

FINITE-TIME CONVERGENCE ANALYSIS OF ACTOR-CRITIC WITH EVOLVING REWARD

Anonymous authors

Paper under double-blind review

ABSTRACT

Many popular practical reinforcement learning (RL) algorithms employ evolving reward functions—through techniques such as reward shaping, entropy regularization, or curriculum learning—yet their theoretical foundations remain underdeveloped. This paper provides the first finite-time convergence analysis of a single-timescale actor-critic algorithm in the presence of an evolving reward function under Markovian sampling. We consider a setting where the reward parameters may change at each time step, affecting both policy optimization and value estimation. Under standard assumptions, we derive non-asymptotic bounds for both actor and critic errors. Our result shows that an $O(1/\sqrt{T})$ convergence rate is achievable, matching the best-known rate for static rewards, provided the reward parameters evolve slowly enough. This rate is preserved when the reward is updated via a gradient-based rule with bounded gradient and on the same timescale as the actor and critic, offering a theoretical foundation for many popular RL techniques. As a secondary contribution, we introduce a novel analysis of distribution mismatch under Markovian sampling, improving the best-known rate by a factor of $\log^2 T$ in the static-reward case.

1 INTRODUCTION

Reinforcement Learning (RL, Sutton et al. 1998) has attracted great research interest in the past decades. On the empirical side, a variety of practical algorithms have been proposed and have demonstrated remarkable success in a wide range of real-world scenarios (Mnih et al., 2013; Lillicrap et al., 2016; Ouyang et al., 2022; DeepSeek-AI et al., 2025). On the theoretical side, great efforts have been made to bridge the gap between empirical practice and theoretical foundations, providing rigorous convergence guarantees and solid theoretical understandings to the empirically powerful algorithms (Agarwal et al., 2021; Mei et al., 2020; Kumar et al., 2023; Wu et al., 2020; Olshevsky & Ghahserifard, 2023).

Still, a very common setting adopted by many practical RL algorithms has been largely overlooked by theoretical analyses. RL theory is built upon Markov Decision Processes (MDPs), which typically assume the existence of a static underlying reward function, and the goal is to learn a policy that maximizes the expected cumulative reward. However, when applying RL in many real-world scenarios, a significant impediment is the difficulty of designing a reward function that is both learnable and aligns with the desired task. This challenge has led to the development of techniques that utilize evolving rewards. These include:

- **Reward Shaping:** Adding auxiliary rewards to guide the policy to the desired goal (Ng et al., 1999; Pathak et al., 2017; Burda et al., 2019; Zheng et al., 2018; Hu et al., 2020; Mahankali et al., 2024; Ma et al., 2024). The auxiliary rewards can come from prior knowledge, or be learned in a self-supervised manner during the training process.
- **Entropy or KL Regularization:** Introducing an entropy or KL regularization term to the optimization objective, which is equivalent to modifying the reward according to the current policy (Hao et al., 2017; 2018a;b; Jaques et al., 2019; Stiennon et al., 2020). The regularization factor can be automatically adjusted during training.
- **Curriculum Learning:** Starting with easier tasks (and their associated rewards) and gradually increasing the difficulty (Narvekar et al., 2020).

Intuitively, a slight change of the reward function does not drastically alter the solution of the underlying MDP, allowing an RL algorithm to remain effective. However, when this change is negligible compared to the under-training policy or value function, the algorithm’s effectiveness becomes questionable, as they are closely interconnected. Therefore, to rigorously support the use of these evolving reward techniques, we must answer the following fundamental question precisely:

How fast can the reward change while still ensuring the convergence of an RL algorithm?

This paper aims to establish a theoretical foundation for this setting by providing the first finite-time convergence analysis for an actor-critic algorithm with an evolving reward. We focus on a single-sample, single-timescale actor-critic algorithm with linear function approximation for the critic under Markovian sampling. This setting is particularly practical yet challenging, as the non-stationarity from the evolving reward affects both the policy gradient (actor) and the value estimation (critic), creating a complex feedback loop.

Specifically, we bound both the expected actor error and the expected critic error in terms of the number of iterations T and the total change of the reward parameters. From this, we derive conditions on the evolving rate of reward parameters necessary to achieve asymptotic convergence and to maintain the $O(1/\sqrt{T})$ convergence rate as in the static-reward case. Moreover, it turns out that $O(1/\sqrt{T})$ convergence can be achieved if the reward parameter follows a gradient-based update with bounded gradient and the same timescale as the actor and the critic, providing theoretical guarantees to a wide range of practical techniques, including curiosity-driven reward shaping (Pathak et al., 2017), random network distillation methods (Burda et al., 2019), and soft actor-critic with automated entropy adjustment (Haarnoja et al., 2018b).

To handle the evolving reward, we exploit the Lipschitz continuity of the objective function and the optimal critic parameter with respect to the reward parameter, which relies on the Lipschitz continuity assumption for the reward itself. In addition, we introduce a novel analysis on the distribution mismatch caused by Markovian sampling, improving the convergence rate by a factor of $\log^2 T$ in the static-reward case.

In summary, our contributions include the following:

- **Important Problem Formulation:** We formalize the problem of Actor-Critic with Evolving Reward, where the reward parameter φ_t (encompassing the true reward and regularization terms) can be updated by an arbitrary oracle at every time step.
- **Novel Non-Asymptotic Results:** Under standard assumptions (Linear function approximated critic, Lipschitz continuity of policy and reward, sufficient exploration), we derive the convergence rate of the single-sample single-timescale actor-critic algorithm with Markovian sampling, and show that it achieves a convergence rate of $O(1/\sqrt{T})$ to a neighborhood of a stationary point under mild condition on the evolving reward, validating a wide range of practical RL techniques.
- **Interesting Technical Tools:** We establish the necessary assumptions and key properties to analyze the effects of the evolving reward on standard actor-critic algorithms. We also provide a novel analysis of the distribution mismatch caused by Markovian sampling, which independently improves the convergence rate for the static-reward case.

2 RELATED WORK

Our work is developed upon the finite-time analysis of policy gradient methods and actor-critic methods.

Finite-time analysis of Policy Gradient Methods. The finite-time convergence guarantees of policy gradient methods have been studied by Agarwal et al. (2021); Mei et al. (2020); Xiao (2022), assuming access to exact gradient oracles. For the stochastic case where the algorithm can only access gradient estimators from sampled trajectories or transitions, convergence results in terms of sample complexity have been established by Liu et al. (2020); Ding et al. (2022; 2025); Fatkhullin et al. (2023); Mondal & Aggarwal (2024).

Finite-time analysis of Actor-Critic Methods. The finite-time analysis of actor-critic methods encompasses several variants of the algorithm, including the double-loop setting (Yang et al., 2019; Kumar et al., 2023; Xu et al., 2020a; Cayci et al., 2024; Gaur et al., 2024; Ganesh et al., 2025), the two-timescale setting (Wu et al., 2020; Xu et al., 2020b; Shen et al., 2023; Hong et al., 2023), and the single-timescale setting (Chen et al., 2021; Olshevsky & Ghahserifard, 2023; Tian et al., 2023; Chen & Zhao, 2023; 2025). The analysis of single-timescale actor-critic methods is particularly relevant to our work. Both Chen et al. (2021) and Olshevsky & Ghahserifard (2023) obtain an $O(1/\sqrt{T})$ convergence rate in discrete spaces, assuming i.i.d. sampling. Chen & Zhao (2023) and Tian et al. (2023) made efforts to tackle the more practical yet challenging Markovian sampling setting. However, the former considers the average-reward scenario instead of the commonly used discounted-reward scenario, while the latter employs Markovian samples for the critic and i.i.d. samples for the actor. Chen & Zhao (2025) ultimately resolves the problem of Markovian sampling, obtaining an $O(\log^2 T/\sqrt{T})$ convergence rate, and further extends the analysis to continuous action spaces. Additionally, Tian et al. (2023) incorporates an analysis of neural network approximated critics, while the other four assume a linear function approximated critic.

While we focus on analyzing the convergence of the single-timescale actor-critic algorithm under evolving rewards, we recognize two related research areas that focus on designing reinforcement learning (RL) algorithms for non-static Markov Decision Processes (MDPs), where both rewards and transitions can change:

Adversarial RL and Non-stationary RL. Tailored for online learning, adversarial RL (Even-Dar et al., 2009; Zimin & Neu, 2013; Jin et al., 2020) and non-stationary RL (Cheung et al., 2023; Feng et al., 2023; Mao et al., 2025) aim to minimize cumulative static or dynamic regret through strategic algorithm design in the presence of adversarial rewards and transitions. However, their formulation does not align with the practical scenarios we consider, where the RL policy is trained before executing, and the evolving reward techniques we previously mentioned focus on enhancing convergence during training and improving performance in execution, rather than adapting in an adversarial manner.

Performative RL. Performative RL (Mandal et al., 2023; Rank et al., 2024; Mandal & Radanovic, 2025) examines situations in which the underlying MDP changes in response to the deployed policy. This setting partially overlaps with ours, as performative rewards represent a specific case within our general evolving reward framework. Under certain mild assumptions regarding the performative rewards and transitions, these studies demonstrate the existence of a stable policy and provide finite-time convergence guarantees for their proposed algorithms, which involve iterating between deployment and retraining. However, their approach necessitates finding exact solutions for a saddle point at each retraining step, which is often infeasible in practice.

3 PRELIMINARIES

3.1 MARKOV DECISION PROCESS

We consider an infinite-horizon discounted Markov Decision Process (MDP), defined by the tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$. Here, \mathcal{S} is the state space, \mathcal{A} is the action space, $\mathcal{P} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is the transition kernel, $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, and $\gamma \in (0, 1)$ is the discount factor.

A policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ maps states to distributions over actions. Starting from an initial state s_0 , at each time step $t = 0, 1, \dots$, an action $a_t \sim \pi(\cdot|s_t)$ is sampled, yielding a reward $r_t = r(s_t, a_t)$ and transitioning to the next state $s_{t+1} \sim \mathcal{P}(\cdot|s_t, a_t)$. The standard value function $V^\pi(s)$ and action-value function $Q^\pi(s, a)$ represent the expected discounted cumulative reward starting from state s (and action a), respectively, and are defined as

$$V^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_t \sim \pi(\cdot|s_t), s_{t+1} \sim \mathcal{P}(\cdot|s_t, a_t) \right], \quad (1)$$

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a, s_{t+1} \sim \mathcal{P}(\cdot|s_t, a_t), a_{t+1} \sim \pi(\cdot|s_{t+1}) \right]. \quad (2)$$

Entropy Regularized MDPs. Policy optimization often benefits from entropy regularization to encourage exploration (Haarnoja et al., 2017; 2018a;b). The entropy of a policy π at a state s is defined as $\mathcal{H}(\cdot|s) = -\int_{\mathcal{A}} \pi(a|s) \log \pi(a|s) da$. Let

$$H^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{H}(\pi(\cdot|s_t)) \mid s_0 = s, a_t \sim \pi(\cdot|s_t), s_{t+1} \sim \mathcal{P}(\cdot|s_t, a_t) \right], \quad (3)$$

the regularized (soft) value function is defined as

$$\tilde{V}^\pi(s) = V^\pi(s) + \alpha H^\pi(s),$$

where $\alpha \geq 0$ is a hyper-parameter, and $\alpha = 0$ recovers the standard, unregularized case. This formulation is equivalent to solving the original MDP with a regularized reward function:

$$\tilde{r}(s, a) = r(s, a) - \alpha \log \pi(a|s). \quad (4)$$

Consequently, $\tilde{V}(s)$ and $\tilde{Q}(s, a)$ can be interpreted as the value functions under this regularized reward \tilde{r} :

$$\tilde{V}^\pi(s) = \mathbb{E}_{\pi, \mathcal{P}} \left[\sum_{t=0}^{\infty} \gamma^t \tilde{r}(s_t, a_t) \mid s_0 = s \right], \quad \tilde{Q}^\pi(s, a) = \mathbb{E}_{\pi, \mathcal{P}} \left[\sum_{t=0}^{\infty} \gamma^t \tilde{r}(s_t, a_t) \mid s_0 = s, a_0 = a \right].$$

Note on KL Regularization. While we focus on entropy regularization for clarity, our analysis also applies to KL regularization against a fixed reference policy π_{ref} . The corresponding regularized reward becomes $\tilde{r}(s, a) = r(s, a) + \alpha \log \pi_{\text{ref}}(a|s) - \alpha \log \pi(a|s)$. Since the term $\alpha \log \pi_{\text{ref}}(a|s)$ can be absorbed into the base reward $r(s, a)$, we will consider only the entropy regularizer in the following analysis without loss of generality.

The goal of reinforcement learning is to find a parameterized policy π_θ that maximizes the objective:

$$J(\theta) = \int_{\mathcal{S}} \rho(s) \tilde{V}^{\pi_\theta}(s) ds, \quad (5)$$

where θ denotes the policy parameter and $\rho \in \Delta(\mathcal{S})$ is the initial distribution.

3.2 ACTOR-CRITIC METHOD

In order to optimize the policy parameter θ , a popular approach is to compute the gradient of the objective $J(\theta)$ and iteratively adjust θ in the direction of $\nabla J(\theta)$. This approach, named the policy gradient method, is based on the policy gradient theorem (Sutton et al., 1999):

$$\begin{aligned} \nabla_\theta J(\theta) &= \frac{1}{1-\gamma} \mathbb{E}_{s \sim \nu_\rho^{\pi_\theta}(\cdot), a \sim \pi_\theta(\cdot|s)} \left[\tilde{Q}^{\pi_\theta}(s, a) \nabla_\theta \log \pi_\theta(a|s) \right] \\ &= \frac{1}{1-\gamma} \mathbb{E}_{s \sim \nu_\rho^{\pi_\theta}(\cdot), a \sim \pi_\theta(\cdot|s)} \left[\left(\tilde{Q}^{\pi_\theta}(s, a) - \tilde{V}^{\pi_\theta}(s) \right) \nabla_\theta \log \pi_\theta(a|s) \right], \end{aligned} \quad (6)$$

where $\nu_\rho^{\pi_\theta} \in \Delta(\mathcal{S})$ is the discounted visitation distribution defined as

$$\nu_\rho^{\pi_\theta}(s) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s \mid s_0 \sim \rho(\cdot), a_t \sim \pi_\theta(\cdot|s_t), s_{t+1} \sim \mathcal{P}(\cdot|s_t, a_t)).$$

Computing this gradient requires an estimator for the \tilde{V} function or \tilde{Q} function associated with the policy π_θ . A Monte-Carlo estimator (used by REINFORCE (Williams, 1992)) suffers from high variance, resulting in slow convergence. Hence, the actor-critic method (Konda & Tsitsiklis, 1999) introduces another trainable model to approximate the true value functions.

Following the setting of prior work (Chen et al., 2021; Olshevsky & Ghahserifard, 2023; Chen & Zhao, 2023; 2025), we assume that the critic approximates the regularized value function linearly as

$$\hat{V}_\omega(s) = \phi(s)^\top \omega,$$

where $\phi : \mathcal{S} \rightarrow \mathbb{R}^d$ is a known feature mapping satisfying $\|\phi(s)\|_2 \leq 1$ for any state s , and $\omega \in \mathbb{R}^d$ is the trainable parameter. Note that $\tilde{Q}^\pi(s, a) = \tilde{r}(s, a) + \gamma \mathbb{E}_{s'}[\tilde{V}^\pi(s')]$, then with $s' \sim \mathcal{P}(\cdot|s, a)$, the temporal difference (TD) error

$$\hat{\delta}(s, a, s') = r(s, a) - \alpha \log \pi_\theta(a|s) + (\gamma \phi(s') - \phi(s))^\top \omega \quad (7)$$

serves as a biased but low-variance gradient estimator for (6), and the update rule for θ will be

$$\theta_{t+1} \leftarrow \theta_t + \eta_t^\theta \hat{\delta}(s_t, a_t, s'_t) \nabla_\theta \log \pi_\theta(a_t|s_t), \quad (8)$$

where η_t^θ denotes the step size for θ at step t .

Also, we need to align \hat{V}_ω with the true value function \tilde{V}^{π_θ} . As proven by Haarnoja et al. (2017), \tilde{V}^π is the unique solution to the soft Bellman equation $V = \mathcal{T}_\alpha^\pi V$ where the soft Bellman operator \mathcal{T}_α^π is defined as

$$\mathcal{T}_\alpha^\pi V(s) = \mathbb{E}_{a \sim \pi(\cdot|s), s' \sim \mathcal{P}(\cdot|a, s)} [r(s, a) - \alpha \log \pi(a|s) + \gamma V(s')].$$

Hence, a common practice is to adjust the value iteration update $V \leftarrow \mathcal{T}_\alpha^\pi V$ into a stochastic semi-gradient TD(0) update:

$$\omega_{t+1} \leftarrow \omega_t + \eta_t^\omega \hat{\delta}(s_t, a_t, s'_t) \phi(s_t), \quad (9)$$

where η_t^ω denotes the step size for ω at step t .

Markovian Sampling. In our single-sample, single-timescale setting, both actor and critic are updated using the same sample tuple (s_t, a_t, s'_t) at each step. Ideally, s_t shall be sampled from the stationary distribution $\nu_\rho^{\pi_\theta}$, but this is impractical. Instead, we adopt a Markovian sampling scheme (Chen & Zhao, 2025):

$$s_t \sim \hat{\mathcal{P}}(\cdot|s_{t-1}, a_{t-1}), a_t \sim \pi_{\theta_t}(\cdot|s_t), s'_t \sim \mathcal{P}(\cdot|s_t, a_t),$$

where the sampling kernel $\hat{\mathcal{P}}(\cdot|s, a) = \gamma \mathcal{P}(\cdot|s, a) + (1 - \gamma) \rho(\cdot)$ ensures ergodicity. Let $\hat{\nu}_t(\cdot)$ denote the probability distribution of s_t induced by this process. When the policy is fixed, $\hat{\nu}_t$ will converge to $\nu_\rho^{\pi_\theta}$ geometrically. Under a slowly changing policy, the distribution mismatch $\|\hat{\nu}_t - \nu_\rho^{\pi_{\theta_t}}\|_1$ can be controlled by the magnitude of the policy updates. Note that this constitutes an off-policy learning setting for the critic.

3.3 ACTOR-CRITIC WITH EVOLVING REWARD

A central focus of this work is the setting where the regularized reward $\tilde{r}(s, a)$ evolves during training. Depending on the algorithmic design, this evolution can arise from modifications to the base reward $r(s, a)$, the regularization factor α , or the policy π_θ itself. To encompass these variables, we introduce a general reward parameter φ to include all factors that determine $\tilde{r}(s, a)$ along with θ , i.e.,

$$\tilde{r}_{\varphi, \theta}(s, a) = r(s, a; \varphi) - \alpha(\varphi) \log \pi_\theta(a|s).$$

We denote the corresponding soft value function as $\tilde{V}_\varphi^{\pi_\theta}(s)$, and the policy objective as $J_\varphi(\theta)$. The reward parameter φ is updated concurrently with θ and ω at each time step via an arbitrary update rule, which can either be a pre-defined schedule or a feedback-driven strategy. For the critic update, we also introduce a projection Proj_{C_ω} to keep the critic norm bounded by C_ω , which is widely adopted in the literature (Wu et al., 2020; Chen et al., 2021; Olshevsky & Ghahserifard, 2023; Chen & Zhao, 2023; 2025). This framework, summarized in Algorithm 1, unifies a wide range of existing techniques, from automated reward shaping (Martin et al., 2017; Pathak et al., 2017; Burda et al., 2019) to adaptive entropy and KL regularization (Haarnoja et al., 2018b). We provide a further literature review of these evolving reward techniques in Appendix A.

4 MAIN RESULTS

This section presents the finite-time convergence guarantees for the Actor-Critic with Evolving Reward algorithm (Algorithm 1). We begin by stating the standard assumptions required for our analysis, then present the main theorem and a key corollary. Finally, we provide an intuitive proof sketch to elucidate the key technical challenges and innovations.

