

---

# Wasserstein Flow Meets Replicator Dynamics: A Mean-Field Analysis of Representation Learning in Actor-Critic

---

**Yufeng Zhang**<sup>†</sup>

Northwestern University  
yufengzhang2023@u.northwestern.edu

**Siyu Chen**<sup>†</sup>

Tsinghua University  
chensy18@mails.tsinghua.edu.cn

**Zhuoran Yang**

Princeton University  
zy6@princeton.edu

**Michael I. Jordan**

UC Berkeley  
jordan@cs.berkeley.edu

**Zhaoran Wang**

Northwestern University  
zhaoranwang@gmail.com

## Abstract

Actor-critic (AC) algorithms, empowered by neural networks, have had significant empirical success in recent years. However, most of the existing theoretical support for AC algorithms focuses on the case of linear function approximations, or linearized neural networks, where the feature representation is fixed throughout training. Such a limitation fails to capture the key aspect of representation learning in neural AC, which is pivotal in practical problems. In this work, we take a mean-field perspective on the evolution and convergence of feature-based neural AC. Specifically, we consider a version of AC where the actor and critic are represented by overparameterized two-layer neural networks and are updated with two-timescale learning rates. The critic is updated by temporal-difference (TD) learning with a larger stepsize while the actor is updated via proximal policy optimization (PPO) with a smaller stepsize. In the continuous-time and infinite-width limiting regime, when the timescales are properly separated, we prove that neural AC finds the globally optimal policy at a sublinear rate. Additionally, we prove that the feature representation induced by the critic network is allowed to evolve within a neighborhood of the initial one.

## 1 Introduction

In reinforcement learning (RL) [56], an agent aims to learn the optimal policy that maximizes the expected total reward by interacting with the environment. Policy-based RL algorithms achieve such a goal by directly optimizing the expected total reward as a function of the policy, which often involves two components: policy evaluation and policy improvement. Specifically, policy evaluation refers to estimating the value function of the current policy, which characterizes the performance of the current policy and reveals the updating direction for finding a better policy, which is known as policy improvement. Algorithms with these two ingredients are also called actor-critic (AC) methods [36], where the actor and the critic refer to the policy and its corresponding value function, respectively.

---

<sup>†</sup> Equal contribution.

Recently, in RL applications with large state spaces, actor-critic empowered by expressive function approximators such as neural networks have achieved striking empirical successes [3, 4, 9, 20, 51, 52, 60]. These successes benefit from the data-dependent representations learned by neural networks. Unfortunately, however, the theoretical understanding of this data-dependent benefit is very limited. The classical theory of AC focuses on the case of linear function approximation, where the actor and critic are represented using linear functions with the feature mapping fixed throughout learning [10, 11, 36]. Meanwhile, a few recent works establish convergence and optimality of AC with overparameterized neural networks [26, 39, 61], where the neural network training is captured by the Neural Tangent Kernel (NTK) [30]. Specifically, with properly designed parameter initialization and stepsizes, and sufficiently large network widths, the neural networks employed by both actor and critic can be assumed to be well approximated by linear functions of a random feature determined by initial parameters. In other words, concerning representation learning, the features induced by these algorithms are by assumption infinitesimally close to the initial featural representation, which is data-independent.

In this work, we make initial steps towards understanding how representation learning comes into play in neural AC. Specifically, we address the following questions:

*Going beyond the NTK regime, does neural AC provably find the globally optimal policy? How does the feature representation associated with the neural network evolve along with neural AC?*

We focus on a version of AC where the critic performs temporal-difference (TD) learning [55] for policy evaluation and the actor improves its policy via proximal policy optimization (PPO) [49], which corresponds to a Kullback-Leibler (KL) divergence regularized optimization problem, with the critic providing the update direction. Moreover, we utilize two-timescale updates where both the actor and critic are updated at each iteration but with the critic having a much larger stepsize. In other words, the critic is updated at a faster timescale. Meanwhile, we represent the critic explicitly as a two-layer overparameterized neural network and parameterize the actor implicitly via the critic and PPO updates. To examine convergence, we study the evolution of the actor and critic in the continuous-time limiting regime with the network width going to infinity. In such a regime, the actor update is closely connected to replicator dynamics [12, 28, 50] and the critic update is captured by a semigradient flow in the Wasserstein space [59]. Moreover, the semigradient flow runs at a faster timescale according to the two-timescale mechanism.

It turns out that the separation of timescales plays an important role in the convergence analysis. In particular, in the continuous-time limit, it enables us to first separately analyze the evolution of actor and critic and then combine these results to get final theoretical guarantees. Specifically, focusing solely on the actor, we prove that the time-averaged suboptimality of the actor converges sublinearly to zero up to the time-averaged policy evaluation error associated with critic updates. Moreover, for the critic, under proper regularity conditions, we connect the Bellman error to the Wasserstein distance and show that the time-averaged policy evaluation error also converges sublinearly to zero. Therefore, we show that neural AC provably achieves global optimality at a sublinear rate. Furthermore, regarding representation learning, we show that the critic induces a data-dependent feature representation within an  $O(1/\alpha)$  neighborhood of the initial representation in terms of the Wasserstein distance, where  $\alpha$  is a sufficiently large scaling parameter.

The key to our technical analysis reposes on three ingredients: (i) infinite-dimensional variational inequalities with a one-point monotonicity [27], (ii) a mean-field perspective on neural networks [19, 41, 42, 53, 54], and (iii) the two-timescale stochastic approximation [13, 37]. In particular, in the infinite-width limit, the neural network and its induced feature representation are identified with a distribution over the parameter space. The mean-field perspective enables us to characterize the evolution of such a distribution within the Wasserstein space via a partial differential equation (PDE) [5, 6, 58, 59]. For policy evaluation, such a PDE is given by the semigradient flow induced by TD learning. We characterize the error of policy evaluation by showing that mean-field Bellman error satisfies a version of one-point monotonicity tailored to the Wasserstein space. Moreover, our actor analysis utilizes the geometry of policy optimization, which shows that the expected total reward,

as a function of the policy, also enjoys the property of one-point monotonicity in the policy space. Finally, the actor and critic errors are connected via two-timescale stochastic approximation. To the best of our knowledge, this is the first time that convergence and global optimality guarantees have been obtained for neural AC.

**Related Work.** AC with linear function approximation has been studied extensively in the literature. In particular, using a two-timescale stochastic approximation via ordinary differential equations, [10, 11, 36] establish asymptotic convergence guarantees in the continuous-time limiting regime. More recently, using more sophisticated optimization techniques, various works [29, 35, 64–66] have established discrete-time convergence guarantees that show that linear AC converges sublinearly to either a stationary point or the globally optimal policy. Furthermore, when overparameterized neural networks are employed, [26, 39, 61] prove that neural AC converges to the global optimum at a sublinear rate. In these works, the initial value of the network parameters and the learning rates are chosen such that both actor and critic updates are captured by the NTK. In other words, when the network width is sufficiently large, such a version of neural AC is well approximated by its linear counterpart via the neural tangent feature. In comparison, we establish a mean-field analysis that has a different scaling than the NTK regime. We also establish finite-time convergence to global optimality, and more importantly, the feature representation induced by the critic is data-dependent and allowed to evolve within a much larger neighborhood around the initial one.

Furthermore, our work is also related to the recent line of research on understanding stochastic gradient descent (SGD) for supervised learning problems involving an overparameterized two-layer neural network under the mean-field regime. See, e.g., [16, 19, 21, 22, 31, 40–42, 53, 54, 63] and the references therein. In the continuous-time and infinite-width limit, these works show that SGD for neural network training is captured by a Wasserstein gradient flow [5, 6, 59] of an energy function that corresponds to the objective function in supervised learning. In contrast, our analysis combines such a mean-field analysis with TD learning and two-timescale stochastic approximation, which are tailored specifically to AC. Moreover, our critic is updated via TD learning, which is a semigradient algorithm and there is no objective functional making TD learning a gradient-based algorithm. Thus, in the mean-field regime, our critic is given by a Wasserstein semigradient flow, which also differs from these existing works.

