.8 \pm 3.7$ | $980.3 \pm 3.4$ | $978.9 \pm 3.6$ |
| Cartpole | Balance | $998.8 \pm 1.2$ | $989.2 \pm 25.8$ | $984.0 \pm 6.0$ | $997.7 \pm 1.5$ | $997.9 \pm 2.1$ |
| Cartpole | BalanceSparse | $998.6 \pm 4.5$ | $899.9 \pm 268.6$ | $872.1 \pm 262.7$ | $997.6 \pm 5.7$ | $1000.0 \pm 0.0$ |
| Cartpole | Swingup | $867.8 \pm 4.7$ | $874.3 \pm 5.8$ | $828.1 \pm 17.2$ | $865.1 \pm 1.6$ | $867.2 \pm 7.5$ |
| Cartpole | SwingupSparse | $544.5 \pm 356.9$ | $797.6 \pm 32.1$ | $385.5 \pm 374.7$ | $234.4 \pm 358.4$ | $157.5 \pm 314.9$ |
| Cheetah | Run | $903.3 \pm 5.9$ | $853.8 \pm 80.0$ | $614.2 \pm 58.2$ | $873.4 \pm 21.5$ | $811.3 \pm 102.2$ |
| Finger | Spin | $988.1 \pm 0.6$ | $982.0 \pm 9.1$ | $940.1 \pm 33.5$ | $966.3 \pm 27.1$ | $947.6 \pm 52.1$ |
| Finger | TurnEasy | $981.1 \pm 5.4$ | $247.4 \pm 133.6$ | $962.6 \pm 34.3$ | $920.0 \pm 91.8$ | $856.5 \pm 109.3$ |
| Finger | TurnHard | $964.8 \pm 27.5$ | $299.2 \pm 266.6$ | $927.3 \pm 99.8$ | $874.1 \pm 100.1$ | $690.2 \pm 167.6$ |
| Fish | Upright | $936.0 \pm 12.1$ | $917.1 \pm 25.6$ | $799.6 \pm 113.8$ | $898.6 \pm 50.4$ | $873.6 \pm 66.7$ |
| Fish | Swim | $496.8 \pm 61.6$ | $526.6 \pm 113.5$ | $159.3 \pm 100.1$ | $342.4 \pm 134.5$ | $251.3 \pm 107.7$ |
| Hopper | Stand | $943.9 \pm 8.9$ | $941.6 \pm 11.4$ | $393.5 \pm 225.8$ | $597.8 \pm 308.8$ | $538.7 \pm 256.2$ |
| Hopper | Hop | $405.0 \pm 91.1$ | $221.8 \pm 68.8$ | $56.8 \pm 36.2$ | $117.4 \pm 82.2$ | $47.8 \pm 46.2$ |
| Humanoid | Stand | $804.5 \pm 39.1$ | $494.9 \pm 145.5$ | $407.2 \pm 336.8$ | $549.6 \pm 201.0$ | $110.6 \pm 206.8$ |
| Humanoid | Walk | $507.0 \pm 37.0$ | $376.4 \pm 182.5$ | $245.0 \pm 222.6$ | $248.4 \pm 220.8$ | $39.3 \pm 101.9$ |
| Humanoid | Run | $145.9 \pm 10.1$ | $115.6 \pm 18.6$ | $70.8 \pm 57.0$ | $83.4 \pm 56.0$ | $18.1 \pm 33.9$ |
| Pendulum | Swingup | $846.6 \pm 14.1$ | $834.0 \pm 30.1$ | $382.6 \pm 297.0$ | $226.2 \pm 228.9$ | $338.0 \pm 232.0$ |
| PointMass | Easy | $882.3 \pm 18.2$ | $793.7 \pm 147.6$ | $880.9 \pm 16.7$ | $889.9 \pm 33.1$ | $838.8 \pm 158.5$ |
| Reacher | Easy | $986.9 \pm 2.3$ | $964.5 \pm 39.5$ | $970.9 \pm 24.4$ | $983.5 \pm 4.2$ | $983.4 \pm 3.7$ |
| Reacher | Hard | $981.8 \pm 1.7$ | $971.8 \pm 5.2$ | $964.1 \pm 24.0$ | $958.6 \pm 40.9$ | $938.2 \pm 63.0$ |
| Swimmer | Swimmer6 | $422.0 \pm 133.2$ | $356.4 \pm 107.5$ | $215.8 \pm 119.0$ | $359.3 \pm 130.9$ | $289.2 \pm 133.6$ |
| Swimmer | Swimmer15 | $295.8 \pm 113.7$ | $222.8 \pm 128.6$ | $178.6 \pm 116.6$ | $264.6 \pm 136.9$ | $236.7 \pm 150.1$ |
| Walker | Stand | $989.6 \pm 1.7$ | $986.3 \pm 4.5$ | $974.0 \pm 12.6$ | $986.8 \pm 2.7$ | $983.4 \pm 4.2$ |
| Walker | Walk | $974.9 \pm 1.6$ | $971.9 \pm 4.8$ | $957.3 \pm 10.6$ | $973.5 \pm 4.4$ | $966.3 \pm 10.4$ |
| Walker | Run | $805.1 \pm 19.4$ | $770.5 \pm 31.0$ | $590.9 \pm 51.6$ | $773.0 \pm 32.9$ | $712.1 \pm 65.6$ |
| Average | Scores | 761.6 | 674.1 | 606.6 | 653.4 | 595.3 |

uate each tested algorithm over ten fixed random seeds. For more implementation details, please refer to the appendix.

## 4.2 Performance Evaluation

Figure 4 shows the learning curves from scratch. Table 1 depicts the average performance after 1M timesteps training. The results demonstrate the following points: i) ERC outperforms the other tested algorithms in a majority (**20** out of **26**) of the environments. Specifically, it outperforms TQC, REDQ, SAC, and TD3 by 13%, 25.6%, 16.6%, and 27.9%, respectively; ii) ERC substantially improves upon its skeleton algorithm, SAC, by the addition of the $\mathcal{R}_{\text{push}}$ regularization term; iii) Additionally, although ERC does not employ complex techniques to eliminate estimation bias, it still substantially outperforms TQC and REDQ

in most environments. Note that both REDQ and TQC leverage an ensemble mechanism to obtain a more accurate and unbiased Q estimation. ERC does not leverage any ensemble critic, distributional value functions, or high UTD ratio, yet it still outperforms them. The results above highlight the potential of utilizing structural information of MDP to improve DRL.