Algorithm 1 Actor Critic with Evolving Reward

```

270 Initialize:  $\theta_0, \omega_0, \varphi_0, \rho, \{\eta_t^\theta\}_{t \geq 0}, \{\eta_t^\omega\}_{t \geq 0}$ 
271 Sample  $s_0 \sim \rho$ 
272
273 for  $t = 0, 1, \dots, T - 1$  do
274   Sample  $a_t \sim \pi_{\theta_t}(\cdot|s_t), s'_t \sim \mathcal{P}(\cdot|s_t, a_t), s_{t+1} \sim \widehat{\mathcal{P}}(\cdot|s_t, a_t)$ 
275    $\hat{\delta}_t \leftarrow \tilde{r}_{\varphi_t, \theta_t}(s_t, a_t) + (\gamma \phi(s'_t) - \phi(s_t))^\top \omega_t$ 
276    $\theta_{t+1} \leftarrow \theta_t + \eta_t^\theta \hat{\delta}_t \nabla_{\theta} \log \pi_{\theta}(a_t|s_t)$ 
277    $\omega_{t+1} \leftarrow \text{Proj}_{C_\omega} \left( \omega_t + \eta_t^\omega \hat{\delta}_t \phi(s_t) \right)$ 
278    $\varphi_{t+1} \leftarrow \text{UpdateReward}(\varphi_t)$ 
279
280 end for

```

4.1 ASSUMPTIONS

Our analysis relies on several standard assumptions in the literature, which we adapt to accommodate the evolving reward setting.

By taking the expectation of ω_{t+1} conditioning on ω_t in (9) with respect to the discounted visitation distribution, we have

$$\mathbb{E}[\omega_{t+1}|\omega_t] = \omega_t + \eta_t^\omega (\mathbf{b}_{\varphi, \theta} - \mathbf{A}_\theta \omega_t),$$

where

$$\mathbf{A}_\theta = \mathbb{E}_{s \sim \nu_\rho^{\pi_\theta}(\cdot), a \sim \pi_\theta(\cdot|s), s' \sim \mathcal{P}(\cdot|s, a)} [\phi(s)(\phi(s) - \gamma \phi(s'))^\top], \quad (10)$$

$$\mathbf{b}_{\varphi, \theta} = \mathbb{E}_{s \sim \nu_\rho^{\pi_\theta}(\cdot), a \sim \pi_\theta(\cdot|s)} [\tilde{r}_{\varphi, \theta}(s, a) \phi(s)]. \quad (11)$$

It has been shown by Sutton et al. (1998) that the TD limiting point $\omega^*(\varphi, \theta)$ satisfies

$$\mathbf{A}_\theta \omega^*(\varphi, \theta) = \mathbf{b}_{\varphi, \theta}. \quad (12)$$

To ensure the existence of ω^* , we need \mathbf{A}_θ to be non-singular. Fortunately, we can show that \mathbf{A}_θ is positive definite given sufficient exploration over the state space.

Assumption 4.1 (Sufficient Exploration) *Let $\Sigma_\theta = \mathbb{E}_{\nu_\rho^{\pi_\theta}} [\phi(s)\phi(s)^\top]$, then for any $\theta \in \Omega(\theta)$, Σ_θ is positive definite with singular values lower-bounded by $\lambda_\Sigma > 0$.*

Proposition 4.2 *For any $\theta \in \Omega(\theta)$, \mathbf{A}_θ is positive definite with singular values lower-bounded by $\lambda = (1 - \sqrt{\gamma})\lambda_\Sigma$.*

Proposition 4.2 is widely adopted as an assumption in analyzing TD-learning and actor-critic with linear function approximation (Wu et al., 2020; Chen et al., 2021; Olshevsky & Ghahserifard, 2023; Chen & Zhao, 2023; 2025). Although the definition of \mathbf{A}_θ here differs from previous literature where the expectation is taken over the stationary distribution instead of the discounted visitation distribution as in our definition, this nice property can still be induced from the fundamental Assumption 4.1 (Bhandari et al., 2018). Equipped with Proposition 4.2, we can now solve from (12) that $\omega^*(\varphi, \theta) = \mathbf{A}_\theta^{-1} \mathbf{b}_{\varphi, \theta}$, and further imply that $\omega^*(\varphi, \theta)$ is bounded by some constant C_ω since both \mathbf{A}_θ^{-1} and $\mathbf{b}_{\varphi, \theta}$ can be shown bounded, which justifies the projection operator introduced in Algorithm 1.

As ω^* is bounded, \widehat{V}_{ω^*} is bounded, the linear function approximation error has a uniform upper bound, denoted as ϵ . Formally,

$$\epsilon := \sup_{\theta, \varphi} \sqrt{\mathbb{E}_{s \sim \nu_\rho^{\pi_\theta}(\cdot)} \left[\left(\phi(s)^\top \omega^*(\varphi, \theta) - \widehat{V}^{\pi_\theta}(s) \right)^2 \right]}. \quad (13)$$

The error ϵ is zero if $\widehat{V}^{\pi_\theta}(\cdot)$ is indeed a linear function for any φ and θ given the feature mapping $\phi(\cdot)$. To capture the bias of the TD-gradient estimator for the actor, we need the following bound that controls the error of TD-errors (refer to Appendix D for a detailed proof):

Proposition 4.3 For any $\theta \in \Omega(\theta)$, $\varphi \in \Omega(\varphi)$,

$$\sqrt{\mathbb{E}_{\nu^{\pi_\theta, \pi_\theta, \mathcal{P}}} \left[\left(\left(\gamma \widehat{V}_{\omega^*}(s') - \widehat{V}_{\omega^*}(s) \right) - \left(\gamma \widetilde{V}^{\pi_\theta}(s') - \widetilde{V}^{\pi_\theta}(s) \right) \right)^2 \right]} \leq 2\sqrt{2}\epsilon.$$

Assumption 4.4 (Lipschitz Continuity of Policy) There exist constants L and S such that for any $\theta \in \Omega(\theta)$, $s \in \mathcal{S}$, $a \in \mathcal{A}$,

$$\|\nabla_{\theta} \log \pi_{\theta}(a|s)\|_2 \leq L, \quad \|\nabla_{\theta}^2 \log \pi_{\theta}(a|s)\|_2 \leq S.$$

Assumption 4.4 is standard in the literature of policy gradient and actor-critic (Wu et al., 2020; Chen et al., 2021; Olshevsky & Ghahesifard, 2023; Chen & Zhao, 2023; Tian et al., 2023), which further implies the following proposition that the policy π_{θ} is Lipschitz continuous with respect to θ (refer to Appendix D for a detailed proof):

Proposition 4.5 For any $\theta_1, \theta_2 \in \Omega(\theta)$, $s \in \mathcal{S}$, $\|\pi_{\theta_1}(\cdot|s) - \pi_{\theta_2}(\cdot|s)\|_1 \leq L\|\theta_1 - \theta_2\|_2$.

Apart from the policy, we also need the Lipschitz continuity for the regularized reward to control the bias caused by the evolving reward.

Assumption 4.6 (Lipschitz Continuity of Regularized Reward) There exist constants $C, D > 0$ such that for any $\theta \in \Omega(\theta)$, $\varphi \in \Omega(\varphi)$, and $s \in \mathcal{S}$, the expected regularized reward satisfies:

1. **Boundedness:** $|\mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[\tilde{r}_{\varphi, \theta}(s, a)]| \leq C$
2. **Bounded Variance:** $\mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[\tilde{r}_{\varphi, \theta}(s, a)^2] \leq C^2$
3. **Lipschitz in θ :** $\|\nabla_{\theta} \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[\tilde{r}_{\varphi, \theta}(s, a)]\|_2 \leq CL$
4. **Smoothness in θ :** $\|\nabla_{\theta}^2 \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[\tilde{r}_{\varphi, \theta}(s, a)]\|_2 \leq C(L^2 + S)$
5. **Lipschitz in φ :** $\|\nabla_{\varphi} \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[\tilde{r}_{\varphi, \theta}(s, a)]\|_2 \leq D$

Parts 1 and 2 of Assumption 4.6 are natural extensions of the standard bounded-reward assumption to the expected regularized reward. The expectation over actions is necessary because the entropy regularization term $-\alpha \log \pi_{\theta}(a|s)$ can be unbounded for individual actions, but its expectation $\alpha \mathcal{H}(\pi_{\theta}(\cdot|s))$ is bounded for finite action spaces. For continuous action spaces, entropy regularization implicitly constrains the policy to have bounded entropy, as unbounded entropy would lead to infinite negative rewards, which is practically avoided.

Parts 3 and 4 of Assumption 4.6 naturally follow from Part 1 and Assumption 4.4. Since

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[\tilde{r}_{\varphi, \theta}(s, a)] &= \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[\tilde{r}_{\varphi, \theta}(s, a) \nabla_{\theta} \log \pi_{\theta}(a|s)], \\ \nabla_{\theta}^2 \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[\tilde{r}_{\varphi, \theta}(s, a)] &= \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[(\tilde{r}_{\varphi, \theta}(s, a) - \alpha) \nabla_{\theta} \log \pi_{\theta}(a|s) \nabla_{\theta} \log \pi_{\theta}(a|s)^{\top}] \\ &\quad + \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[\tilde{r}_{\varphi, \theta}(s, a) \nabla_{\theta}^2 \log \pi_{\theta}(a|s)], \end{aligned}$$

the Lipschitz continuity and smoothness w.r.t. θ can be derived from the boundedness of $\mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[\tilde{r}_{\varphi, \theta}(s, a)]$, $\nabla_{\theta} \log \pi_{\theta}(a|s)$, $\nabla_{\theta}^2 \log \pi_{\theta}(a|s)$ and α .

Part 5 of Assumption 4.6 is the most critical for handling evolving rewards. It guarantees that small changes in the reward parameters φ (e.g., from reward shaping or entropy adjustment) lead to proportionally small changes in the expected reward. This allows the algorithm to track the evolving learning objective rather than being destabilized by it.

In essence, Assumption 4.6 ensures that the regularized reward function changes in a controlled and predictable manner as the policy and reward parameters evolve. This is crucial for analyzing non-stationary learning dynamics.