Additionally, our work is closely related to [1, 69], who provide mean-field analyses for neural TD-learning and Q-learning [62]. In comparison, we focus on AC, which is a two-timescale policy optimization algorithm. Finally, [2] studies softmax policy gradient with neural network policies in the mean-field regime, where policy gradient is cast as a Wasserstein gradient flow with respect to the expected total reward. The algorithm assumes that the critic directly gets the desired value function and thus the algorithm is single-timescale. Moreover, the convergence guarantee in [2] is asymptotic. In comparison, our AC is two-timescale and we establish non-asymptotic sublinear convergence guarantees to global optimality.

**Notation.** We denote by  $\mathcal{P}(\mathcal{X})$  the set of probability measures over the measurable space  $\mathcal{X}$ . Given a curve  $\rho : \mathbb{R} \rightarrow \mathcal{X}$ , we denote by  $\dot{\rho}_s = \partial_t \rho_t |_{t=s}$  its derivative with respect to the time. For an operator  $F : \mathcal{X} \rightarrow \mathcal{X}$  and a measure  $\mu \in \mathcal{P}(\mathcal{X})$ , we denote by  $F_{\#}\mu = \mu \circ F^{-1}$  the push forward of  $\mu$  through  $F$ . We denote by  $\chi^2(\rho \| \mu)$  the chi-squared divergence between probability measures  $\rho$  and  $\mu$ , which is defined as  $\chi^2(\rho \| \mu) = \int (\rho/\mu - 1)^2 d\mu$ . Given two probability measures  $\rho$  and  $\mu$ , we denote the Kullback-Leibler divergence or the relative entropy from  $\mu$  to  $\rho$  by  $\text{KL}(\rho \| \mu) = \int \log(\rho/\mu) d\rho$ . For  $\nu_1, \nu_2, \mu \in \mathcal{P}(\mathcal{X})$ , we define the  $\dot{H}^{-1}(\mu)$  weighted homogeneous Sobolev norm as  $\|\nu_1 - \nu_2\|_{\dot{H}^{-1}(\mu)} = \sup \{ |\langle f, \nu_1 - \nu_2 \rangle| \mid \|f\|_{\dot{H}^1(\mu)} \leq 1 \}$ . We denote by  $\|f(x)\|_{p,\mu} = (\int |f(x)|^p \mu(dx))^{1/p}$  the  $\ell_p$ -norm with respect to probability measure  $\mu$ . We denote by  $\otimes$  the semidirect product, i.e.,  $\mu \otimes K = K(y|x)\mu(x)$  for  $\mu \in \mathcal{P}(\mathcal{X})$  and transition kernel  $K : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{Y})$ . For a function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , we denote by  $\text{Lip}(f) = \sup_{x,y \in \mathcal{X}, x \neq y} |f(x) - f(y)| / \|x - y\|$  its Lipschitz constant. We denote a normal distribution on  $\mathbb{R}^D$  by  $\mathcal{N}(\mu, \Sigma)$ , where  $\mu$  is the mean value and  $\Sigma$  is the covariance matrix.

## 2 Background

In this section, we first introduce the policy optimization problem and the actor-critic method. We then present the definition of the Wasserstein space.

### 2.1 Policy Optimization and Actor-Critic

We consider a Markov decision process (MDP) given by  $(\mathcal{S}, \mathcal{A}, \gamma, P, r, \mathcal{D}_0)$ , where  $\mathcal{S} \subseteq \mathbb{R}^{d_1}$  is the state,  $\mathcal{A} \subseteq \mathbb{R}^{d_2}$  is the action space,  $\gamma \in (0, 1)$  is the discount factor,  $P : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$  is the transition kernel,  $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}_+$  is the reward function, and  $\mathcal{D}_0 \in \mathcal{P}(\mathcal{S})$  is the initial state distribution. Without loss of generality, we assume that  $\mathcal{S} \times \mathcal{A} \subseteq \mathbb{R}^d$  and  $\|(s, a)\|_2 \leq 1$ , where  $d = d_1 + d_2$ . We remark that as long as the state-action space is bounded, we can normalize the space to be within the unit sphere. Given a policy  $\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$ , at the  $m$ th step, the agent takes an action  $a_m$  at state  $s_m$  according to  $\pi(\cdot | s_m)$  and observes a reward  $r_m = r(s_m, a_m)$ . The environment then transits to the next state  $s_{m+1}$  according to the transition kernel  $P(\cdot | s_m, a_m)$ . Note that the policy  $\pi$  induces Markov chains on both  $\mathcal{S}$  and  $\mathcal{S} \times \mathcal{A}$ . Considering the Markov chain on  $\mathcal{S}$ , we denote the induced Markov transition kernel by  $P^\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{S})$ , which is defined as  $P^\pi(s' | s) = \int_{\mathcal{A}} P(s' | s, a) \pi(da | s)$ . Likewise, we denote the Markov transition kernel on  $\mathcal{S} \times \mathcal{A}$  by  $\tilde{P}^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S} \times \mathcal{A})$ , which is defined as  $\tilde{P}^\pi(s', a' | s, a) = \pi(a' | s') P(s' | s, a)$ . Let  $\tilde{\mathcal{D}}$  be a probability measure on  $\mathcal{S} \times \mathcal{A}$ . We then define the visitation measure induced by policy  $\pi$  and starting from  $\tilde{\mathcal{D}}$  as

$$\tilde{\mathcal{E}}_{\tilde{\mathcal{D}}}^\pi(d(s, a)) = (1 - \gamma) \cdot \sum_{m \geq 0} \gamma^m \cdot \mathbb{P}((s_m, a_m) \in d(s, a) | (s_0, a_0) \sim \tilde{\mathcal{D}}), \quad (2.1)$$

where  $(s_m, a_m)$  is the trajectory generated by starting from  $(s_0, a_0) \sim \tilde{\mathcal{D}}$  and following policy  $\pi$  thereafter. If  $\tilde{\mathcal{D}} = \mathcal{D} \otimes \pi$  holds, we then denote such a visitation measure by  $\tilde{\mathcal{E}}_{\mathcal{D}}^\pi$ . Furthermore, we denote by  $\mathcal{E}(ds) = \int_{\mathcal{A}} \tilde{\mathcal{E}}(ds, da)$  the marginal distribution of visitation measure  $\tilde{\mathcal{E}}$  with respect to  $\mathcal{S}$ . In particular, when  $(s_0, a_0) \sim \mathcal{D} \otimes \pi$  holds in (2.1), it follows that  $\tilde{\mathcal{E}}_{\mathcal{D}}^\pi = \mathcal{E}_{\mathcal{D}}^\pi \otimes \pi$ . In policy optimization, we aim to maximize the expected total rewards  $J(\pi)$  defined as follows,

$$J(\pi) = \mathbb{E}^\pi \left[ \sum_{m \geq 0} \gamma^m \cdot r(s_m, a_m) \mid s_0 \sim \mathcal{D}_0 \right],$$

where we denote by  $\mathbb{E}^\pi$  the expectation with respect to  $a_m \sim \pi(\cdot | s_m)$  and  $s_{m+1} \sim P(\cdot | s_m, a_m)$  for  $m \geq 0$ . We define the action value function  $Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  and the state value function  $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$  as follows,

$$Q^\pi(s, a) = \mathbb{E}^\pi \left[ \sum_{m \geq 0} \gamma^m \cdot r(s_m, a_m) \mid s_0 = s, a_0 = a \right], \quad V^\pi(s) = \langle Q^\pi(s, \cdot), \pi(\cdot | s) \rangle_{\mathcal{A}}, \quad (2.2)$$

where we denote by  $\langle \cdot, \cdot \rangle_{\mathcal{A}}$  the inner product on the action space  $\mathcal{A}$ . Correspondingly, the advantage function  $A^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is defined as