### 4.3   Variance and Approximation Error

We select four distinct environments 'acrobot-swing', 'humanoid-stand', 'finger-turn_easy', and 'fish-swim'. This selection is made because ERC has been demonstrated to perform well in the first two environments, while its performance is not as strong in the latter two. Thus, by evaluating ERC across these four environments, a comprehensive assessment of ERC can be obtained.
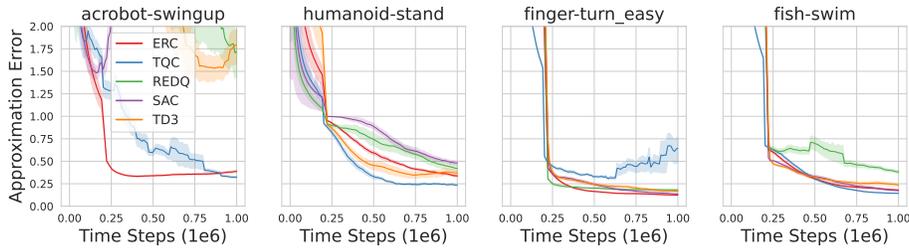


Fig. 5: Approximation error curves. The results demonstrate that the approximation error of ERC is empirically minimal when compared to other algorithms (such as TQC and REDQ) that are specifically designed to obtain accurate unbiased value estimation.
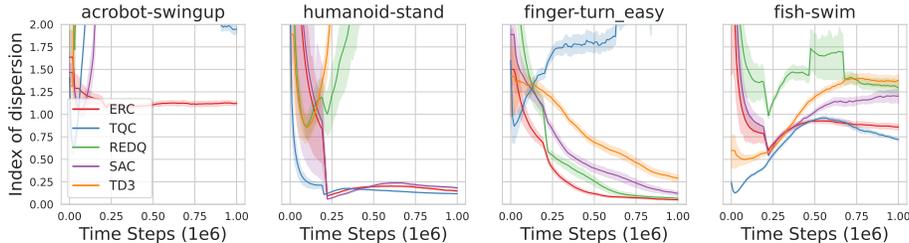


Fig. 6: Index of dispersion curves. The results demonstrate that ERC effectively controls the variance. Furthermore, the variance of ERC is observed to be minimal on the selected tasks, even in situations where the performance of ERC is not as good as TQC and REDQ.

**Approximation Error.** ERC aims to push the Bellman error to 1-eigensubspace to obtain an efficient and stable value approximation as discussed in Section 3.2. Thus we study the approximation error of ERC on the DMControl suite. We obtain the true value $Q^*$ using the Monte Carlo method and compare it to the estimated Q. We normalize the approximation error by the estimation value to

Table 2: Average Return after 1M timesteps training on DMControl suite. $\beta = 5e-3$-trunc means the truncation mechanism is used in the evaluation.

| 1M Steps Scores | $\beta$ =1e-4 | $\beta$ =5e-4 | $\beta$ =1e-3 | $\beta$ =5e-3 | $\beta$ =5e-3-trunc | $\beta$ =1e-2 | $\beta$ =5e-2 | SAC |
|---|---|---|---|---|---|---|---|---|
| Acrobot, Swingup | $151.2 \pm 70.0$ | $152.7 \pm 43.4$ | $91.4 \pm 54.3$ | $151.9 \pm 69.3$ | $151.0 \pm 36.8$ | $129.9 \pm 58.7$ | $114.7 \pm 31.4$ | $26.9 \pm 47.7$ |
| Humanoid, Stand | $418.2 \pm 255.3$ | $755.8 \pm 88.3$ | $692.1 \pm 200.2$ | $699.0 \pm 129.3$ | $804.5 \pm 39.1$ | $742.4 \pm 130.8$ | $550.0 \pm 217.2$ | $549.6 \pm 201.0$ |
| Finger, TurnEasy | $940.4 \pm 49.5$ | $919.8 \pm 80.5$ | $959.2 \pm 37.3$ | $979.4 \pm 5.5$ | $981.1 \pm 5.4$ | $979.9 \pm 5.3$ | $959.7 \pm 43.7$ | $920.0 \pm 91.8$ |
| Fish, Swim | $448.5 \pm 102.8$ | $436.9 \pm 71.4$ | $468.3 \pm 130.5$ | $390.0 \pm 65.8$ | $496.8 \pm 61.6$ | $414.4 \pm 55.0$ | $434.9 \pm 61.4$ | $342.4 \pm 134.5$ |

eliminate the differences in scale, and the absolute value of the approximation error is used in each sample to eliminate inaccuracies caused by offsetting positive and negative errors. The results, presented in Figure 5, indicate that ERC exhibits minimal approximation error compared to other methods that obtain accurate and unbiased value estimates.

**Variance Reduction.** Above analysis shows that ERC reduces the value function variance. To investigate this claim in Section 3.2, we empirically analyze the variance of ERC. To eliminate the effect of the size of the value function on the variance, we use *index of dispersion* [16], which can be considered as a variance normalized by its mean, to evaluate the variance of algorithms. Our results, presented in Figure 6, indicate that i) the variance of the value function of ERC is the lowest or equally low as other methods (TQC, and REDQ) specifically designed for variance reduction. And ii) ERC effectively controls variance even in environments where its performance is not as strong as other methods such as TQC and REDQ. These findings provide strong evidence for the theoretical claim that ERC reduces variance of the value function.

### 4.4 Ablation

The component introduced by ERC, as outlined in Equation (10), is sole $\mathcal{R}_{\text{push}}$, in which $\beta$ regulates the efficiency of pushing the Bellman error towards the 1-eigensubspace. Therefore, investigating the selection of $\beta$ can aid in comprehending the impact of hyper-parameter on ERC. We vary $\beta$ and eliminate the truncation mechanism outlined in Equation (16). Experiments are conducted on the same four environments as described in Section 4.3. The results, presented in Table 2, indicate the following: i) the value of $\beta$ influences the empirical performance of the ERC, yet ERC with various $\beta$ values consistently outperforms its skeleton algorithm, SAC, in most settings. That demonstrates the effectiveness of our proposed value approximation method; ii) Truncating parameters can improve the performance of ERC on specific tasks (Humanoid-Stand, Fish-Swim); iii) Carefully tuning both the hyper-parameter $\beta$ and the truncation mechanism ensures optimal performance.

## 5   Related Work

### 5.1   Value Function Approximation

Mnih et al. [23] utilizes neural networks (NN) to approximate value function by TD learning [30, 36]. Combined with NN, several derivative value approximation

methods [3, 6, 7, 23] exist. Some   [3, 6, 7, 17] use a distributional view to obtain a better value function. Some reduce bias and variance of value function approximation [5, 9, 12, 13, 18, 34]. Others [8, 11, 24] develop more robust value function approximation methods. ERC utilizes structural information from MDP to improve value approximation, which distinguishes it from previous work.