4.2 MAIN RESULTS

With the assumptions above, we are ready to present our finite-time analysis of Algorithm 1. We measure the performance of Algorithm 1 using the following time-averaged errors over the second half of the T iterations:

- **Actor Error:** $G_T = \frac{1}{T/2} \sum_{t=T/2}^{T-1} \mathbb{E} \|\nabla_{\theta} J_{\varphi_t}(\theta_t)\|_2^2$
- **Critic Error:** $W_T = \frac{1}{T/2} \sum_{t=T/2}^{T-1} \mathbb{E} \|\omega_t - \omega_t^*\|_2^2$, where $\omega_t^* = \omega^*(\varphi_t, \theta_t)$
- **Reward Variation:** $F_T = \frac{1}{T/2} \sum_{t=T/2}^{T-1} \mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2^2$

Theorem 4.7 Consider Algorithm 1 with $\eta_t^\theta = \frac{c_\theta}{\sqrt{t}}$ and $\eta_t^\omega = \frac{c_\omega}{\sqrt{t}}$, where the ratio $\frac{c_\theta}{c_\omega}$ is chosen to be sufficiently small such that $\frac{c_\theta}{c_\omega} \leq \frac{\lambda}{LS_\omega} \wedge \frac{1}{16LL_\omega}$. Under Assumption 4.1, 4.4 and 4.6, the following bounds hold:

$$G_T = O\left(\frac{1}{\sqrt{T}}\right) + O\left(F_T \sqrt{T}\right) + O\left(\sqrt{\frac{F_T}{T}}\right) + O(\epsilon)$$

$$W_T = O\left(\frac{1}{\sqrt{T}}\right) + O\left(F_T \sqrt{T}\right) + O\left(\sqrt{\frac{F_T}{T}}\right) + O(\epsilon)$$

Interpretation of Theorem 4.7:

- **Static-Reward Case** ($F_T \equiv 0$): The algorithm achieves the canonical $O(1/\sqrt{T})$ convergence rate for both actor and critic, matching the best-known rate for single-timescale actor-critic methods under i.i.d. sampling (Chen et al., 2021; Olshevsky & Ghahserifard, 2023). Our analysis, by carefully handling the Markovian sampling, improves upon previous works by eliminating a $\log^2 T$ factor compared to Tian et al. (2023); Chen & Zhao (2023; 2025).
- **Evolving-Reward Case** ($F_T > 0$): The convergence rate depends critically on the total variation of the reward parameters, F_T . For the errors to converge to zero asymptotically, we require $F_T = o(1/\sqrt{T})$. To preserve the $O(1/\sqrt{T})$ rate, we need the stronger condition $F_T = O(1/T)$. This means the reward function must change slowly enough for the actor-critic algorithm to track it effectively.

The following corollary shows that a common class of reward update rules satisfies this stringent condition.

Corollary 4.8 If the reward parameter adopts a gradient-based update rule, i.e.

$$\varphi_{t+1} \leftarrow \varphi_t + \eta_t^\varphi h_\varphi(t),$$

then given $\mathbb{E} \|h_\varphi(t)\|_2^2 \leq C_\varphi^2$ and $\eta_t^\varphi = \frac{c_\varphi}{\sqrt{t}}$ where C_φ and c_φ are constants, we have $F_T = O\left(\frac{1}{T}\right)$, and hence

$$G_T = O\left(\frac{1}{\sqrt{T}}\right) + O(\epsilon), \quad W_T = O\left(\frac{1}{\sqrt{T}}\right) + O(\epsilon).$$

The proof of Corollary 4.8 can be found in Appendix D. Here, the step size for updating the reward parameter is of the same order as the actor’s and the critic’s, and hence the requirements can be achieved by applying gradient clipping, a technique that is very common in practice. Therefore, Corollary 4.8 provides a solid theoretical foundation for a wide range of empirical practice of RL.

4.3 PROOF SKETCH OF THE MAIN THEOREM

The proof of Theorem 4.7 proceeds in three interconnected steps, which we outline below. A rigorous proof is provided in Appendix C. The key innovations lie in (1) rigorously analyzing the impact of evolving rewards by establishing Lipschitz continuity properties of policy objective $J_\varphi(\theta)$ and optimal critic parameter $\omega^*(\varphi, \theta)$ w.r.t φ , and (2) providing a novel analysis on the distribution mismatch induced by Markovian sampling through deriving the following key proposition (refer to Appendix D for a detailed proof):

Proposition 4.9 *Following the Markovian sampling strategy described in Algorithm 1, we have*

$$\mathbb{E}\|\hat{\nu}_t - \nu_\rho^{\pi_{\theta_t}}\|_1 \leq LC_\delta L_\nu \sum_{k=0}^{t-1} \gamma^{t-1-k} \eta_k^\theta + \gamma^t \|\rho - \nu_\rho^{\pi_{\theta_0}}\|_1$$

for any $t \geq 0$, where C_δ and L_ν are constants (refer to Appendix B for a formal definition).

This analysis does not rely on the mixing time of the ergodic Markov chain. Instead, it directly utilizes the contraction properties of the induced operator acting on state distributions, which is a stronger property than the ergodicity, hence resulting in a tighter bound on the distribution mismatch.

Step 1: Bounding the Actor Error. The primary challenge introduced by an evolving reward is that the policy optimization objective $J_{\varphi_t}(\theta_t)$ changes at every time step t . To address this, we first show that $J_\varphi(\theta)$ is D_J -Lipschitz with respect to the reward parameter φ (Lemma B.1), which allows us to bound the change in the objective function by the change in φ :

$$\mathbb{E}[J_{\varphi_{t+1}}(\theta_{t+1})] - J_{\varphi_t}(\theta_t) \geq -D_J \mathbb{E}\|\varphi_{t+1} - \varphi_t\|_2 + \mathbb{E}[J_{\varphi_t}(\theta_{t+1})] - J_{\varphi_t}(\theta_t)$$

We then analyze the improvement in the objective for a fixed reward parameter. A Taylor expansion of $J_{\varphi_t}(\theta)$ around θ_t yields a bound on the squared policy gradient norm, $\|\nabla_\theta J_{\varphi_t}(\theta_t)\|_2^2$. This bound involves several error terms:

- I_1 (**Approximation Error**): This term arises from the bias introduced by linear function approximation and is bounded by $O(\epsilon)$.
- I_2 (**Critic Error**): This term captures the error from using an estimated critic ω_t instead of the optimal critic ω_t^* and is bounded by $O(\|\omega_t - \omega_t^*\|_2)$.
- I_3 (**Markovian Noise**): This term quantifies the error due to sampling states from the Markovian distribution $\hat{\nu}_t$ rather than the true stationary distribution $\nu_\rho^{\pi_{\theta_t}}$. Proposition 4.9 provides a tighter bound on this distribution mismatch, which is crucial for improving the overall convergence rate.

After summing over iterations and applying a telescoping series, we obtain the following inequality for the actor error (Theorem C.1):

$$(1 - \gamma)G_T \leq 2L\sqrt{G_T W_T} + O\left(\sqrt{\frac{F_T}{T}}\right) + O\left(\frac{1}{\sqrt{T}}\right) + O(\epsilon) \quad (14)$$

Step 2: Bounding the Critic Error. The critic update must track a moving target: the optimal parameter $\omega_t^* = \omega^*(\varphi_t, \theta_t)$ changes with both the policy parameter θ_t and the reward parameter φ_t . We analyze the evolution of the critic error $\|\omega_{t+1} - \omega_{t+1}^*\|_2^2$. A central decomposition yields:

$$\begin{aligned} \mathbb{E}\|\omega_{t+1} - \omega_{t+1}^*\|_2^2 &\leq \|\omega_t - \omega_t^*\|_2^2 + 2\mathbb{E}\left\|\hat{\delta}(s_t, a_t, s'_t)\phi(s_t)\right\|_2^2 + 2\mathbb{E}\|\omega_t^* - \omega_{t+1}^*\|_2^2 \\ &\quad + 2\eta_t^\omega \mathbb{E}\left\langle \omega_t - \omega_t^*, \hat{\delta}(s_t, a_t, s'_t)\phi(s_t) \right\rangle + 2\mathbb{E}\left\langle \omega_t - \omega_t^*, \omega_t^* - \omega_{t+1}^* \right\rangle, \end{aligned}$$

where $\mathbb{E}\left\langle \omega_t - \omega_t^*, \hat{\delta}(s_t, a_t, s'_t)\phi(s_t) \right\rangle$ is further decomposed into three components:

- J_1 : This term is zero by the definition of ω_t^* .
- J_2 (**Contraction**): This term provides a negative contribution $-\lambda \|\omega_t - \omega_t^*\|_2^2$, ensuring the critic error contracts towards zero.
- J_3 (**Markovian Noise**): Similar to I_3 in the actor analysis, this term is bounded using Proposition 4.9.

The critical difference from the static-reward case is the presence of terms involving $\omega_t^* - \omega_{t+1}^*$. We bound these by establishing the Lipschitz continuity of ω^* with respect to both θ and φ (Lemma B.5). This introduces terms proportional to $\mathbb{E}\|\varphi_{t+1} - \varphi_t\|_2^2$ and $\mathbb{E}\|\varphi_{t+1} - \varphi_t\|_2$ into the bound. After summation, we derive the following inequality for the critic error (Theorem C.2):

$$\frac{1}{1 - \gamma} W_T \leq 2L\omega \frac{c_\theta}{c_\omega} \sqrt{G_T W_T} + O\left(F_T \sqrt{T}\right) + O\left(\sqrt{\frac{F_T}{T}}\right) + O\left(\frac{1}{\sqrt{T}}\right) + O(\epsilon) \quad (15)$$

Step 3: Solving the System of Inequalities Steps 1 and 2 result in a system of two inequalities (14 and 15) that couple the actor error G_T and the critic error W_T . To solve this system, we use the algebraic inequality

$$2\sqrt{G_T W_T} \leq \frac{1-\gamma}{2L} G_T + \frac{2L}{1-\gamma} W_T.$$

By substituting this into the inequalities and choosing the step-size ratio $\frac{c_e}{c_w}$ to be sufficiently small, we can decouple the two errors and obtain a bound of $O(1/\sqrt{T}) + O(F_T\sqrt{T}) + O(\sqrt{F_T/T}) + O(\epsilon)$ for both G_T and W_T , thus completing the proof.