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s).$$

It is known that the action value function  $Q^\pi$  is the unique global minimizer to the following mean-squared Bellman error (MSBE),

$$\text{MSBE}(Q; \pi) = \frac{1}{2} \mathbb{E}_{(s, a) \sim \tilde{\Phi}^\pi} \left[ (Q(s, a) - r(s, a) - \gamma \mathbb{E}_{(s', a') \sim \tilde{P}^\pi(\cdot | s, a)} [Q(s', a')])^2 \right], \quad (2.3)$$

where  $\tilde{\Phi}^\pi$  is a weighting distribution depending on policy  $\pi$  and is with full support, i.e.,  $\text{supp}(\tilde{\Phi}^\pi) = \mathcal{S} \times \mathcal{A}$ . Therefore, the policy optimization problem can be written as the following bilevel optimization problem,

$$\max_{\pi} J(\pi) = \mathbb{E}_{s \sim \mathcal{D}_0} \left[ \langle Q^\pi(s, \cdot), \pi(\cdot | s) \rangle_{\mathcal{A}} \right], \quad \text{subject to } Q^\pi = \underset{Q}{\text{argmin}} \text{MSBE}(Q; \pi). \quad (2.4)$$

The inner problem in (2.4) is known as a policy evaluation subproblem, while the outer problem is the policy improvement subproblem. One of the most popular way to solve the policy optimization problem is actor-critic (AC) methods [56], where the job of the critic is to evaluate current policy and then the actor updates its policy according to the critic's evaluation.

## 2.2 Wasserstein Space

Let  $\Theta \subseteq \mathbb{R}^D$  be a Polish space. We denote by  $\mathcal{P}_2(\Theta) \subseteq \mathcal{P}(\Theta)$  the set of probability measures with finite second moments. Then, the Wasserstein-2 ( $W_2$ ) distance between  $\mu, \nu \in \mathcal{P}_2(\Theta)$  is defined as follows,

$$W_2(\mu, \nu) = \inf \left\{ \mathbb{E}[\|X - Y\|^2]^{1/2} \mid \text{law}(X) = \mu, \text{law}(Y) = \nu \right\},$$

where the infimum is taken over the random variables  $X$  and  $Y$  on  $\Theta$  and we denote by  $\text{law}(X)$  the distribution of a random variable  $X$ . We call  $\mathcal{M} = (\mathcal{P}_2(\Theta), W_2)$  the Wasserstein ( $W_2$ ) space, which is an infinite-dimensional manifold [59]. See §A.1 for more details.

## 3 Algorithm

**Two-timescale Actor-critic.** We consider a two-timescale Actor-critic (AC) algorithm [34, 45] for the policy optimization problem in (2.4). For policy evaluation, we parameter the critic  $Q$  with a neural network and update the parameter via temporal-difference (TD) learning [55]. For policy improvement, we update the actor policy  $\pi$  via proximal policy optimization (PPO) [49]. Our algorithm is two-timescale since both the actor and critic are updated at each iteration with different stepsizes. Specifically, we parameterize the critic  $Q$  by the following neural network with width  $M$  and parameter  $\hat{\theta} = (\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(M)}) \in \mathbb{R}^{D \times M}$ ,

$$Q_{\hat{\theta}}(s, a) = \frac{\alpha}{M} \sum_{i=1}^M \sigma(s, a; \hat{\theta}^{(i)}). \quad (3.1)$$

Here  $\sigma(s, a; \theta) : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^D \rightarrow \mathbb{R}$  is the activation function and  $\alpha > 0$  is the scaling parameter. Such a structure also appears in [17, 18, 41]. In a discrete-time finite-width (DF) scenario, at the  $k$ th iteration, the critic and actor are updated as follows,

$$\text{DF-TD: } \hat{\theta}_{k+1}^{(i)} = \hat{\theta}_k^{(i)} - \frac{\varepsilon'}{\alpha} (Q_{\hat{\theta}_k}(s_k, a_k) - r(s_k, a_k) - \gamma Q_{\hat{\theta}_k}(s'_k, a'_k)) \nabla_{\theta} \sigma(s, a; \hat{\theta}_k^{(i)}), \quad (3.2)$$

$$\text{DF-PPO: } \hat{\pi}_{k+1}(\cdot | s) = \underset{\pi}{\operatorname{argmax}} \left\{ \langle Q_{\hat{\theta}_k}(s, \cdot), \pi(\cdot | s) \rangle_{\mathcal{A}} - \varepsilon^{-1} \cdot \text{KL}(\pi(\cdot | s) \parallel \hat{\pi}_k(\cdot | s)) \right\}, \quad (3.3)$$

where  $(s_k, a_k) \sim \tilde{\Phi}^{\hat{\pi}_k}$  and  $(s'_k, a'_k) \sim \tilde{P}^{\hat{\pi}_k}(\cdot | s_k, a_k)$ . Here  $\hat{\pi}_k$  is the policy for the actor at the  $k$ th iteration,  $\tilde{\Phi}^{\hat{\pi}_k}$  is the corresponding weighting distribution,  $\varepsilon$  and  $\varepsilon'$  are the stepsizes for the DF-PPO update and the DF-TD update, respectively. In (3.2), the scaling of  $\alpha^{-1}$  arises since our update falls into the lazy-training regime [18]. In the sequel, we denote by  $\eta = \varepsilon'/\varepsilon$  the relative TD timescale. Note that in a double-loop AC algorithm, the critic can usually be solved with high precision. In the two-timescale AC however, even with the KL-divergence term in (3.3) which regularizes the policy update and helps to improve the local estimation quality of the TD update, the critic  $Q_{\hat{\theta}_k}$  for updating the actor's policy  $\hat{\pi}_k$  can still be far from the true action value function  $Q^{\hat{\pi}_k}$ . Since the policy evaluation problem is not fully solved at each iteration, the two-timescale AC can be more efficient in computation while more challenging to establish a theoretical guarantee.

**Mean-field (MF) Limit.** To analyze the convergence of the two-timescale AC with neural networks, we employ the analysis that studies the mean-field limit regime [41, 42]. Here, by saying the mean-field limit, we refer to the infinite-width limit, i.e.,  $M \rightarrow \infty$  for the neural network width  $M$  in (3.1), and the continuous-time limit, i.e.,  $t = k\varepsilon$  where  $\varepsilon \rightarrow 0$  for the stepsize in (3.2) and (3.3). For  $\hat{\theta} = \{\hat{\theta}^{(i)}\}_{i=1}^M$  independently sampled from a distribution  $\rho$ , we can write the infinite-width limit of (3.1) as

$$Q(s, a; \rho) = \alpha \int \sigma(s, a; \theta) \rho(d\theta). \quad (3.4)$$

In the sequel, we denote by  $\hat{\rho}_k$  the distribution of  $\hat{\theta}_k^{(i)}$  for the infinite-width limit of the neural network at the  $k$ th iteration. We further let  $\rho_t$  and  $\pi_t$  be the continuous-time limits of  $\hat{\rho}_k$  and  $\hat{\pi}_k$ , respectively.