### 5.2   Using Structural Information of MDPs

Recent work [2, 26, 33, 39] focus on the property of low-rank MDP. They study learning an optimal policy assuming a low-rank MDP given the optimal representation. Nevertheless, such a setting is not practical as the optimal representation is computationally intractable. Ren et al. [25] proposed a practical algorithm under low-rank MDP assumption, which, different from ERC, involves estimating the dynamics of MDP. Some other work [1, 19, 21] discuss the architecture of MDP from a matrix decomposition perspective. Lyle et al. [21] establishes a connection between the spectral decomposition of the transition operator and the representations of $V$ function induced by a variety of auxiliary tasks. ERC differs from Lyle et al. [21] in the following aspects: i) We analyze the dynamics of the Q approximation error, which is more commonly studied in DRL literature; ii) We consider the learning process of $Q$ function, whereas Lyle et al. [21] considers the dynamics of representations.

## 6   Conclusion

In this work, we examine the eigensubspace of the TD dynamics and its potential use in improving value approximation in DRL. We begin by analyzing the matrix form of the Bellman equation and subsequently derive the dynamics of the approximation error through the solution of a differential equation. This solution depends on the transition kernel of the MDP, which motivates us to perform eigenvalue decomposition, resulting in the inherent path of value approximation in TD. To the best of our knowledge, this inherent path has not been leveraged to design DRL algorithms in previous work. Our insight is to improve value approximation by directing the approximation error towards the 1-eigensubspace, resulting in a more efficient and stable path. Thus, we propose the ERC algorithm with a theoretical convergence guarantee. Theoretical analysis and experiments demonstrate ERC results in a variance reduction, which validates our insight. Extensive experiments on the DMControl suite demonstrate that ERC outperforms state-of-the-art algorithms in the majority of tasks. The limitation is that ERC is evaluated on the DMControl suite. Verifying the effectiveness of ERC on other suites is left for future work. Our contributions represent a significant step forward in leveraging the rich inherent structure of MDP to improve value approximation and ultimately enhance performance in RL.

## Ethical Statement

This work investigates the eigensubspace of the Temporal-Difference algorithm of reinforcement learning and how it improves DRL. The experimental part only uses a physical simulation engine DMControl. Therefore this work does not involve any personal information. This work may be used in areas such as autonomous vehicles, games, and robot control.

# Bibliography

[1] Agarwal, A., Jiang, N., Kakade, S.M., Sun, W.: Reinforcement learning: Theory and algorithms. CS Dept., UW Seattle, Seattle, WA, USA, Tech. Rep pp. 10–4 (2019)

[2] Agarwal, A., Kakade, S.M., Krishnamurthy, A., Sun, W.: FLAMBE: structural complexity and representation learning of low rank mdps. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (eds.) Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual (2020)

[3] Bellemare, M.G., Dabney, W., Munos, R.: A distributional perspective on reinforcement learning. In: Precup, D., Teh, Y.W. (eds.) Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017, Proceedings of Machine Learning Research, vol. 70, pp. 449–458, PMLR (2017)

[4] Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., Zaremba, W.: Openai gym (2016)

[5] Chen, X., Wang, C., Zhou, Z., Ross, K.W.: Randomized ensembled double q-learning: Learning fast without a model. In: International Conference on Learning Representations (2020)

[6] Dabney, W., Ostrovski, G., Silver, D., Munos, R.: Implicit quantile networks for distributional reinforcement learning. In: International conference on machine learning, pp. 1096–1105, PMLR (2018)

[7] Dabney, W., Rowland, M., Bellemare, M., Munos, R.: Distributional reinforcement learning with quantile regression. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol. 32 (2018)

[8] Espeholt, L., Soyer, H., Munos, R., Simonyan, K., Mnih, V., Ward, T., Doron, Y., Firoiu, V., Harley, T., Dunning, I., Legg, S., Kavukcuoglu, K.: IMPALA: scalable distributed deep-rl with importance weighted actor-learner architectures. In: Dy, J.G., Krause, A. (eds.) Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018, Proceedings of Machine Learning Research, vol. 80, pp. 1406–1415, PMLR (2018)

[9] Fujimoto, S., van Hoof, H., Meger, D.: Addressing function approximation error in actor-critic methods. In: Dy, J.G., Krause, A. (eds.) Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018, Proceedings of Machine Learning Research, vol. 80, pp. 1582–1591, PMLR (2018)

[10] Ghosh, D., Bellemare, M.G.: Representations for stable off-policy reinforcement learning. In: International Conference on Machine Learning, pp. 3556–3565, PMLR (2020)

[11] Haarnoja, T., Zhou, A., Abbeel, P., Levine, S.: Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In:

Dy, J.G., Krause, A. (eds.) Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018, Proceedings of Machine Learning Research, vol. 80, pp. 1856–1865, PMLR (2018)

[12] Hasselt, H.: Double q-learning. Advances in neural information processing systems **23**, 2613–2621 (2010)

[13] He, Q., Hou, X.: Wd3: Taming the estimation bias in deep reinforcement learning. In: 2020 IEEE 32nd International Conference on Tools with Artificial Intelligence (ICTAI), pp. 391–398, IEEE (2020)

[14] He, Q., Su, H., Zhang, J., Hou, X.: Representation gap in deep reinforcement learning. CoRR **abs/2205.14557** (2022)

[15] He, Q., Su, H., Zhang, J., Hou, X.: Frustratingly easy regularization on representation can boost deep reinforcement learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 20215–20225 (2023)

[16] Hoel, P.G.: On indices of dispersion. The Annals of Mathematical Statistics **14**(2), 155–162 (1943)

[17] Kuznetsov, A., Shvechikov, P., Grishin, A., Vetrov, D.: Controlling overestimation bias with truncated mixture of continuous distributional quantile critics. In: Proceedings of the 37th International Conference on Machine Learning, pp. 5556–5566 (2020)

[18] Lan, Q., Pan, Y., Fyshe, A., White, M.: Maxmin q-learning: Controlling the estimation bias of q-learning. In: International Conference on Learning Representations (2019)

[19] Lyle, C., Rowland, M., Dabney, W.: Understanding and preventing capacity loss in reinforcement learning. In: The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022, OpenReview.net (2022)

[20] Lyle, C., Rowland, M., Dabney, W., Kwiatkowska, M., Gal, Y.: Learning dynamics and generalization in deep reinforcement learning. In: Chaudhuri, K., Jegelka, S., Song, L., Szepesvári, C., Niu, G., Sabato, S. (eds.) International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA, Proceedings of Machine Learning Research, vol. 162, pp. 14560–14581, PMLR (2022)

[21] Lyle, C., Rowland, M., Ostrovski, G., Dabney, W.: On the effect of auxiliary tasks on representation dynamics. In: International Conference on Artificial Intelligence and Statistics, pp. 1–9, PMLR (2021)

[22] Meyer, C.D.: Matrix analysis and applied linear algebra, vol. 71. Siam (2000)