5 CONCLUSION

In this work, we have undertaken a systematic theoretical investigation of actor-critic methods in the presence of evolving rewards—a setting that mirrors the reality of many practical RL algorithms but has been largely overlooked by theoretical analyses. We formulated the problem, established necessary assumptions, and provided the first finite-time convergence guarantees for a single-timescale actor-critic algorithm under Markovian sampling.

Our analysis demonstrates that the single-timescale actor-critic algorithm is remarkably robust to reward non-stationarity. The canonical $O(1/\sqrt{T})$ convergence rate can be maintained for both the actor and critic, provided the reward parameters evolve at a controlled pace. A key corollary confirms that gradient-based reward updates—a common pattern in algorithms that learn intrinsic rewards or adapt regularization strengths—satisfy this condition, thereby providing a solid theoretical foundation for their empirical success. Furthermore, our novel technique for bounding distribution mismatch under Markovian sampling yields a tighter analysis, improving upon prior rates by a factor of $\log^2 T$ even when the reward is static.

This work opens several avenues for future research. Extending the analysis to nonlinear function approximation, particularly with neural networks, is a critical next step. Furthermore, exploring the implications of our theoretical findings for the design of more effective and provably stable reward-shaping algorithms presents an exciting direction for both theoretical and applied work. Finally, this analysis lays a foundational stone for a deeper theoretical understanding of reinforcement learning with dynamic objectives due to evolving reward, shifting initial distribution or transition probabilities.

REFERENCES

- Alekh Agarwal, Sham M. Kakade, Jason D. Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *J. Mach. Learn. Res.*, 22: 98:1–98:76, 2021. URL <https://jmlr.org/papers/v22/19-736.html>.
- Zafarali Ahmed, Nicolas Le Roux, Mohammad Norouzi, and Dale Schuurmans. Understanding the impact of entropy on policy optimization. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 151–160. PMLR, 2019. URL <http://proceedings.mlr.press/v97/ahmed19a.html>.
- John Asmuth, Michael L. Littman, and Robert Zinkov. Potential-based shaping in model-based reinforcement learning. In Dieter Fox and Carla P. Gomes (eds.), *Proceedings of the Twenty-Third AAAI Conference on Artificial Intelligence, AAAI 2008, Chicago, Illinois, USA, July 13-17, 2008*, pp. 604–609. AAAI Press, 2008. URL <http://www.aaai.org/Library/AAAI/2008/aaai08-096.php>.
- Jalaj Bhandari, Daniel Russo, and Raghav Singal. A finite time analysis of temporal difference learning with linear function approximation. In Sébastien Bubeck, Vianney Perchet, and Philippe Rigollet (eds.), *Conference On Learning Theory, COLT 2018, Stockholm, Sweden, 6-9 July 2018*, volume 75 of *Proceedings of Machine Learning Research*, pp. 1691–1692. PMLR, 2018. URL <http://proceedings.mlr.press/v75/bhandari18a.html>.

- 540 Yuri Burda, Harrison Edwards, Amos J. Storkey, and Oleg Klimov. Exploration by random network
541 distillation. In *7th International Conference on Learning Representations, ICLR 2019, New Or-*
542 *leans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL [https://openreview.net](https://openreview.net/forum?id=H1lJjN5Ym)
543 [/forum?id=H1lJjN5Ym](https://openreview.net/forum?id=H1lJjN5Ym).
- 544 Semih Cayci, Niao He, and R. Srikant. Finite-time analysis of entropy-regularized neural natural
545 actor-critic algorithm. *Trans. Mach. Learn. Res.*, 2024, 2024. URL [https://openreview](https://openreview.net/forum?id=BkEqk7pS1I)
546 [.net/forum?id=BkEqk7pS1I](https://openreview.net/forum?id=BkEqk7pS1I).
- 547 Tianyi Chen, Yuejiao Sun, and Wotao Yin. Closing the gap: Tighter analysis of alternating stochastic
548 gradient methods for bilevel problems. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N.
549 Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information*
550 *Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021,*
551 *NeurIPS 2021, December 6-14, 2021, virtual*, pp. 25294–25307, 2021. URL [https://proc](https://proceedings.neurips.cc/paper/2021/hash/d4dd111a4fd973394238aca5c05bebe3-Abstract.html)
552 [eedings.neurips.cc/paper/2021/hash/d4dd111a4fd973394238aca5c05be](https://proceedings.neurips.cc/paper/2021/hash/d4dd111a4fd973394238aca5c05bebe3-Abstract.html)
553 [be3-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/d4dd111a4fd973394238aca5c05bebe3-Abstract.html).
- 554 Xuyang Chen and Lin Zhao. Finite-time analysis of single-timescale actor-critic. In Alice Oh, Tris-
555 tan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances*
556 *in Neural Information Processing Systems 36: Annual Conference on Neural Information Pro-*
557 *cessing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
558 URL [http://papers.nips.cc/paper_files/paper/2023/hash/160adf2dc](http://papers.nips.cc/paper_files/paper/2023/hash/160adf2dc118a920e7858484b92a37d8-Abstract-Conference.html)
559 [118a920e7858484b92a37d8-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/160adf2dc118a920e7858484b92a37d8-Abstract-Conference.html).
- 560 Xuyang Chen and Lin Zhao. On the convergence of continuous single-timescale actor-critic. In
561 *Forty-second International Conference on Machine Learning*, 2025. URL [https://openre](https://openreview.net/forum?id=pV7hSmGJXP)
562 [view.net/forum?id=pV7hSmGJXP](https://openreview.net/forum?id=pV7hSmGJXP).
- 563 Wang Chi Cheung, David Simchi-Levi, and Ruihao Zhu. Nonstationary reinforcement learning: The
564 blessing of (more) optimism. *Manage. Sci.*, 69(10):5722–5739, October 2023. ISSN 0025-1909.
565 doi: 10.1287/mnsc.2023.4704. URL <https://doi.org/10.1287/mnsc.2023.4704>.
- 566 DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu,
567 Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu,
568 Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao
569 Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,
570 Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao,
571 Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding,
572 Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang
573 Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai
574 Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang,
575 Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang,
576 Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang,
577 Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang,
578 R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye,
579 Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, and S. S. Li. Deepseek-r1: Incentivizing
580 reasoning capability in llms via reinforcement learning. *CoRR*, abs/2501.12948, 2025. doi: 10.4
581 8550/ARXIV.2501.12948. URL <https://doi.org/10.48550/arXiv.2501.12948>.
- 582 Sam Devlin and Daniel Kudenko. Dynamic potential-based reward shaping. In *Proceedings of*
583 *the 11th International Conference on Autonomous Agents and Multiagent Systems - Volume 1,*
584 *AAMAS ’12*, pp. 433–440, Richland, SC, 2012. International Foundation for Autonomous Agents
585 and Multiagent Systems. ISBN 0981738117.
- 586 Yuhao Ding, Junzi Zhang, and Javad Lavaei. On the global optimum convergence of momentum-
587 based policy gradient. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera (eds.),
588 *International Conference on Artificial Intelligence and Statistics, AISTATS 2022, 28-30 March*
589 *2022, Virtual Event*, volume 151 of *Proceedings of Machine Learning Research*, pp. 1910–1934.
590 PMLR, 2022. URL <https://proceedings.mlr.press/v151/ding22a.html>.
- 591 Yuhao Ding, Junzi Zhang, Hyunin Lee, and Javad Lavaei. Beyond exact gradients: Convergence
592 of stochastic soft-max policy gradient methods with entropy regularization. *IEEE Trans. Autom.*
593

- 594 *Control.*, 70(8):5129–5144, 2025. doi: 10.1109/TAC.2025.3540965. URL <https://doi.org/10.1109/TAC.2025.3540965>.
- 595
- 596
- 597 Eyal Even-Dar, Sham. M. Kakade, and Yishay Mansour. Online markov decision processes. *Mathematics of Operations Research*, 34(3):726–736, 2009. doi: 10.1287/moor.1090.0396. URL <https://doi.org/10.1287/moor.1090.0396>.
- 598
- 599
- 600 Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Reinforcement learning for fine-tuning text-to-image diffusion models. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/fc65fab891d83433bd3c8d966edde311-Abstract-Conference.html.
- 601
- 602
- 603
- 604
- 605
- 606
- 607
- 608 Ilyas Fatkhullin, Anas Barakat, Anastasia Kireeva, and Niao He. Stochastic policy gradient methods: Improved sample complexity for fisher-non-degenerate policies. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 9827–9869. PMLR, 2023. URL <https://proceedings.mlr.press/v202/fatkhullin23a.html>.
- 609
- 610
- 611
- 612
- 613
- 614
- 615 Songtao Feng, Ming Yin, Ruiquan Huang, Yu-Xiang Wang, Jing Yang, and Yingbin Liang. Non-stationary reinforcement learning under general function approximation. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 9976–10007. PMLR, 2023. URL <https://proceedings.mlr.press/v202/feng23e.html>.
- 616
- 617
- 618
- 619
- 620
- 621 Swetha Ganesh, Jiayu Chen, Washim Uddin Mondal, and Vaneet Aggarwal. Order-optimal global convergence for actor-critic with general policy and neural critic parametrization. In *The 41st Conference on Uncertainty in Artificial Intelligence*, 2025. URL <https://openreview.net/forum?id=HPxE1IejA5>.
- 622
- 623
- 624
- 625 Mudit Gaur, Amrit Bedi, Di Wang, and Vaneet Aggarwal. Closing the gap: Achieving global convergence (Last iterate) of actor-critic under Markovian sampling with neural network parametrization. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 15153–15179. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/gaur24a.html>.
- 626
- 627
- 628
- 629
- 630
- 631
- 632 Dhawal Gupta, Yash Chandak, Scott M. Jordan, Philip S. Thomas, and Bruno C. da Silva. Behavior alignment via reward function optimization. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/a5357781c204d4412e44ed9cbcd08d5-Abstract-Conference.html.
- 633
- 634
- 635
- 636
- 637
- 638
- 639
- 640 Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with deep energy-based policies. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1352–1361. PMLR, 2017. URL <http://proceedings.mlr.press/v70/haarnoja17a.html>.
- 641
- 642
- 643
- 644
- 645 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In Jennifer G. Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings*
- 646
- 647