As studied in [69], the mean-field limit of the DF-TD update in (3.2) is

$$\text{MF-TD: } \partial_t \rho_t = -\eta \operatorname{div}(\rho_t \cdot g(\cdot; \rho_t, \pi_t)), \quad (3.5)$$

where  $\eta$  is the relative TD timescale and

$$g(\theta; \rho, \pi) = -\mathbb{E}_{\tilde{\Phi}^\pi} \left\{ [Q(s, a; \rho) - r(s, a) - \gamma \cdot Q(s', a'; \rho)] \cdot \alpha^{-1} \nabla_\theta \sigma(s, a; \theta) \right\} \quad (3.6)$$

is a vector field. Here  $\mathbb{E}_{\tilde{\Phi}^\pi}$  is taken with respect to  $(s, a) \sim \tilde{\Phi}^\pi$  and  $(s', a') \sim \tilde{P}^\pi(\cdot | s, a)$ . It remains to characterize the mean-field limit of the DF-PPO update in (3.3). By solving the maximization problem in (3.3), the infinite-width limit of DF-PPO update can be written in closed form as

$$\varepsilon^{-1} \cdot \left\{ \log [\hat{\pi}_{k+1}(a | s)] - \log [\hat{\pi}_k(a | s)] \right\} = Q(s, a; \hat{\rho}_k) - \hat{Z}_k(s),$$

where  $\hat{Z}_k(s)$  is the normalizing factor such that  $\int \hat{\pi}_k(da | s) = 1$  for any  $s \in \mathcal{S}$ . By letting  $t = k\varepsilon$  and  $\varepsilon \rightarrow 0$ , we have  $\partial_t \log \pi_t = Q_t - Z_t$ , which can be further written as  $\partial_t \pi_t = \pi_t \cdot (Q_t - Z_t)$ . Here we have  $Q_t(a, s) = Q(a, s; \rho_t)$  and  $Z_t$  is the continuous-time limit of  $\hat{Z}_k$ . Furthermore, noting that  $\partial_t \int \pi_t(da | s) = 0$ , the mean-field limit of the DF-PPO update in (3.3) is

$$\text{MF-PPO: } \frac{d}{dt} \pi_t = \pi_t \cdot A_t, \quad \text{where } A_t(s, a) = Q_t(s, a) - \int Q_t(s, a) \pi_t(da | s). \quad (3.7)$$

The two updates (3.5) and (3.7) correspond to the mean-field limits of (3.2) and (3.3), respectively, and together serve as the mean-field limit of the two-timescale AC. In particular, we remark that the MF-TD update in (3.5) for the critic is captured by a semigradient flow in the Wasserstein space [59] while the MF-PPO update in (3.7) for the actor resembles the replicator dynamics [12, 28, 50]. Note that such a framework is applicable to continuous state and action space. In this paper, we aim to provide a theoretical analysis of the mean-field limit of the two-timescale AC.

## 4 Main Result

In this section, we first establish the convergence of the MF-PPO update in §4.1. Then, with additional assumptions, we establish the optimality and convergence of the mean-field two-timescale AC in §4.2.

### 4.1 Convergence of Mean-field PPO

For the MF-PPO update in (3.7), we establish the following theorem on its global optimality and convergence rate.

**Theorem 4.1** (Convergence of MF-PPO). Let  $\pi^* = \operatorname{argmax}_\pi J(\pi)$  be the optimal policy and  $\pi_0$  be the initial policy. Then, it holds that

$$\frac{1}{T} \int_0^T (J(\pi^*) - J(\pi_t)) dt \leq \frac{\zeta}{T} + 4\kappa \cdot \underbrace{\frac{1}{T} \int_0^T \|Q_t - Q^{\pi_t}\|_{2, \tilde{\phi}^{\pi_t}} dt}_{\text{policy evaluation error}}, \quad (4.1)$$

where  $\tilde{\phi}^{\pi_t} \in \mathcal{P}(\mathcal{S} \times \mathcal{A})$  is an evaluation distribution for the policy evaluation error and  $\zeta = \mathbb{E}_{s \sim \mathcal{E}_{\mathcal{D}_0}^{\pi^*}} \left[ \operatorname{KL}(\pi^*(\cdot | s) \| \pi_0(\cdot | s)) \right]$  is the expected KL-divergence between  $\pi^*$  and  $\pi_0$ . Furthermore, letting  $\tilde{\phi}^{\pi_t} = \frac{1}{2} \tilde{\phi}_0 + \frac{1}{2} \phi_0 \otimes \pi_t$ , where  $\tilde{\phi}_0 \in \mathcal{P}(\mathcal{S} \times \mathcal{A})$  is a base distribution and  $\phi_0 = \int_{\mathcal{A}} \tilde{\phi}_0(\cdot, da)$ , the concentrability coefficient  $\kappa$  is then given by

$$\kappa = \left\| \frac{\tilde{\mathcal{E}}_{\mathcal{D}_0}^{\pi^*}}{\tilde{\phi}_0} \right\|_\infty.$$

*Proof.* See §B.1 for a detailed proof. □

The concentrability coefficient commonly appears in the reinforcement learning literature [7, 23, 24, 38, 39, 43, 48, 57, 61]. In contrast to a more standard concentrability coefficient form, note that  $\kappa$  is irrelevant to the update of the algorithm. To show the convergence of the MF-PPO, our condition here is much weaker since we only need a given base distribution  $\tilde{\phi}_0$  such that  $\kappa < \infty$ .

Theorem 4.1 shows that the MF-PPO converges to the globally optimal policy at a rate of  $\mathcal{O}(T^{-1})$  up to the policy evaluation error. Such a theorem implies the global optimality and convergence of a double-loop AC algorithm, where the critic  $Q_t$  is solved to high precision and the policy evaluation error is sufficiently small. In the sequel, we consider a more challenging setting, where the critic  $Q_t$  is updated simultaneously along with the update of the actor’s policy  $\pi_t$ .

## 4.2 Global Optimality and Convergence of Two-timescale AC

In what follows, we aim to characterize the upper bound of the policy evaluation error when the critic and the actor are updated simultaneously. Specifically, the actor is updated via MF-PPO in (3.7) and the critic  $Q_t = Q(\cdot; \rho_t)$  is updated via the MF-TD in (3.5). For the smooth function  $\sigma$  in the parameterization of the Q function in (3.4), we consider it to be the following two-layer neural network,

$$\sigma(s, a; \theta) = B_\beta \cdot \beta(b) \cdot \tilde{\sigma}(w^\top(s, a, 1)), \quad (4.2)$$

where  $\tilde{\sigma} : \mathbb{R} \rightarrow \mathbb{R}$  is the activation function,  $\theta = (b, w)$  is the parameter, and  $\beta : \mathbb{R} \rightarrow (-1, 1)$  is an odd and invertible function with scaling hyper-parameter  $B_\beta > 0$ . It then holds that  $D = d + 2$ , where  $d$  and  $D$  are the dimensions of  $(s, a)$  and  $\theta$ , respectively. It is worth noting that the function class of  $\int \sigma(s, a; \theta) \rho(d\theta)$  for  $\rho \in \mathcal{P}_2(\mathbb{R}^D)$  is the same as

$$\mathcal{F} = \left\{ \int \beta' \cdot \tilde{\sigma}(w^\top(s, a, 1)) \nu(d\beta', dw) \mid \nu \in \mathcal{P}_2((-B_\beta, B_\beta) \times \mathbb{R}^{d+1}) \right\}, \quad (4.3)$$

which captures a vast function class because of the universal function approximation theorem [8, 47]. We remark that we introduce the rescaling function  $\beta$  in (4.2) to avoid the study of the space of probability measures over  $(-B_\beta, B_\beta) \times \mathbb{R}^{d+1}$  in (4.3), which has boundary and thus lacks the regularity in the study of optimal transport. Furthermore, note that we introduce a hyper-parameter  $\alpha > 1$  in the Q function in (3.4). Thus, we are using  $\alpha \cdot \mathcal{F}$  to represent  $\mathcal{F}$ , which causes an “over-representation” when  $\alpha > 1$ . Such over-representation appears to be essential for our analysis. For a brief peek, we remark that  $\alpha$  actually controls the gap in the average total reward over time when the relative time-scale  $\eta$  is properly selected according to Theorem 4.6. Furthermore, such an influence is imposed through Lemma 4.4, which shows that the Wasserstein distance between  $\rho_0$  and  $\rho_{\pi_t}$  is upper bounded by  $O(1/\alpha)$ . In what follows, we consider the initialization of the TD update to be  $\rho_0 = \mathcal{N}(0, I_D)$ , which implies that  $Q(s, a; \rho_0) = 0$ . We next impose the following assumption on the two-layer neural network  $\sigma$ .