[23] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., et al.: Human-level control through deep reinforcement learning. nature **518**(7540), 529–533 (2015)

[24] Munos, R., Stepleton, T., Harutyunyan, A., Bellemare, M.G.: Safe and efficient off-policy reinforcement learning. In: Lee, D.D., Sugiyama, M., von Luxburg, U., Guyon, I., Garnett, R. (eds.) Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain, pp. 1046–1054 (2016)

[25] Ren, T., Zhang, T., Lee, L., Gonzalez, J.E., Schuurmans, D., Dai, B.: Spectral decomposition representation for reinforcement learning. In: International Conference on Learning Representations (2023)

[26] Sekhari, A., Dann, C., Mohri, M., Mansour, Y., Sridharan, K.: Agnostic reinforcement learning with low-rank mdps and rich observations. In: Ranzato, M., Beygelzimer, A., Dauphin, Y.N., Liang, P., Vaughan, J.W. (eds.) Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 19033–19045 (2021)

[27] Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L., Kumaran, D., Graepel, T., et al.: A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. Science **362**(6419), 1140–1144 (2018)

[28] Singh, S., Jaakkola, T., Littman, M.L., Szepesvári, C.: Convergence results for single-step on-policy reinforcement-learning algorithms. Machine learning **38**, 287–308 (2000)

[29] Sutton, R.S., Barto, A.G.: Reinforcement learning: An introduction. MIT press (2018)

[30] Tesauro, G.: Temporal difference learning and td-gammon. Commun. ACM **38**(3), 58–68 (1995)

[31] Todorov, E., Erez, T., Tassa, Y.: Mujoco: A physics engine for model-based control. In: 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2012, Vilamoura, Algarve, Portugal, October 7-12, 2012, pp. 5026–5033, IEEE (2012)

[32] Tunyasuvunakool, S., Muldal, A., Doron, Y., Liu, S., Bohez, S., Merel, J., Erez, T., Lillicrap, T., Heess, N., Tassa, Y.: dm-control: Software and tasks for continuous control. Software Impacts **6**, 100022 (2020), ISSN 2665-9638

[33] Uehara, M., Zhang, X., Sun, W.: Representation learning for online and offline RL in low-rank mdps. In: The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022, OpenReview.net (2022)

[34] Van Hasselt, H., Guez, A., Silver, D.: Deep reinforcement learning with double q-learning. In: Thirtieth AAAI conference on artificial intelligence (2016)

[35] Vinyals, O., Babuschkin, I., Czarnecki, W.M., Mathieu, M., Dudzik, A., Chung, J., Choi, D.H., Powell, R., Ewalds, T., Georgiev, P., et al.: Grandmaster level in starcraft ii using multi-agent reinforcement learning. Nature **575**(7782), 350–354 (2019)

[36] Watkins, C.J., Dayan, P.: Q-learning. Machine learning **8**(3), 279–292 (1992)

[37] Yang, G., Ajay, A., Agrawal, P.: Overcoming the spectral bias of neural value approximation. In: The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022, OpenReview.net (2022)

[38] Yarats, D., Kostrikov, I.: Soft actor-critic (sac) implementation in pytorch (2020)

[39] Zhang, W., He, J., Zhou, D., Zhang, A., Gu, Q.: Provably efficient representation learning in low-rank markov decision processes. arXiv preprint arXiv:2106.11935 (2021)

## A   Notations

Table 3: Notations used in this work.

| Symbol | Description |
|---|---|
| $\mathcal{S}$ | State space |
| $\mathcal{A}$ | Action space |
| $R$ | Reward function |
| $R_t$ | Return $R_t = \sum_{i=t}^{T} \gamma^{i-t} r(s_i, a_i)$ |
| $r$ | A reward |
| $P$ | Transition kernel (or transition probability), $P : \mathcal{S} \times \mathcal{A} \to p(s)$ |
| $\gamma$ | Discount factor, $\gamma \in [0, 1)$ |
| $\mathcal{D}$ | Replay buffer |
| $\rho_0$ | Initial state distribution |
| $p(s)$ | A state distribution |
| $Q^\pi(s, a)$ | Action value function |
| $\tau$ | Trajectory, a sequence of state and action |
| $\pi$ | Policy |
| $V^\pi(s)$ | State value function given policy $\pi$ |
| $\alpha$ | Temperature for MaxEnt RL [11] |
| $\mathcal{H}$ | Pre-selected target entropy. |
| $\theta$ | neural network parameters for value function |
| $\phi$ | neural network parameters for policy function |
| $P^\pi_{s,a,s',a'}$ | Transition matrix, $P^\pi_{s,a,s',a'} := P(s'|s,a)\pi(a'|s')$. |
| $\mathbf{V}^\pi$ | Vector of all V value with length $|\mathcal{S}|$ |
| $\mathbf{Q}^\pi$ | Vector of all Q value with length $|\mathcal{S}| \cdot |\mathcal{A}|$ |
| $\mathbf{r}$ | Vector of all reward with length $|\mathcal{S}| \cdot |\mathcal{A}|$ |
| $\mathbf{e}$ | A column of all ones |
| $I$ | Identity matrix |
| $(\lambda_i, H_i)$ | Eigenpair |
| $(\mathfrak{B}, \|\cdot\|)$ | A Banach space equipped with a norm $\|\cdot\|$ |
| $\beta$ | Hyper-parameter for ERC, controlling trend to 1-eigensubspace |

## B   Theoretical Derivations

In this section, theorems and lemmas in the main text are restated, and are given the related proof.

### B.1   An Inherent Path of Value Approximation

**Lemma 4 (Dynamics of approximation error).** *Consider a continuous sequence $\{Q_t | t \geq 0\}$, satisfy Equation (7) with initial condition $\mathbf{Q}_0$ at time step $t = 0$, then*

$$\mathbf{Q}_t - \mathbf{Q}^* = \exp\{-t(I - \gamma P^\pi)\}(\mathbf{Q}_0 - \mathbf{Q}^*). \tag{17}$$

*Proof.* This ordinary differential equation can be solved directly with the help of Remark 2. □

**Theorem 3.** *Assumption 1 holds, then* $\mathbf{Q}_t - \mathbf{Q}^* = \alpha_1 \exp\{t(\gamma\lambda_1 - 1\}H_1 + \sum_{i=2}^{|\mathcal{S}|\cdot|\mathcal{A}|} \alpha_i \exp\{t(\gamma\lambda_i - 1)\}H_i = \alpha_1 \exp\{t(\gamma\lambda_1 - 1\}H_1 + o\Big(\alpha_1 \exp\{t(\gamma\lambda_1 - 1\}\Big)$

*Proof.* We have

$$\mathbf{Q}_0 - \mathbf{Q}^* = \sum_i \alpha_i H_i, \tag{18}$$

where $\alpha_i$ is some constants. Note that $\mathbf{Q}_t$ is in the column space of $P^\pi$ and $\{H_i\}$ can is a basis of $P^\pi$.