- 648 *of Machine Learning Research*, pp. 1856–1865. PMLR, 2018a. URL [http://proceedings](http://proceedings.mlr.press/v80/haarnoja18b.html)
649 [s.mlr.press/v80/haarnoja18b.html](http://proceedings.mlr.press/v80/haarnoja18b.html).
- 650
- 651 Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash
652 Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Soft actor-critic algo-
653 rithms and applications. *CoRR*, abs/1812.05905, 2018b. URL [http://arxiv.org/abs/18](http://arxiv.org/abs/1812.05905)
654 [12.05905](http://arxiv.org/abs/1812.05905).
- 655 Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang. A two-timescale stochastic algorithm
656 framework for bilevel optimization: Complexity analysis and application to actor-critic. *SIAM*
657 *Journal on Optimization*, 33(1):147–180, 2023. doi: 10.1137/20M1387341. URL [https:](https://doi.org/10.1137/20M1387341)
658 [//doi.org/10.1137/20M1387341](https://doi.org/10.1137/20M1387341).
- 659 Yujing Hu, Weixun Wang, Hangtian Jia, Yixiang Wang, Yingfeng Chen, Jianye Hao, Feng Wu,
660 and Changjie Fan. Learning to utilize shaping rewards: A new approach of reward shaping.
661 In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien
662 Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural*
663 *Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL
664 [https://proceedings.neurips.cc/paper/2020/hash/b710915795b9e9c02](https://proceedings.neurips.cc/paper/2020/hash/b710915795b9e9c02cf10d6d2bdb688c-Abstract.html)
665 [cf10d6d2bdb688c-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/b710915795b9e9c02cf10d6d2bdb688c-Abstract.html).
- 666 Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Àgata Lapedriza, Noah
667 Jones, Shixiang Gu, and Rosalind W. Picard. Way off-policy batch deep reinforcement learning
668 of implicit human preferences in dialog. *CoRR*, abs/1907.00456, 2019. URL [http://arxiv.](http://arxiv.org/abs/1907.00456)
669 [org/abs/1907.00456](http://arxiv.org/abs/1907.00456).
- 670 Chi Jin, Tiancheng Jin, Haipeng Luo, Suvrit Sra, and Tiancheng Yu. Learning adversarial Markov
671 decision processes with bandit feedback and unknown transition. In Hal Daumé III and Aarti
672 Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume
673 119 of *Proceedings of Machine Learning Research*, pp. 4860–4869. PMLR, 13–18 Jul 2020.
674 URL <https://proceedings.mlr.press/v119/jin20c.html>.
- 675
- 676 Vijay R. Konda and John N. Tsitsiklis. Actor-critic algorithms. In Sara A. Solla, Todd K. Leen,
677 and Klaus-Robert Müller (eds.), *Advances in Neural Information Processing Systems 12, [NIPS*
678 *Conference, Denver, Colorado, USA, November 29 - December 4, 1999]*, pp. 1008–1014. The
679 MIT Press, 1999. URL [http://papers.nips.cc/paper/1786-actor-critic-a](http://papers.nips.cc/paper/1786-actor-critic-algorithms)
680 [lgorithms](http://papers.nips.cc/paper/1786-actor-critic-algorithms).
- 681 Harshat Kumar, Alec Koppel, and Alejandro Ribeiro. On the sample complexity of actor-critic
682 method for reinforcement learning with function approximation. *Mach. Learn.*, 112(7):2433–
683 2467, 2023. doi: 10.1007/S10994-023-06303-2. URL [https://doi.org/10.1007/s1](https://doi.org/10.1007/s10994-023-06303-2)
684 [0994-023-06303-2](https://doi.org/10.1007/s10994-023-06303-2).
- 685 Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,
686 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In Yoshua
687 Bengio and Yann LeCun (eds.), *4th International Conference on Learning Representations, ICLR*
688 *2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. URL [http:](http://arxiv.org/abs/1509.02971)
689 [//arxiv.org/abs/1509.02971](http://arxiv.org/abs/1509.02971).
- 690 Yanli Liu, Kaiqing Zhang, Tamer Basar, and Wotao Yin. An improved analysis of (variance-reduced)
691 policy gradient and natural policy gradient methods. In Hugo Larochelle, Marc’Aurelio Ranzato,
692 Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information*
693 *Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020,*
694 *NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL [https://proceedings.neur](https://proceedings.neurips.cc/paper/2020/hash/56577889b3c1cd083b6d7b32d32f99d5-Abstract.html)
695 [ips.cc/paper/2020/hash/56577889b3c1cd083b6d7b32d32f99d5-Abstrac](https://proceedings.neurips.cc/paper/2020/hash/56577889b3c1cd083b6d7b32d32f99d5-Abstract.html)
696 [t.html](https://proceedings.neurips.cc/paper/2020/hash/56577889b3c1cd083b6d7b32d32f99d5-Abstract.html).
- 697 Sam Lobel, Akhil Bagaria, and George Konidaris. Flipping coins to estimate pseudocounts for
698 exploration in reinforcement learning. In Andreas Krause, Emma Brunskill, Kyunghyun Cho,
699 Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Ma-*
700 *chine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Pro-*
701 *ceedings of Machine Learning Research*, pp. 22594–22613. PMLR, 2023. URL [https:](https://proceedings.mlr.press/v202/lobel23a.html)
[//proceedings.mlr.press/v202/lobel23a.html](https://proceedings.mlr.press/v202/lobel23a.html).

- 702 Haozhe Ma, Kuankuan Sima, Thanh Vinh Vo, Di Fu, and Tze-Yun Leong. Reward shaping for
703 reinforcement learning with an assistant reward agent. In *Forty-first International Conference on*
704 *Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. URL
705 <https://openreview.net/forum?id=a3XFF0PGLU>.
706
- 707 Marlos C. Machado, Marc G. Bellemare, and Michael Bowling. Count-based exploration with
708 the successor representation. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence,*
709 *AAAI 2020, New York, NY, USA, February 7-12, 2020*, pp. 5125–5133. AAAI Press, 2020. doi:
710 10.1609/AAAI.V34I04.5955. URL <https://doi.org/10.1609/aaai.v34i04.5955>.
711
- 712 Srinath Mahankali, Zhang-Wei Hong, Ayush Sekhari, Alexander Rakhlin, and Pulkit Agrawal. Ran-
713 dom latent exploration for deep reinforcement learning. In *Forty-first International Conference*
714 *on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024.
715 URL <https://openreview.net/forum?id=Y9qzwNlKVU>.
716
- 717 Debmalya Mandal and Goran Radanovic. Performative reinforcement learning with linear markov
718 decision process. In Yingzhen Li, Stephan Mandt, Shipra Agrawal, and Emtiyaz Khan (eds.),
719 *Proceedings of The 28th International Conference on Artificial Intelligence and Statistics*, volume
720 258 of *Proceedings of Machine Learning Research*, pp. 3232–3240. PMLR, 03–05 May 2025.
721 URL <https://proceedings.mlr.press/v258/mandal25a.html>.
722
- 723 Debmalya Mandal, Stelios Triantafyllou, and Goran Radanovic. Performative reinforcement learn-
724 ing. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato,
725 and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29*
726 *July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*,
727 pp. 23642–23680. PMLR, 2023. URL [https://proceedings.mlr.press/v202/man-](https://proceedings.mlr.press/v202/mandal23a.html)
728 [dal23a.html](https://proceedings.mlr.press/v202/mandal23a.html).
729
- 730 Weichao Mao, Kaiqing Zhang, Ruihao Zhu, David Simchi-Levi, and Tamer Basar. Model-
731 free nonstationary reinforcement learning: Near-optimal regret and applications in multiagent
732 reinforcement learning and inventory control. *Manag. Sci.*, 71(2):1564–1580, 2025. doi:
733 10.1287/MNSC.2022.02533. URL <https://doi.org/10.1287/mnsc.2022.02533>.
734
- 735 Jarryd Martin, S. Suraj Narayanan, Tom Everitt, and Marcus Hutter. Count-based exploration in fea-
736 ture space for reinforcement learning. In *Proceedings of the 26th International Joint Conference*
737 *on Artificial Intelligence, IJCAI’17*, pp. 2471–2478. AAAI Press, 2017. ISBN 9780999241103.
738
- 739 Jincheng Mei, Chenjun Xiao, Csaba Szepesvári, and Dale Schuurmans. On the global convergence
740 rates of softmax policy gradient methods. In *Proceedings of the 37th International Conference*
741 *on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings*
742 *of Machine Learning Research*, pp. 6820–6829. PMLR, 2020. URL [http://proceedings.](http://proceedings.mlr.press/v119/mei20b.html)
743 [mlr.press/v119/mei20b.html](http://proceedings.mlr.press/v119/mei20b.html).
744
- 745 Farzan Memarian, Wonjoon Goo, Rudolf Lioutikov, Scott Niekum, and Ufuk Topcu. Self-supervised
746 online reward shaping in sparse-reward environments. In *IEEE/RSJ International Conference on*
747 *Intelligent Robots and Systems, IROS 2021, Prague, Czech Republic, September 27 - Oct. 1,*
748 *2021*, pp. 2369–2375. IEEE, 2021. doi: 10.1109/IROS51168.2021.9636020. URL <https://doi.org/10.1109/IROS51168.2021.9636020>.
749
- 750 David Mguni, Taher Jafferjee, Jianhong Wang, Nicolas Perez Nieves, Wenbin Song, Feifei Tong,
751 Matthew E. Taylor, Tianpei Yang, Zipeng Dai, Hui Chen, Jiangcheng Zhu, Kun Shao, Jun Wang,
752 and Yaodong Yang. Learning to shape rewards using a game of two partners. In Brian Williams,
753 Yiling Chen, and Jennifer Neville (eds.), *Thirty-Seventh AAAI Conference on Artificial Intelli-*
754 *gence, AAAI 2023, Washington, DC, USA, February 7-14, 2023*, pp. 11604–11612. AAAI Press,
755 2023. doi: 10.1609/AAAI.V37I10.26371. URL [https://doi.org/10.1609/aaai.v37](https://doi.org/10.1609/aaai.v37i10.26371)
[i10.26371](https://doi.org/10.1609/aaai.v37i10.26371).
756
- 757 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
758 Wierstra, and Martin A. Riedmiller. Playing atari with deep reinforcement learning. *CoRR*,
759 abs/1312.5602, 2013. URL <http://arxiv.org/abs/1312.5602>.