**Assumption 4.2** (Regularity of the Neural Network). For the two-layer neural network  $\sigma$  defined in (4.2), we assume that the following properties hold.

- (i) The rescaling function  $\beta : \mathbb{R} \rightarrow (-1, 1)$  is odd,  $L_{0,\beta}$ -Lipschitz continuous,  $L_{1,\beta}$ -smooth, and invertible. Meanwhile, the inverse  $\beta^{-1}$  is locally Lipschitz continuous. In particular, we assume that  $\beta^{-1}$  is  $\ell_\beta$ -Lipschitz continuous in  $[-2/3, 2/3]$ .
- (ii) The activation function  $\tilde{\sigma} : \mathbb{R} \rightarrow \mathbb{R}$  is odd,  $B_{\tilde{\sigma}}$ -bounded,  $L_{0,\tilde{\sigma}}$ -Lipschitz continuous, and  $L_{1,\tilde{\sigma}}$ -smooth.

We remark that Assumption 4.2 is not restrictive and can be satisfied by a large family of neural networks, e.g.,  $\tilde{\sigma}(x) = \tanh(x)$  and  $\beta(b) = \tanh(b)$ . Noting that  $\|(s, a)\|_2 \leq 1$ , Assumption 4.2 implies that the function  $\sigma(s, a; \theta)$  in (4.2) is odd with respect to  $w$  and  $b$  and is also bounded, Lipschitz continuous, and smooth in the parameter domain, that is,

$$|\nabla_\theta \sigma(s, a; \theta)| < B_1, \quad |\nabla_{\theta\theta}^2 \sigma(s, a; \theta)| < B_2. \quad (4.4)$$

We then impose the following assumption on the MDP.

**Assumption 4.3** (Regularity of the MDP). For the MDP  $(\mathcal{S}, \mathcal{A}, \gamma, P, r, \mathcal{D}_0)$ , we assume the following properties hold.

- (i) The reward function  $r$  and the transition kernel  $P$  admit the following representations with respect to the activation function  $\tilde{\sigma}$ ,

$$r(s, a) = B_r \cdot \int \tilde{\sigma}((s, a, 1)^\top w) \mu(dw), \quad (4.5)$$

$$P(s' | s, a) = \int \tilde{\sigma}((s, a, 1)^\top w) \varphi(s') \psi(s'; dw), \quad (4.6)$$

where  $\mu$  and  $\psi(s'; \cdot)$  are probability measures in  $\mathcal{P}_2(\mathbb{R}^{d+1})$  for any  $s' \in \mathcal{S}$ ,  $B_r$  is a positive scaling parameter, and  $\varphi(s') : \mathcal{S} \rightarrow \mathbb{R}_+$  is a nonnegative function.

- (ii) The reward function  $r$  satisfies that  $r(s, a) \geq 0$  for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . For the representation of  $r$  in (4.5) and the representation of the transition kernel  $P$  in (4.6), we assume that

$$\begin{aligned} \chi^2(\mu \| \rho_{w,0}) &< M_\mu, & \chi^2(\psi(s; \cdot) \| \rho_{w,0}) &< M_\psi, & \forall s \in \mathcal{S}, \\ \int \varphi(s) ds &\leq M_{1,\varphi}, & \int \varphi(s)^2 ds &\leq M_{2,\varphi}, \end{aligned}$$

where  $\rho_{w,0}$  is the marginal distribution of  $\rho_0$  with respect to  $w$ , i.e.,  $\rho_{w,0} = \int \rho_0(db, \cdot)$ ,  $\chi^2$  is the chi-squared divergence, and  $M_\mu, M_\psi, M_{1,\varphi}, M_{2,\varphi}$  are absolute constants.

- (iii) We assume that there exists an absolute constant  $\mathcal{G}$  such that

$$\begin{aligned} \|\psi(s; \cdot) - \psi(s'; \cdot)\|_{\dot{H}^{-1}(\mu)} &< \mathcal{G}, & \|\psi(s; \cdot) - \mu\|_{\dot{H}^{-1}(\mu)} &< \mathcal{G}, \\ \|\psi(s; \cdot) - \mu\|_{\dot{H}^{-1}(\psi(s'; \cdot))} &< \mathcal{G}, & \|\psi(s; \cdot) - \psi(s'; \cdot)\|_{\dot{H}^{-1}(\psi(s''; \cdot))} &< \mathcal{G}, \end{aligned} \quad \forall s, s', s'' \in \mathcal{S},$$

where  $\|\cdot\|_{\dot{H}^{-1}(\cdot)}$  is the weighted homogeneous Sobolev norm.

We remark that by assuming  $\psi$  to be a probability measure and that  $\varphi(s') \geq 0$  in (4.6), the representation of the transition kernel does not lose generality. Specifically, the function class of (4.6) is the same as

$$\mathcal{P} = \left\{ \int \tilde{\sigma}((s, a, 1)^\top w) \tilde{\psi}(s'; dw) \mid \tilde{\psi}(s'; \cdot) \text{ is a signed measure for any } s' \in \mathcal{S} \right\}.$$

See §C.1 for a detailed proof. Assumption 4.3 generalizes the linear MDP in [14, 32, 67, 68]. In contrast, our representation of the reward function and the transition kernel benefits from the universal function approximation theorem and is thus not as restrictive as the original linear MDP assumption. Note that the infinite-width neural network has a two-layer structure by (4.2). We establish the following lemma on the regularity of the representation of the action value function  $Q^\pi$  by such a neural network.

**Lemma 4.4** (Regularity of Representation of  $Q^\pi$ ). Suppose that Assumptions 4.2 and 4.3 hold. For any policy  $\pi$ , there exists a probability measure  $\rho_\pi \in \mathcal{P}_2(\mathbb{R}^D)$  for the representation of  $Q^\pi$  with the following properties.

- (i) For function  $Q(s, a; \rho_\pi)$  defined by (3.4) with  $\rho = \rho_\pi$  and the action value function  $Q^\pi(s, a)$  defined by (2.2), we have  $Q(s, a; \rho_\pi) = Q^\pi(s, a)$  for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .
- (ii) By letting  $B_\beta \geq 2(B_r + \gamma(1 - \gamma)^{-1} B_r M_{1,\varphi})$  for the neural network defined in (4.2) and  $\rho_0 \sim \mathcal{N}(0, I_D)$  for the initial distribution, we have  $\widetilde{W}_2(\rho_\pi, \rho_0) \leq \bar{D}$  for any policy  $\pi$ , where we define  $\widetilde{W}_2(\cdot, \cdot) = \alpha W_2(\cdot, \cdot)$  as the scaled  $W_2$  metric. Here constant  $\bar{D}$  depends on the discount factor  $\gamma$  and the absolute constants  $L_{0,\beta}, L_{1,\beta}, l_\beta, B_r, M_\mu, M_\psi, M_{1,\varphi}, M_{2,\varphi}$  defined in Assumptions 4.2 and 4.3.