Recall Lemma 1, we have

$$
\begin{aligned}
\mathbf{Q}_t - \mathbf{Q}^* &= \exp\{-t(I - \gamma P^\pi)\}(\mathbf{Q}_0 - \mathbf{Q}^*) \\
&= \exp\{-t(I - \gamma P^\pi)\} \sum_i \alpha_i H_i \\
&= \sum_i \alpha_i \exp\{-t(I - \gamma P^\pi)\} H_i \\
&= \sum_i \alpha_i \exp\{-t(1 - \gamma\lambda_i)\} H_i \\
&= \alpha_1 \exp\{t(\gamma\lambda_1 - 1)\}H_1 + \sum_{i=2} \alpha_i \exp\{t(\gamma\lambda_i - 1)\}H_i \\
&= \alpha_1 \exp\{t(\gamma - 1)\}H_1 + o\Big(\alpha_1 \exp\{t(\gamma - 1\}\Big).
\end{aligned}
\tag{19}
$$

The fourth equation holds because $\exp\{-t(I - \gamma P^\pi)\}$ is also diagonalizable under the same basis $\{H_i\}$, with the eigenvalue $\exp\{t(\gamma\lambda_i - 1)\}$, $i \in \{1, 2, \cdots, |\mathcal{S}| \cdot |\mathcal{A}|\}$. □

**Lemma 5.** *Consider a Banach space* $(\mathfrak{B}, \|\cdot\|)$ *of dimension N, and let the N-dimensional Bellman error at timestep t, represented by* $\mathbf{B}^t$, *have coordinates* $(B_1, B_2, \cdots, B_N)$ *in* $(\mathfrak{B}, \|\cdot\|)$. *Within this Banach space, the projected point in the 1-eigensubspace, which is closest to* $B^t$, *is* $Z^t$ *whose coordinates are* $(z^t, z^t, \cdots, z^t)$ *where* $z^t = \frac{1}{N}\sum_{j=1}^N B_i^t$.

*Proof.* Assume $Z$ is in 1-eigensubspace whose coordinates are $(z^t, z^t, \cdots, z^t)$ with dimension N. The distance from $Z^t$ to $N^t$ is $\sum_i^N (B_i^t - z^t)^2$. To find the coordinates of $Z^t$, we can build the following optimization problem:

$$\min_z \sum_i^N (B_i^t - z^t)^2, \tag{20}$$

which is a convex optimization problem. The solution is $z^t = \frac{1}{N}\sum_{j=1}^N B_i^t$. □

## B.2   Theoretical Analysis

**Lemma 6.** *Given the ERC update rules in Equation* (10)*, ERC updates the value function in tabular form in the following way*

$$Q_{t+1} = (1 - \alpha_t(1+\beta))Q_t + \alpha_t(1+\beta)\mathcal{B}Q_t - \alpha_t\beta C_t, \tag{21}$$

*where* $\mathcal{B}Q_t(s_t, a_t) = r_t + \gamma\mathbb{E}_{s_{t+1},a_{t+1}}Q_t(s_{t+1}, a_{t+1})$, $C_t = 2\mathbb{NG}(\mathbb{E}[\mathcal{B}Q_t - Q_t])$, *and* $\mathbb{NG}$ *means stopping gradient.*

*Proof.* For $\mathcal{R}_{\text{push}}$, we have

$$
\begin{aligned}
\mathcal{R}_{\text{push}} &= \mathbb{E}_{s,a}\Big((Q - \mathbb{NG}(\mathcal{B}Q)) - \mathbb{NG}(\mathbb{E}[Q - \mathcal{B}Q])\Big)^2 \\
&= \mathbb{E}_{s,a}((Q - \mathbb{NG}(\mathcal{B}Q)))^2 - 2\mathbb{E}[\mathbb{NG}(\mathcal{B}Q) - Q]\mathbb{NG}(\mathbb{E}[\mathcal{B}Q - Q]) + (\mathbb{NG}(\mathbb{E}[Q - \mathcal{B}Q]))^2 \\
&= \mathcal{L}_{\text{PE}} - 2\mathbb{E}[\mathbb{NG}(\mathcal{B}Q) - Q]\mathbb{NG}(\mathbb{E}[\mathcal{B}Q - Q]) + (\mathbb{NG}(\mathbb{E}[\mathcal{B}Q - Q]))^2 \\
&= \mathcal{L}_{\text{PE}} - C_t\mathbb{E}[\mathbb{NG}(\mathcal{B}Q) - Q] + \underbrace{(\mathbb{NG}(\mathbb{E}[Q - \mathcal{B}Q]))^2}_{\text{No contribution to gradient}},
\end{aligned}
\tag{22}
$$

where $C_t = 2\mathbb{NG}(\mathbb{E}[\mathcal{B}Q_t - Q_t])$.

Now we consider the gradient of ERC w.r.t. $Q$.

$$
\begin{aligned}
\nabla_Q\mathcal{L}_{\text{ERC}} &= \nabla_Q\mathcal{L}_{\text{PE}} + \nabla_Q\beta\mathcal{R}_{\text{push}} \\
&= (1+\beta)\nabla_Q\mathcal{L}_{\text{PE}} - \beta C_t\nabla_Q\mathbb{E}[\mathbb{NG}(\mathcal{B}Q) - Q] \\
&= (1+\beta)(Q - \mathcal{B}Q) + \beta C_t.
\end{aligned}
\tag{23}
$$

Thus we have the following update rule

$$
\begin{aligned}
Q_{t+1} &= Q_t - \alpha_t\nabla_{Q_t}\mathcal{L}_{\text{ERC}} \\
&= Q_t - \alpha_t(1+\beta)(Q_t - \mathcal{B}Q_t) - \alpha_t\beta C_t \\
&= (1 - \alpha_t(1+\beta))Q_t + \alpha_t(1+\beta)\mathcal{B}Q_t - \alpha_t\beta C_t.
\end{aligned}
\tag{24}
$$

**Theorem 4 (Convergence of 1-Eigensubspace Regularized Value Approximation).** *Consider the Bellman backup operator* $\mathcal{B}$ *and a mapping* $Q : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$, *and* $Q^k$ *is updated with Equation* (21)*. Then the sequence* $\{Q^k\}_{k=0}^{\infty}$ *will converge to 1-eigensubspace regularized optimal* $Q$ *value of* $\pi$ *as* $k \to \infty$.