- 756 Washim Uddin Mondal and Vaneet Aggarwal. Improved sample complexity analysis of natural
757 policy gradient algorithm with general parameterization for infinite horizon discounted reward
758 markov decision processes. In Sanjoy Dasgupta, Stephan Mandt, and Yingzhen Li (eds.), *Inter-
759 national Conference on Artificial Intelligence and Statistics, 2-4 May 2024, Palau de Congres-
760 sos, Valencia, Spain*, volume 238 of *Proceedings of Machine Learning Research*, pp. 3097–3105.
761 PMLR, 2024. URL <https://proceedings.mlr.press/v238/u-mondal24a.html>.
- 762 Sanmit Narvekar, Bei Peng, Matteo Leonetti, Jivko Sinapov, Matthew E. Taylor, and Peter Stone.
763 Curriculum learning for reinforcement learning domains: A framework and survey. *J. Mach.
764 Learn. Res.*, 21:181:1–181:50, 2020. URL [https://jmlr.org/papers/v21/20-212](https://jmlr.org/papers/v21/20-212.html)
765 [.html](https://jmlr.org/papers/v21/20-212.html).
- 767 Andrew Y. Ng, Daishi Harada, and Stuart J. Russell. Policy invariance under reward transformations:
768 Theory and application to reward shaping. In *Proceedings of the Sixteenth International Confer-
769 ence on Machine Learning, ICML ’99*, pp. 278–287, San Francisco, CA, USA, 1999. Morgan
770 Kaufmann Publishers Inc. ISBN 1558606122.
- 771 Alex Olshevsky and Bahman Ghahserifard. A small gain analysis of single timescale actor critic.
772 *SIAM J. Control. Optim.*, 61(2):980–1007, 2023. doi: 10.1137/22M1483335. URL <https://doi.org/10.1137/22m1483335>.
- 774 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin,
775 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser
776 Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan
777 Leike, and Ryan Lowe. Training language models to follow instructions with human feedback.
778 In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Ad-
779 vances in Neural Information Processing Systems 35: Annual Conference on Neural Information
780 Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9,
781 2022*, 2022. URL [http://papers.nips.cc/paper_files/paper/2022/hash/b](http://papers.nips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html)
782 [1efde53be364a73914f58805a001731-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html).
- 783 Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven exploration
784 by self-supervised prediction. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the
785 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11
786 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2778–2787. PMLR,
787 2017. URL <http://proceedings.mlr.press/v70/pathak17a.html>.
- 789 Deepak Pathak, Dhiraj Gandhi, and Abhinav Gupta. Self-supervised exploration via disagreement.
790 In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International
791 Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*,
792 volume 97 of *Proceedings of Machine Learning Research*, pp. 5062–5071. PMLR, 2019. URL
793 <http://proceedings.mlr.press/v97/pathak19a.html>.
- 794 Aditya A. Ramesh, Louis Kirsch, Sjoerd van Steenkiste, and Jürgen Schmidhuber. Explor-
795 ing through random curiosity with general value functions. In Sanmi Koyejo, S. Mohamed,
796 A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information
797 Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022,
798 NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL [http://papers.nips.cc/paper_files/paper/2022/hash/76e57c3c6b3e06f332a](http://papers.nips.cc/paper_files/paper/2022/hash/76e57c3c6b3e06f332a4832ddd6a9a12-Abstract-Conference.html)
799 [4832ddd6a9a12-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/76e57c3c6b3e06f332a4832ddd6a9a12-Abstract-Conference.html).
- 801 Ben Rank, Stelios Triantafyllou, Debmalya Mandal, and Goran Radanovic. Performative reinforce-
802 ment learning in gradually shifting environments. In Negar Kiyavash and Joris M. Mooij (eds.),
803 *Proceedings of the Fortieth Conference on Uncertainty in Artificial Intelligence*, volume 244
804 of *Proceedings of Machine Learning Research*, pp. 3041–3075. PMLR, 15–19 Jul 2024. URL
805 <https://proceedings.mlr.press/v244/rank24a.html>.
- 806 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li,
807 Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open
808 language models. *CoRR*, abs/2402.03300, 2024. doi: 10.48550/ARXIV.2402.03300. URL
809 <https://doi.org/10.48550/arXiv.2402.03300>.

- 810 Han Shen, Kaiqing Zhang, Mingyi Hong, and Tianyi Chen. Towards understanding asynchronous
811 advantage actor-critic: Convergence and linear speedup. *Trans. Sig. Proc.*, 71:2579–2594, January
812 2023. ISSN 1053-587X. doi: 10.1109/TSP.2023.3268475. URL [https://doi.org/10.1](https://doi.org/10.1109/TSP.2023.3268475)
813 [109/TSP.2023.3268475](https://doi.org/10.1109/TSP.2023.3268475).
- 814
815 Bradly C. Stadie, Lunjun Zhang, and Jimmy Ba. Learning intrinsic rewards as a bi-level optimization
816 problem. In Ryan P. Adams and Vibhav Gogate (eds.), *Proceedings of the Thirty-Sixth Conference*
817 *on Uncertainty in Artificial Intelligence, UAI 2020, virtual online, August 3-6, 2020*, volume 124
818 of *Proceedings of Machine Learning Research*, pp. 111–120. AUAI Press, 2020. URL [http:](http://proceedings.mlr.press/v124/stadie20a.html)
819 [//proceedings.mlr.press/v124/stadie20a.html](http://proceedings.mlr.press/v124/stadie20a.html).
- 820 Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,
821 Dario Amodei, and Paul F. Christiano. Learning to summarize from human feedback. *CoRR*,
822 abs/2009.01325, 2020. URL <https://arxiv.org/abs/2009.01325>.
- 823
824 Hao Sun, Lei Han, Rui Yang, Xiaoteng Ma, Jian Guo, and Bolei Zhou. Exploit reward shifting
825 in value-based deep-rl: Optimistic curiosity-based exploration and conservative exploitation via
826 linear reward shaping. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho,
827 and A. Oh (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on*
828 *Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November*
829 *28 - December 9, 2022, 2022*. URL [http://papers.nips.cc/paper_files/paper](http://papers.nips.cc/paper_files/paper/2022/hash/f600d1a3f6a63f782680031f3ce241a7-Abstract-Conference.html)
830 [/2022/hash/f600d1a3f6a63f782680031f3ce241a7-Abstract-Conference.](http://papers.nips.cc/paper_files/paper/2022/hash/f600d1a3f6a63f782680031f3ce241a7-Abstract-Conference.html)
831 [html](http://papers.nips.cc/paper_files/paper/2022/hash/f600d1a3f6a63f782680031f3ce241a7-Abstract-Conference.html).
- 832
833 Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1. MIT
834 press Cambridge, 1998.
- 835
836 Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods
837 for reinforcement learning with function approximation. In S. Solla, T. Leen, and K. Müller
838 (eds.), *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 1999. URL
839 [https://proceedings.neurips.cc/paper_files/paper/1999/file/464d8](https://proceedings.neurips.cc/paper_files/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf)
840 [28b85b0bed98e80ade0a5c43b0f-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf).
- 841
842 Haoxing Tian, Alex Olshevsky, and Yannis Paschalidis. Convergence of actor-critic with multi-
843 layer neural networks. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine
844 (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 9279–9321. Curran
845 Associates, Inc., 2023. URL [https://proceedings.neurips.cc/paper_files/p](https://proceedings.neurips.cc/paper_files/paper/2023/file/1dc9fbd6b6b4d9955ad377cb983232c9f-Paper-Conference.pdf)
846 [aper/2023/file/1dc9fbd6b6b4d9955ad377cb983232c9f-Paper-Conferenc](https://proceedings.neurips.cc/paper_files/paper/2023/file/1dc9fbd6b6b4d9955ad377cb983232c9f-Paper-Conference.pdf)
847 [e.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/1dc9fbd6b6b4d9955ad377cb983232c9f-Paper-Conference.pdf).
- 848
849 Alexander Trott, Stephan Zheng, Caiming Xiong, and Richard Socher. Keeping your distance:
850 Solving sparse reward tasks using self-balancing shaped rewards. In Hanna M. Wallach, Hugo
851 Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.),
852 *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Informa-*
853 *tion Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp.
854 10376–10386, 2019. URL [https://proceedings.neurips.cc/paper/2019/hash](https://proceedings.neurips.cc/paper/2019/hash/64c26b2a2dcf068c49894bd07e0e6389-Abstract.html)
855 [/64c26b2a2dcf068c49894bd07e0e6389-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/64c26b2a2dcf068c49894bd07e0e6389-Abstract.html).
- 856
857 Eric Wiewiora. Potential-based shaping and q-value initialization are equivalent. *J. Artif. Int. Res.*,
858 19(1):205–208, September 2003. ISSN 1076-9757.
- 859
860 Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
861 learning. *Mach. Learn.*, 8:229–256, 1992. doi: 10.1007/BF00992696. URL [https://doi.](https://doi.org/10.1007/BF00992696)
862 [org/10.1007/BF00992696](https://doi.org/10.1007/BF00992696).
- 863
864 Yue Wu, Weitong Zhang, Pan Xu, and Quanquan Gu. A finite-time analysis of two time-scale actor-
865 critic methods. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan,
866 and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Con-*
867 *ference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020,*
868 *virtual*, 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/cc9](https://proceedings.neurips.cc/paper/2020/hash/cc9b3c69b56df284846bf2432f1cba90-Abstract.html)
869 [b3c69b56df284846bf2432f1cba90-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/cc9b3c69b56df284846bf2432f1cba90-Abstract.html).