*Proof.* See §B.2 for a detailed proof. □

Property (i) of Lemma 4.4 shows that the action value function  $Q^\pi$  can be parameterized with the infinite-width two-layer neural network  $Q(\cdot; \rho_\pi)$  in (3.4). Note that a larger  $B_\beta$  captures a larger function class in (4.3). Without loss of generality, we consider that  $B_\beta \geq 2(B_r + \gamma(1-\gamma)^{-1}B_r M_{1,\varphi})$  holds in the sequel. Hence, by Property (ii), it holds that  $\widetilde{W}_2(\rho_\pi, \rho_0) \leq O(1)$  for any policy  $\pi$ . In particular, it holds by Property (i) of Lemma 4.4 that  $\|Q_t - Q^{\pi_t}\|_{2, \widetilde{\phi}^{\pi_t}} = \|Q(\cdot; \rho_t) - Q(\cdot; \rho_{\pi_t})\|_{2, \widetilde{\phi}^{\pi_t}}$  and we have the following theorem to characterize such an error with regard to the  $W_2$  space.

**Theorem 4.5** (Upper Bound of Policy Evaluation Error). Suppose that Assumptions 4.2 and 4.3 hold and  $\rho_0 \sim \mathcal{N}(0, I_D)$  is the initial distribution. We specify the weighting distribution  $\widetilde{\Phi}^{\pi_t}$  in MF-TD (3.5) as  $\widetilde{\Phi}^{\pi_t} = \widetilde{\mathcal{E}}_{\widetilde{\phi}^{\pi_t}}^{\pi_t}$ , where  $\widetilde{\phi}^{\pi_t} \in \mathcal{P}(\mathcal{S} \times \mathcal{A})$  is the evaluation distribution for the policy evaluation error in Theorem 4.1. Then, it holds that

$$(1 - \sqrt{\gamma}) \cdot \|Q_t - Q^{\pi_t}\|_{2, \widetilde{\phi}^{\pi_t}}^2 \leq -\frac{d}{dt} \frac{\widetilde{W}_2^2(\rho_t, \rho_{\pi_t})}{2\eta} + \Delta_t, \quad (4.7)$$

where

$$\begin{aligned} \Delta_t = & 2\alpha^{1/2}\eta^{-1}\mathcal{B}B_1 \cdot \widetilde{W}_2(\rho_t, \rho_0)\widetilde{W}_2(\rho_t, \rho_{\pi_t}) \\ & + \alpha^{-1}B_2 \cdot \left(4B_1 \max\{\widetilde{W}_2(\rho_{\pi_t}, \rho_0), \widetilde{W}_2(\rho_t, \rho_0)\} + B_r\right) \widetilde{W}_2(\rho_t, \rho_{\pi_t})^2. \end{aligned}$$

Here  $B_1$  and  $B_2$  are defined in (4.4) of Assumption 4.2,  $\eta$  is the relative TD timescale,  $\alpha$  is the scaling parameter of the neural network, and  $\widetilde{W}_2 = \alpha W_2$  is the scaled  $W_2$  metric. Moreover, constant  $\mathcal{B}$  depends on the discount factor  $\gamma$ , the scaling parameter  $B_\beta$  in (4.2), and the absolute constants  $l_\beta, B_r, M_{1,\varphi}, \mathcal{G}$  defined in Assumptions 4.2 and 4.3.

*Proof.* See §B.3 for a detailed proof.  $\square$

Here we give a nonrigorous discussion on how to upper bound  $\Delta_t$  in (4.7). If  $\widetilde{W}_2(\rho_t, \rho_0) \leq O(1)$  holds for any  $t \in [0, T]$ , by  $\widetilde{W}_2(\rho_{\pi_t}, \rho_0) \leq O(1)$  in Lemma 4.4 and the triangle inequality of  $W_2$  distance [59], it follows that  $\widetilde{W}_2(\rho_t, \rho_{\pi_t}) \leq O(1)$  and  $\Delta_t \leq O(\alpha^{1/2}\eta^{-1} + \alpha^{-1})$ . Taking a time average of integration on both sides of (4.7), the policy evaluation error  $\frac{1}{T} \int_0^T \|Q_t - Q^{\pi_t}\|_{2, \widetilde{\phi}^{\pi_t}} dt$  is then upper bounded by  $O(\eta^{-1}T^{-1} + \alpha^{1/2}\eta^{-1} + \alpha^{-1})$ . Inspired by such a fact, we introduce the following restarting mechanism to ensure  $\widetilde{W}_2(\rho_t, \rho_0) \leq O(1)$ .

**Restarting Mechanism.** Let  $\widetilde{W}_0 = \lambda\bar{D}$  be a threshold, where  $\bar{D}$  is the upper bound for  $\widetilde{W}_2(\rho_\pi, \rho_0)$  by Lemma 4.4,  $\lambda \geq 3$  is a constant scaling parameter for the restarting threshold,  $\rho_t$  is the distribution of the parameters in the neural network at time  $t$ , and  $\rho_0$  is the initial distribution. Whenever we detect that  $\widetilde{W}_2(\rho_t, \rho_0)$  reaches  $\widetilde{W}_0$  in the update, we pause and reset  $\rho_t$  to  $\rho_0$  by resampling the parameters from  $\rho_0$ . Then, we reset the critic with the newly sampled parameters while keeping the actor's policy  $\pi_t$  unchanged and continue the update.

The restarting mechanism guarantees  $\widetilde{W}_2(\rho_t, \rho_0) \leq \lambda\bar{D}$  by restricting the distribution  $\rho_t$  of the parameters to be close to  $\rho_0$ . Moreover, by letting  $\lambda \geq 3$ , we ensure that  $\rho_{\pi_t}$  is realizable by  $\rho_t$  since  $\widetilde{W}_2(\rho_{\pi_t}, \rho_0) \leq \bar{D} \leq \lambda\bar{D}$ , which means that the neural network is capable of capturing the representation of the action value function  $Q^{\pi_t}$ . We remark that by letting  $\widetilde{W}_0 = O(1)$ , we allow  $\rho_t$  to deviate from  $\rho_0$  up to  $W_2(\rho_t, \rho_0) \leq O(\alpha^{-1})$  in the restarting mechanism. In contrast, the NTK regime [15] which corresponds to letting  $\alpha = \sqrt{M}$  in (3.1) only allows  $\rho_t$  to deviate from  $\rho_0$  by the chi-squared divergence  $\chi^2(\rho_t \| \rho_0) \leq O(M^{-1}) = o(1)$ . That is, the NTK regime fails to induce a feature representation significantly different from the initial one. Before moving on, we summarize the construction of the weighting distribution  $\widetilde{\Phi}^{\pi_t}$  in Theorem 4.1 and 4.5 as follows,

$$\widetilde{\Phi}^{\pi_t} = \widetilde{\mathcal{E}}_{\widetilde{\phi}^{\pi_t}}^{\pi_t}, \quad \widetilde{\phi}^{\pi_t} = \frac{1}{2}\widetilde{\phi}_0 + \frac{1}{2}\phi_0 \otimes \pi_t, \quad \phi_0 = \int_{\mathcal{A}} \widetilde{\phi}_0(\cdot, da), \quad (4.8)$$

where  $\widetilde{\phi}_0$  is the base distribution. Now we have the following theorem that characterizes the global optimality and convergence of the two-timescale AC with restarting mechanism.