*Proof.* According to Lemma 6, the update rule of ERC can be rewritten as

$$
\begin{aligned}
Q^{k+1} &\leftarrow (1 - \alpha_t(1+\beta))Q^k + \alpha_t(1+\beta)\Big(\mathcal{B}Q^k - \frac{\beta}{1+\beta}C^k\Big) \\
&\leftarrow (1 - \alpha_t(1+\beta))Q^k + \alpha_t(1+\beta)\Big(r(s,a) + \gamma\mathbb{E}_{s',a'}Q_t(s',a') - \frac{\beta}{1+\beta}C^k\Big) \\
&\leftarrow (1 - \alpha_t(1+\beta))Q^k + \alpha_t(1+\beta)\Big(r(s,a) + \frac{\beta}{1+\beta}\mathbb{E}_{s,a}[Q^k - \mathcal{B}Q^k] \\
&\quad + \gamma\mathbb{E}_{s',a'}Q^k(s',a')\Big)
\end{aligned}
\tag{25}
$$

From Equation (25), the update target in ERC is $r(s, a) + \gamma\mathbb{E}_{s',a'}Q_t(s', a') - \frac{\beta}{1+\beta}C^k$. Define the 1-eigensubspace regularized reward as $r_{erc}(s_t, a_t) := r(s, a) + \frac{\beta}{1+\beta}\mathbb{E}_{s,a}[Q^k - \mathcal{B}^\pi Q^k]$, and the update rule of ERC can be further rewritten as

$$Q(s, a) \leftarrow r_{erc}(s, a) + \gamma\mathbb{E}[Q(s', a')], \tag{26}$$

and apply the standard convergence results for policy evaluation [28, 29], which concludes the proof.                                                              □

### B.3   ERC*

In the text, ERC*, which leverages optimal $Q^*$ to push approximation error to 1-eigensubspace, is compared to ERC and TD in Figure 3. We give a detailed description of the ERC* in this section. Following Lemma 3, the update rule of ERC* can be given as

$$Q_{t+1} = (1 - \alpha_t(1 + \beta))Q_t + \alpha_t\mathcal{T}Q_t + \alpha_t\beta(Q^* + \mathbb{E}(Q_t - Q^*)). \tag{27}$$

The procedure of the proof is the same as in Lemma 3.

## C   Additional Details Regarding Experiments

In this section, we provide a detailed description of the experimental procedures and configurations used to generate the tables and figures presented in the main text.

**Implementations.** To ensure the validity and reproducibility of the experiments, we have fixed all random seeds, including but not limited to those used in PyTorch, Numpy, Gym, Random, and CUDA packages, across all experiments. For random seeds, unless otherwise specified, we evaluate each tested algorithm over 10 fixed random seeds. For more implementation details, please refer to our code.

**Comparison with REDQ and TQC.** Both REDQ and TQC leverage an ensemble mechanism to obtain a more accurate and unbiased Q estimation. Specifically, REDQ uses 10 critics and the UTD ratio is 20, i.e., REDQ performs gradient updates twenty times for every interaction with the environment. As for TQC, it uses 5 distributional critics. However, ERC does not leverage any ensemble critic, distributional value functions, or high UTD ratio, yet it still outperforms them.

**Figure 3.** The experiments are conducted on FrozenLake-v1 environment. The true Q values are computed through the Monte Carlo method. The shaded area represents a standard deviation over trials. In the case of no additional description, the shaded areas carried in all subsequent figures indicate a standard deviation over ten random seeds. The environment used in the experiment has 16 states and 4 actions. The learning rate is 0.01. The discounted factor is 0.9. A fixed $\beta$ value of 0.3 is used for both ERC and ERC*.

**Figure 4.** The performance curves for OpenAI gym continuous control tasks on the DeepMind Control suite. The shaded region represents a 50% standard deviation of the average evaluation over 10 seeds. and the curves are smoothed with a moving average window of size 10. The evaluation of each algorithm is conducted over a period of 1 million timesteps, with the average return of the algorithm being evaluated every 10k timesteps over ten episodes.

**Figure 6.** The shaded region in these figures also represents a 50% standard deviation of the average evaluation over 10 seeds and the curves are smoothed with a moving average window of size 10. The evaluation of each algorithm is conducted over a period of 1 million timesteps, with the index of dispersion of the algorithm being evaluated every 10k timesteps over ten episodes. For a better visualization effect, the y-axis is clipped to $[-1, 5]$.

**Figure 5.** The experimental settings are consistent with those of Figure 6.

**Table 1.** Average Return after 1M timesteps training on DMControl suite. $\pm$ indicates a standard deviation over ten trials. The DMControl suite benchmarks contain 28 tasks, and our evaluation of the algorithms was conducted on 26 tasks, with the tasks 'acrobot-swingup_sparse' and 'manipulator-bring_ball' excluded due to the inability of all tested DRL algorithms to produce meaningful results on these challenging tasks.

**Table 2.** The experimental settings are consistent with those of Table 1.

# D    Additional Experimental Results

In this section, we present additional experimental results.

## D.1    The Inherent Path in Practice

Section 3.1 discusses the inherent path of approximation error. Does this occur in the practical scene? To empirically show the path of the approximation error, we visualize the path of approximation error given a fixed policy. The results are given in Figure 7, where we perform experiments on FrozenLake-v1 environment since the true value $Q^*$ can be evaluated by the MC Method. For comparison, we also visualize the path of value function learning by the Monte Carlo method. The approximation error of the TD method is optimized toward the 1-eigensubspace before ultimately converging to zero rather than directly toward the optimum that the Monte Carlo (MC) method acts. The inherent path in Figure 7a is consistent with Theorem 1. Thus, such a path is non-trivial and motivates us to improve the value approximation by guiding the Bellman error to 1-eigensubspace, resulting in an efficient and robust path.

## D.2    Value Approximation on Additional Environment

In the text, we show the value approximation process for various algorithms in Figure 3 on FrozenLake-v1 environment. Now We offer the Value function approximation process for multiple methods on CliffWalking-v0 environment in Figure 8.

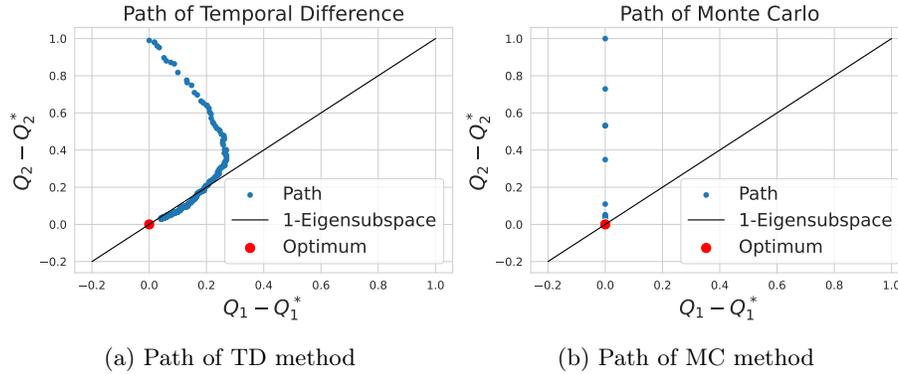(a) Path of TD method                (b) Path of MC method

Fig. 7: Value approximation path for TD and MC methods. (a) illustrates that there exists an inherent path that approximation error approaches 1-eigensubspace before converging to zero. The empirical fact is consistent with our theoretical analysis in Theorem 1. (b) demonstrates the path of the Value approximation error for the MC method. The value function approaches the true value with the shortest path because the MC method obtains the unbiased true value of the objective by a large number of samples.

### D.3   Additional Figures

Due to space limitations, we present only selective figures in the text. In this section, we give the remaining figures for the 26 tasks.

Figure 9 presents the learning curves from scratch. Figure 10 shows approximation error curves. Figure 11 is index of dispersion curves. We note that TQC has the best approximation error on 'cartpole-swingup_sparse' task. This is because TQC outperforms the other algorithms substantially on this task as well.
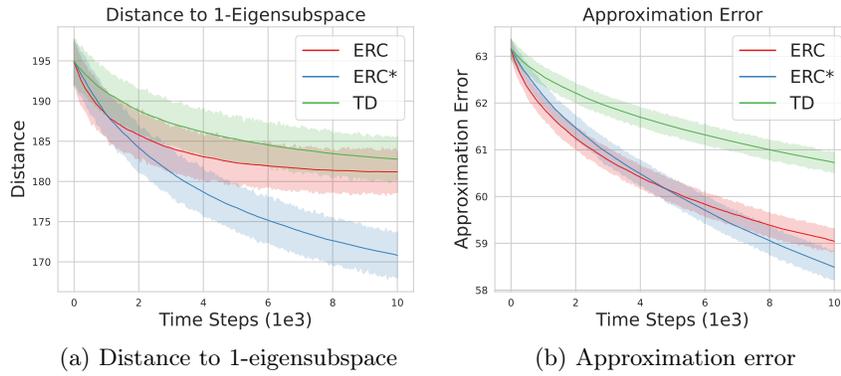
(a) Distance to 1-eigensubspace          (b) Approximation error