**Theorem 4.6** (Global Optimality and Convergence Rate of Two-timescale AC with Restarting Mechanism). Suppose that (4.8) and Assumptions 4.2 and 4.3 hold. With the restarting mechanism, it holds that

$$\frac{1}{T} \int_0^T (J(\pi^*) - J(\pi_t)) dt \leq \underbrace{\frac{\zeta}{T}}_{(a)} + 4\kappa \underbrace{\sqrt{\alpha^{-1}S_1 + \alpha^{1/2}\eta^{-1}S_2 + \frac{\eta^{-1}\bar{D}^2}{2T(1-\sqrt{\gamma})}}}_{(b)}, \quad (4.9)$$

where we have

$$\zeta = \mathbb{E}_{s \sim \mathcal{E}_{\mathcal{D}_0}^*} \left[ \text{KL}(\pi^*(\cdot | s) \| \pi_0(\cdot | s)) \right], \quad \kappa = \left\| \frac{\tilde{\mathcal{E}}_{\mathcal{D}_0}^*}{\tilde{\phi}_0} \right\|_{\infty},$$

$$S_1 = \frac{(1 + \lambda)^2 \bar{D}^2 B_2 (4B_1 \lambda \bar{D} + B_r)}{1 - \sqrt{\gamma}}, \quad S_2 = \frac{2\mathcal{B}B_1 \lambda (1 + \lambda) \bar{D}^2}{1 - \sqrt{\gamma}}.$$

Here  $B_r$ ,  $B_1$  and  $B_2$  are defined in Assumption 4.2 and 4.3,  $\bar{D}$  is the upper bound for  $\tilde{W}_2(\rho_\pi, \rho_0)$  in Lemma 4.4,  $\mathcal{B}$  depends on the discount factor  $\gamma$  and the absolute constants defined in Assumption 4.2 and 4.3, and  $\lambda$  is the scaling parameter for the restarting threshold. Besides, it holds for the total restarting number  $N$  that

$$N \leq (\lambda - 2)^{-1} ((\alpha^{-1}\eta S_1 + \alpha^{1/2}S_2)2T\bar{D}^{-2}(1 - \sqrt{\gamma}) + 1).$$

*Proof.* See §B.4 for a detailed proof.  $\square$

Note that for a given MDP with starting distribution  $\mathcal{D}_0$ , the expected KL-divergence  $\zeta$  and the concentrability coefficient  $\kappa$  are both independent of the two-timescale update. We remark that our condition for (4.9) to be bounded is not restrictive. Specifically, we only need a given  $\pi_0$  and  $\tilde{\phi}_0$  such that the KL-divergence  $\zeta < \infty$  and the concentrability coefficient  $\kappa < \infty$ , which is weaker than the concentrability coefficient used in [7, 23, 24, 38, 39, 43, 48, 57, 61].

The first term (a) on the right-hand side of (4.9) diminishes as  $T \rightarrow \infty$ . The second term (b) corresponds to the policy evaluation error. We give an example to demonstrate the convergence of the two-time AC. We let the scaling parameter  $\lambda = 3$  for the restarting threshold. By letting  $\eta = \alpha^{3/2}$ , it holds that (b) =  $O(\alpha^{-1/2})$  as  $\alpha \rightarrow \infty$ . Thus, we have that  $\frac{1}{T} \int_0^T (J(\pi^*) - J(\pi_t)) dt$  descends at a rate of  $O(T^{-1} + O(\alpha^{-1/2}) + O(\alpha^{-3/4}T^{-1/2}))$ . Note that  $\eta = \alpha^{3/2}$  shows that the critic has a larger relative TD timescale in (3.5). As for the total number of restartings  $N$ , it holds that  $N \leq O(\alpha^{1/2}T)$  as  $\alpha \rightarrow \infty$ , which induces a tradeoff, i.e., a larger  $\alpha$  guarantees a smaller gap in  $\frac{1}{T} \int_0^T (J(\pi^*) - J(\pi_t)) dt$  but yields in more restartings and a larger relative TD timescale.

## 5 Acknowledgement

Zhaoran Wang acknowledges National Science Foundation (Awards 2048075, 2008827, 2015568, 1934931), Simons Institute (Theory of Reinforcement Learning), Amazon, J.P. Morgan, and Two Sigma for their supports. Zhuoran Yang acknowledges Simons Institute (Theory of Reinforcement Learning). This work was also supported in part by the Vannevar Bush Faculty Fellowship program under grant number N00014-21-1-2941.

## References

- [1] Agazzi, A. and Lu, J. (2019). Temporal-difference learning for nonlinear value function approximation in the lazy training regime. *arXiv preprint arXiv:1905.10917*.
- [2] Agazzi, A. and Lu, J. (2020). Global optimality of softmax policy gradient with single hidden layer neural networks in the mean-field regime. *arXiv preprint arXiv:2010.11858*.

- [3] Agostinelli, F., McAleer, S., Shmakov, A. and Baldi, P. (2019). Solving the rubik’s cube with deep reinforcement learning and search. *Nature Machine Intelligence*, **1** 356–363.
- [4] Akkaya, I., Andrychowicz, M., Chociej, M., Litwin, M., McGrew, B., Petron, A., Paino, A., Plappert, M., Powell, G., Ribas, R. et al. (2019). Solving rubik’s cube with a robot hand. *arXiv preprint arXiv:1910.07113*.
- [5] Ambrosio, L. and Gigli, N. (2013). A user’s guide to optimal transport. In *Modelling and Optimisation of Flows on Networks*. Springer, 1–155.
- [6] Ambrosio, L., Gigli, N. and Savaré, G. (2008). *Gradient flows: In metric spaces and in the space of probability measures*. Springer.
- [7] Antos, A., Szepesvári, C. and Munos, R. (2008). Learning near-optimal policies with Bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, **71** 89–129.
- [8] Barron, A. R. (1993). Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information Theory*, **39** 930–945.
- [9] Berner, C., Brockman, G., Chan, B., Cheung, V., Debiak, P., Dennison, C., Farhi, D., Fischer, Q., Hashme, S., Hesse, C. et al. (2019). Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*.
- [10] Bhatnagar, S., Ghavamzadeh, M., Lee, M. and Sutton, R. S. (2008). Incremental natural actor-critic algorithms. In *Advances in Neural Information Processing Systems*.
- [11] Bhatnagar, S., Sutton, R. S., Ghavamzadeh, M. and Lee, M. (2009). Natural actor-critic algorithms. *Automatica*, **45** 2471–2482.
- [12] Börgers, T. and Sarin, R. (1997). Learning through reinforcement and replicator dynamics. *Journal of economic theory*, **77** 1–14.
- [13] Borkar, V. S. (2009). *Stochastic approximation: A dynamical systems viewpoint*. Springer.
- [14] Cai, Q., Yang, Z., Jin, C. and Wang, Z. (2019). Provably efficient exploration in policy optimization. *arXiv preprint arXiv:1912.05830*.
- [15] Cai, Q., Yang, Z., Lee, J. D. and Wang, Z. (2019). Neural temporal-difference learning converges to global optima. In *Advances in Neural Information Processing Systems*.
- [16] Chen, M., Wang, Y., Liu, T., Yang, Z., Li, X., Wang, Z. and Zhao, T. (2020). On computation and generalization of generative adversarial imitation learning. *arXiv preprint arXiv:2001.02792*.
- [17] Chen, Z., Cao, Y., Gu, Q. and Zhang, T. (2020). Mean-field analysis of two-layer neural networks: Non-asymptotic rates and generalization bounds. *arXiv preprint arXiv:2002.04026*.
- [18] Chizat, L. and Bach, F. (2018). A note on lazy training in supervised differentiable programming. *arXiv preprint arXiv:1812.07956*.
- [19] Chizat, L. and Bach, F. (2018). On the global convergence of gradient descent for over-parameterized models using optimal transport. In *Advances in Neural Information Processing Systems*.
- [20] Duan, Y., Chen, X., Houthoofd, R., Schulman, J. and Abbeel, P. (2016). Benchmarking deep reinforcement learning for continuous control. In *International Conference on Machine Learning*.
- [21] Fang, C., Dong, H. and Zhang, T. (2019). Over parameterized two-level neural networks can learn near optimal feature representations. *arXiv preprint arXiv:1910.11508*.