Fig. 8: Value function approximation process for various methods on CliffWalking-v0 environment. (a) illustrates the distance between the approximation error and the 1-eigensubspace for various methods, where ERC* denotes ERC utilizing an oracle $Q*$ to push the approximation error towards the 1-eigensubspace. The results demonstrate that the ERC method is closer to the 1-eigensubspace at the same time compared to the TD method. Besides, the ERC* is closer to 1-eigensubspace because it leverages oracle true $Q^*$. (b) represents the absolute approximation error for various algorithms. The result illustrates that the ERC method has a smaller approximation error at the same time than the TD method. ERC is close to ERC* in this metric, which means that although ERC* is more efficient near the 1-eigensubspace, the performance of the ERC method is also almost the same as ERC* in terms of the approximation error metric that we really care about. The shaded area represents a standard deviation over ten trials.
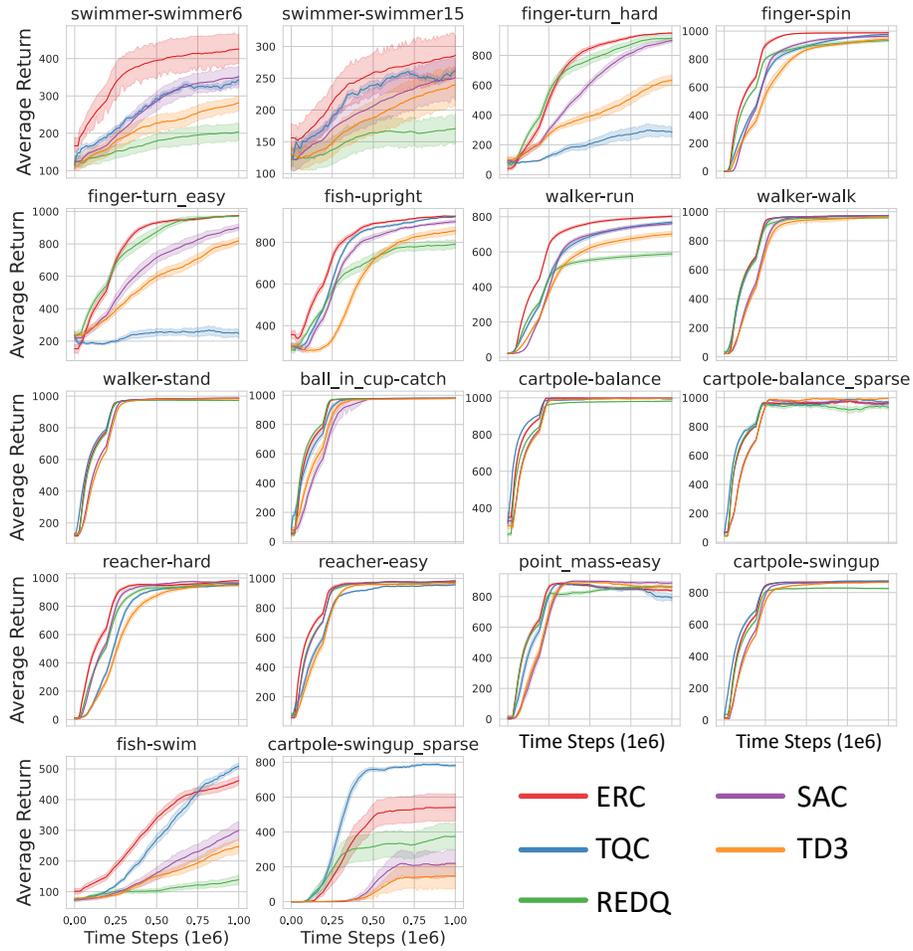
Fig. 9: Performance curves for OpenAI gym continuous control tasks on Deep-Mind Control suite. The proposed algorithm, ERC, is observed to significantly outperform the other tested algorithms. The shaded region represents half of the standard deviation of the average evaluation over 10 seeds. The curves are smoothed with a moving average window of size ten.
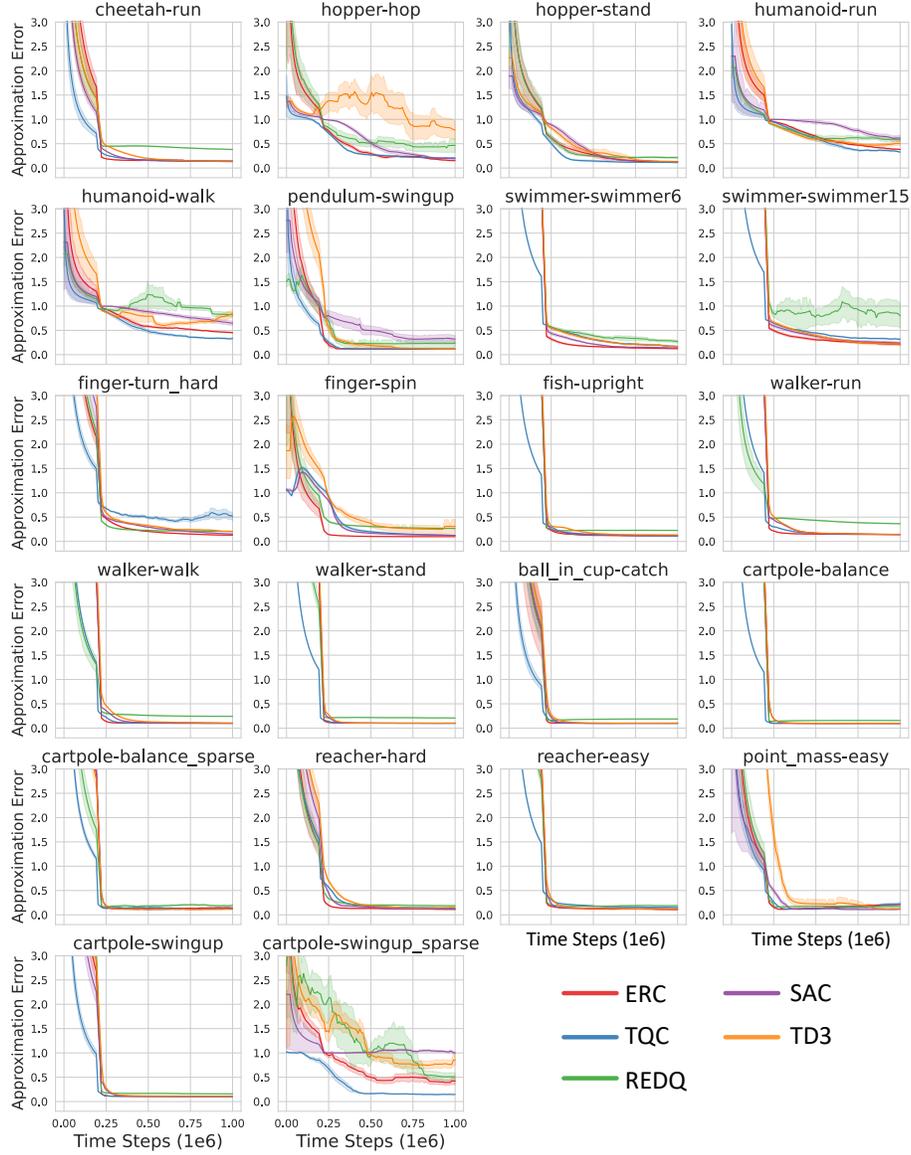
Fig. 10: Approximation error curves. The results demonstrate that the approximation error of ERC is empirically minimal when compared to other algorithms (such as TQC and REDQ) that are specifically designed to obtain accurate unbiased value estimation.
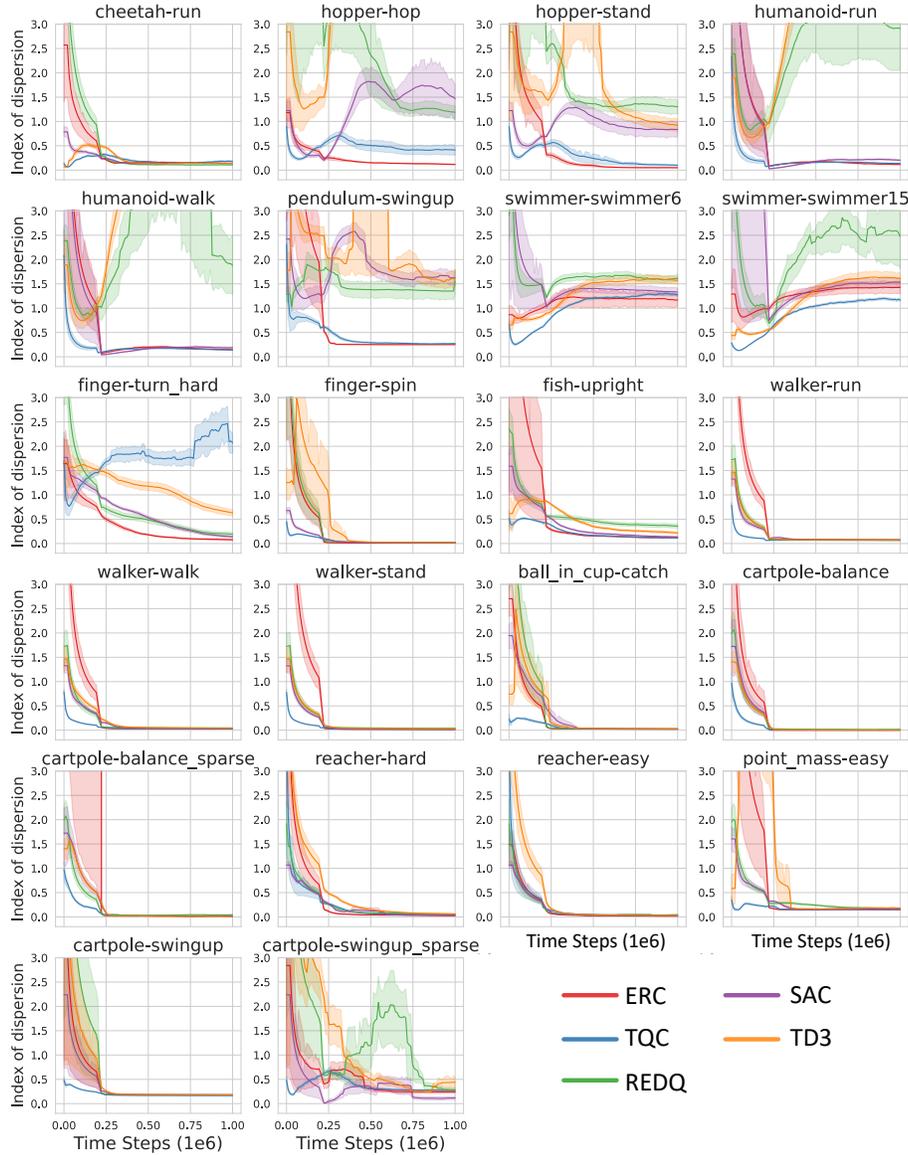
Fig. 11: Index of dispersion curves. The results demonstrate that ERC effectively controls the variance. Furthermore, the variance of ERC is observed to be minimal on the selected tasks, even in situations where the performance of ERC is not as good as TQC and REDQ