- [22] Fang, C., Gu, Y., Zhang, W. and Zhang, T. (2019). Convex formulation of overparameterized deep neural networks. *arXiv preprint arXiv:1911.07626*.
- [23] Farahmand, A.-m., Ghavamzadeh, M., Szepesvári, C. and Mannor, S. (2016). Regularized policy iteration with nonparametric function spaces. *The Journal of Machine Learning Research*, **17** 4809–4874.
- [24] Farahmand, A.-m., Szepesvári, C. and Munos, R. (2010). Error propagation for approximate policy and value iteration. In *Advances in Neural Information Processing Systems*.
- [25] Friedrichs, K. O. (1944). The identity of weak and strong extensions of differential operators. *Transactions of the American Mathematical Society*, **55** 132–151.
- [26] Fu, Z., Yang, Z. and Wang, Z. (2020). Single-timescale actor-critic provably finds globally optimal policy. *arXiv preprint arXiv:2008.00483*.
- [27] Harker, P. T. and Pang, J.-S. (1990). Finite-dimensional variational inequality and nonlinear complementarity problems: A survey of theory, algorithms and applications. *Mathematical Programming*, **48** 161–220.
- [28] Hennes, D., Morrill, D., Omidshafiei, S., Munos, R., Perolat, J., Lanctot, M., Gruslys, A., Lespiau, J.-B., Parmas, P., Duéñez-Guzmán, E. et al. (2020). Neural replicator dynamics: Multi-agent learning via hedging policy gradients. In *International Conference on Autonomous Agents and MultiAgent Systems*.
- [29] Hong, M., Wai, H.-T., Wang, Z. and Yang, Z. (2020). A two-timescale framework for bilevel optimization: Complexity analysis and application to actor-critic. *arXiv preprint arXiv:2007.05170*.
- [30] Jacot, A., Gabriel, F. and Hongler, C. (2018). Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*.
- [31] Javanmard, A., Mondelli, M. and Montanari, A. (2019). Analysis of a two-layer neural network via displacement convexity. *arXiv preprint arXiv:1901.01375*.
- [32] Jin, C., Yang, Z., Wang, Z. and Jordan, M. I. (2020). Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*.
- [33] Kakade, S. and Langford, J. (2002). Approximately optimal approximate reinforcement learning. In *International Conference on Machine Learning*, vol. 2.
- [34] Kakade, S. M. (2002). A natural policy gradient. In *Advances in Neural Information Processing Systems*.
- [35] Khodadadian, S., Doan, T. T., Maguluri, S. T. and Romberg, J. (2021). Finite sample analysis of two-time-scale natural actor-critic algorithm. *arXiv preprint arXiv:2101.10506*.
- [36] Konda, V. R. and Tsitsiklis, J. N. (2000). Actor-critic algorithms. In *Advances in Neural Information Processing Systems*.
- [37] Kushner, H. and Yin, G. G. (2003). *Stochastic approximation and recursive algorithms and applications*. Springer.
- [38] Lazaric, A., Ghavamzadeh, M. and Munos, R. (2016). Analysis of classification-based policy iteration algorithms. *The Journal of Machine Learning Research*, **17** 583–612.
- [39] Liu, B., Cai, Q., Yang, Z. and Wang, Z. (2019). Neural proximal/trust region policy optimization attains globally optimal policy. *arXiv preprint arXiv:1906.10306*.

- [40] Lu, Y., Ma, C., Lu, Y., Lu, J. and Ying, L. (2020). A mean field analysis of deep resnet and beyond: Towards provably optimization via overparameterization from depth. In *International Conference on Machine Learning*. PMLR.
- [41] Mei, S., Misiakiewicz, T. and Montanari, A. (2019). Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. In *Conference on Learning Theory*. PMLR.
- [42] Mei, S., Montanari, A. and Nguyen, P.-M. (2018). A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, **115** E7665–E7671.
- [43] Munos, R. and Szepesvári, C. (2008). Finite-time bounds for fitted value iteration. *The Journal of Machine Learning Research*, **9** 815–857.
- [44] Otto, F. and Villani, C. (2000). Generalization of an inequality by Talagrand and links with the logarithmic Sobolev inequality. *Journal of Functional Analysis*, **173** 361–400.
- [45] Peters, J. and Schaal, S. (2008). Natural actor-critic. *Neurocomputing*, **71** 1180–1190.
- [46] Peyre, R. (2011). Comparison between  $w_2$  distance and  $\dot{H}^{-1}$  norm, and localisation of wasserstein distance. *arXiv preprint arXiv:1104.4631*.
- [47] Pinkus, A. (1999). Approximation theory of the MLP model in neural networks. *Acta Numerica*, **8** 143–195.
- [48] Scherrer, B., Ghavamzadeh, M., Gabillon, V., Lesner, B. and Geist, M. (2015). Approximate modified policy iteration and its application to the game of Tetris. *The Journal of Machine Learning Research*, **16** 1629–1676.
- [49] Schulman, J., Wolski, F., Dhariwal, P., Radford, A. and Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- [50] Schuster, P. and Sigmund, K. (1983). Replicator dynamics. *Journal of theoretical biology*, **100** 533–538.
- [51] Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M. et al. (2016). Mastering the game of Go with deep neural networks and tree search. *Nature*, **529** 484–489.
- [52] Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A. et al. (2017). Mastering the game of Go without human knowledge. *Nature*, **550** 354.
- [53] Sirignano, J. and Spiliopoulos, K. (2020). Mean field analysis of neural networks: A central limit theorem. *Stochastic Processes and their Applications*, **130** 1820–1852.
- [54] Sirignano, J. and Spiliopoulos, K. (2020). Mean field analysis of neural networks: A law of large numbers. *SIAM Journal on Applied Mathematics*, **80** 725–752.
- [55] Sutton, R. S. (1988). Learning to predict by the methods of temporal differences. *Machine Learning*, **3** 9–44.
- [56] Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.
- [57] Szepesvári, C. and Munos, R. (2005). Finite time bounds for sampling based fitted value iteration. In *International Conference on Machine Learning*. ACM.
- [58] Villani, C. (2003). *Topics in optimal transportation*. American Mathematical Society.
- [59] Villani, C. (2008). *Optimal transport: Old and new*. Springer.

- [60] Vinyals, O., Babuschkin, I., Czarnecki, W. M., Mathieu, M., Dudzik, A., Chung, J., Choi, D. H., Powell, R., Ewalds, T., Georgiev, P. et al. (2019). Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature*, **575** 350–354.
- [61] Wang, L., Cai, Q., Yang, Z. and Wang, Z. (2019). Neural policy gradient methods: Global optimality and rates of convergence. *arXiv preprint arXiv:1909.01150*.
- [62] Watkins, C. and Dayan, P. (1992). Q-learning. *Machine Learning*, **8** 279–292.
- [63] Wei, C., Lee, J. D., Liu, Q. and Ma, T. (2019). Regularization matters: Generalization and optimization of neural nets vs their induced kernel. In *Advances in Neural Information Processing Systems*.
- [64] Wu, Y., Zhang, W., Xu, P. and Gu, Q. (2020). A finite time analysis of two time-scale actor critic methods. *arXiv preprint arXiv:2005.01350*.
- [65] Xu, T., Wang, Z. and Liang, Y. (2020). Improving sample complexity bounds for actor-critic algorithms. *arXiv preprint arXiv:2004.12956*.
- [66] Xu, T., Wang, Z. and Liang, Y. (2020). Non-asymptotic convergence analysis of two time-scale (natural) actor-critic algorithms. *arXiv preprint arXiv:2005.03557*.
- [67] Yang, L. F. and Wang, M. (2019). Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. *arXiv preprint arXiv:1905.10389*.
- [68] Yang, L. F. and Wang, M. (2019). Sample-optimal parametric q-learning using linearly additive features. *arXiv preprint arXiv:1902.04779*.
- [69] Zhang, Y., Cai, Q., Yang, Z., Chen, Y. and Wang, Z. (2020). Can temporal-difference and q-learning learn representation? a mean-field theory. *arXiv preprint arXiv:2006.04761*.