- 864 Lin Xiao. On the convergence rates of policy gradient methods. *J. Mach. Learn. Res.*, 23:282:1–
865 282:36, 2022. URL <https://jmlr.org/papers/v23/22-0056.html>.
866
- 867 Tengyu Xu, Zhe Wang, and Yingbin Liang. Improving sample complexity bounds for (natural)
868 actor-critic algorithms. In Hugo Larochelle, Marc Aurelio Ranzato, Raia Hadsell, Maria-Florina
869 Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual
870 Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12,
871 2020, virtual*, 2020a. URL <https://proceedings.neurips.cc/paper/2020/hash/2e1b24a664f5e9c18f407b2f9c73e821-Abstract.html>.
872
- 873 Tengyu Xu, Zhe Wang, and Yingbin Liang. Non-asymptotic convergence analysis of two time-scale
874 (natural) actor-critic algorithms. *CoRR*, abs/2005.03557, 2020b. URL <https://arxiv.org/abs/2005.03557>.
875
- 876 Kai Yang, Jian Tao, Jiafei Lyu, and Xiu Li. Exploration and anti-exploration with distributional
877 random network distillation. In *Forty-first International Conference on Machine Learning, ICML
878 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=rIrpzmqRBk>.
879
- 880 Zhuoran Yang, Yongxin Chen, Mingyi Hong, and Zhaoran Wang. Provably global convergence
881 of actor-critic: A case for linear quadratic regulator with ergodic cost. In Hanna M. Wallach,
882 Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett
883 (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural
884 Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC,
885 Canada*, pp. 8351–8363, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/9713faa264b94e2bf346a1bb52587fd8-Abstract.html>.
886
- 887 Zeyu Zheng, Junhyuk Oh, and Satinder Singh. On learning intrinsic rewards for policy gradient
888 methods. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-
889 Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual
890 Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8,
891 2018, Montréal, Canada*, pp. 4649–4659, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/51de85ddd068f0bc787691d356176df9-Abstract.html>.
892
- 893 Alexander Zimin and Gergely Neu. Online learning in episodic markovian decision processes by
894 relative entropy policy search. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q.
895 Weinberger (eds.), *Advances in Neural Information Processing Systems*, volume 26. Curran As-
896 sociates, Inc., 2013. URL https://proceedings.neurips.cc/paper_files/paper/2013/file/68053af2923e00204c3ca7c6a3150cf7-Paper.pdf.
897

900 A LITERATURE REVIEW OF EVOLVING REWARD TECHNIQUES

901 The idea of modifying rewards to improve learning is long-standing. Potential-based reward shap-
902 ing, introduced and developed by Ng et al. (1999); Wiewiora (2003); Asmuth et al. (2008); Devlin
903 & Kudenko (2012), defines the shaped reward as $\gamma\Phi(s') - \Phi(s)$ where $\Phi(\cdot)$ is a potential func-
904 tion to guarantee policy invariance. More recent work, however, modifies the reward to balance the
905 exploration and exploitation behavior of the policy. This sorts of work include randomly perturbing
906 the reward function (Mahankali et al., 2024), various design of explicit exploration bonus like
907 count-based methods (Martin et al., 2017; Machado et al., 2020; Lobel et al., 2023), curiosity-driven
908 methods (Pathak et al., 2017; Ramesh et al., 2022; Sun et al., 2022) and random network
909 distillation (Burda et al., 2019; Yang et al., 2024), as well as fully self-supervised intrinsic rewards
910 (Zheng et al., 2018; Stadie et al., 2020; Memarian et al., 2021; Mguni et al., 2023; Ma et al., 2024) or
911 incorporation of prior knowledge (Trott et al., 2019; Hu et al., 2020; Gupta et al., 2023) that enhance
912 the performance of the resulting policy in terms of the original reward.
913

914 Another series of work that result in evolving rewards is the entropy or KL regularization. Entropy
915 regularization (Haarnoja et al., 2018a;b; Ahmed et al., 2019) is commonly used technique to en-
916 courage exploration and avoid near-deterministic suboptimal policy. KL regularization is common
917 in fine-tuning RL policies, especially in training Large Language Models (Ouyang et al., 2022; Shao

et al., 2024) and Diffusion Models (Fan et al., 2023). These methods result in evolving reward because the regularization term $-\alpha\mathcal{H}(\pi(\cdot|s))$ (or $-\alpha d_{\text{KL}}(\pi(\cdot|s)||\pi_{\text{ref}}(\cdot|s))$) is equivalent to a penalty term $-\alpha \log \pi(a|s)$ (or $-\alpha \log \frac{\pi(a|s)}{\pi_{\text{ref}}(a|s)}$) added to the reward function $r(s, a)$. Hence, the reward function will change because of the under-training policy $\pi(\cdot|s)$, the adaptive regularization factor α , or the change of the reference policy π_{ref} .

Besides, curriculum learning (Narvekar et al., 2020) is also closely related, as it inherently involve a sequence of evolving learning objectives (and thus rewards).

B PRELIMINARY LEMMAS

Lemma B.1 *There exist constants C_J, L_J, S_J and D_J such that for any $\theta \in \Omega(\theta), s \in \mathcal{S}, \tilde{V}_\varphi^{\pi_\theta}(s)$ is C_J -bounded, L_J -Lipschitz and S_J -smooth w.r.t. θ , and D_J -Lipschitz w.r.t. φ , where $C_J = O((1-\gamma)^{-1}), L_J = O((1-\gamma)^{-2}), S_J = O((1-\gamma)^{-3}), D_J = O((1-\gamma)^{-1})$.*

Corollary B.2 *There exist constants L_ν and S_ν such that for any $\theta \in \Omega(\theta), \nu_\rho^{\pi_\theta}(\cdot)$ is L_ν -Lipschitz and S_ν -smooth w.r.t. θ in terms of $\|\cdot\|_1$, where $L_\nu = O((1-\gamma)^{-1}), S_\nu = O((1-\gamma)^{-2})$.*

Lemma B.3 *There exist constants L_A and S_A such that A_θ is L_A -Lipschitz and S_A -smooth w.r.t. θ , where $L_A = O((1-\gamma)^{-1}), S_A = O((1-\gamma)^{-2})$.*

Lemma B.4 *There exist constants C_b, L_b, S_b and D_b such that $b_{\theta, \varphi}$ is C_b -bounded, L_b -Lipschitz and S_b -smooth w.r.t. θ and D_b -Lipschitz w.r.t. φ , where $C_b = O(1), L_b = O((1-\gamma)^{-1}), S_b = O((1-\gamma)^{-2}), D_b = O(1)$.*

Lemma B.5 *There exist constants C_ω, L_ω and S_ω and D_ω such that $\omega^*(\varphi, \theta)$ is C_ω -bounded, L_ω -Lipschitz and S_ω -smooth w.r.t. θ and D_ω -Lipschitz w.r.t. φ , where $C_\omega = O(\lambda^{-1}), L_\omega = O((1-\gamma)^{-1}\lambda^{-2}), S_\omega = O((1-\gamma)^{-2}\lambda^{-3}), D_\omega = O(\lambda^{-1})$.*

Lemma B.6 *There exists a constant C_δ such that for any $\theta \in \Omega(\theta), \varphi \in \Omega(\varphi)$ and $\|\omega\|_2 \leq C_\omega$,*

$$\mathbb{E}_{\nu, \pi_\theta, \mathcal{P}}[\hat{\delta}(s, a, s')^2] \leq C_\delta^2,$$

where $C_\delta = O(\lambda^{-1})$.

C PROOF OF MAIN THEOREM

C.1 STEP 1: BOUNDING THE ACTOR ERROR

Theorem C.1 (Actor Update) *Tate $\eta_t^\theta = \frac{c_\theta}{\sqrt{t}}$ where c_θ is a constant, then*

$$G_T \leq \frac{2L}{1-\gamma} \sqrt{G_T W_T} + O\left(\sqrt{\frac{F_T}{T}}\right) + O\left(\frac{1}{\sqrt{T}}\right) + O(\epsilon).$$

Proof We first bound the change of the objective function by the change of the reward parameter:

$$\begin{aligned} \mathbb{E}[J_{\varphi_{t+1}}(\theta_{t+1})] - J_{\varphi_t}(\theta_t) &= \mathbb{E}[J_{\varphi_{t+1}}(\theta_{t+1}) - J_{\varphi_t}(\theta_{t+1})] + \mathbb{E}[J_{\varphi_t}(\theta_{t+1})] - J_{\varphi_t}(\theta_t) \\ &\geq \mathbb{E}[-|J_{\varphi_{t+1}}(\theta_{t+1}) - J_{\varphi_t}(\theta_{t+1})|] + \mathbb{E}[J_{\varphi_t}(\theta_{t+1})] - J_{\varphi_t}(\theta_t) \\ &\geq -D_J \mathbb{E}\|\varphi_{t+1} - \varphi_t\|_2 + \mathbb{E}[J_{\varphi_t}(\theta_{t+1})] - J_{\varphi_t}(\theta_t). \end{aligned}$$

For simplicity, we denote $\xi = (s, a, s')$ and

$$\begin{aligned} h_\theta(\theta, \omega, \varphi, \xi) &= (\tilde{r}_{\varphi, \theta}(s, a) + (\gamma\phi(s') - \phi(s))^\top \omega) \nabla_\theta \log \pi_\theta(a|s), \\ \bar{h}_\theta(\theta, \omega, \varphi, \nu) &= \mathbb{E}_{\nu, \pi_\theta, \mathcal{P}}[h_\theta(\theta, \omega, \varphi, \xi)], \end{aligned}$$

thereby

$$\theta_{t+1} = \theta_t + \eta_t^\theta h_\theta(\theta_t, \omega_t, \varphi_t, \xi_t), \quad \mathbb{E}[\theta_{t+1}|\theta_t] = \theta_t + \eta_t^\theta \bar{h}_\theta(\theta_t, \omega_t, \varphi_t, \hat{\nu}_t).$$

Then we apply the Taylor expansion on $J_{\varphi_t}(\boldsymbol{\theta})$ around $\boldsymbol{\theta}_t$. The second-order term is proportional to $\eta_t^{\theta^2}$ as the gradient has bounded variance, while the first-order term can be decomposed into three parts, associated with approximation error, critic error and Markovian noise, respectively.

$$\begin{aligned}
\mathbb{E}[J_{\varphi_t}(\boldsymbol{\theta}_{t+1})] - J_{\varphi_t}(\boldsymbol{\theta}_t) &\geq \mathbb{E} \langle \nabla_{\boldsymbol{\theta}} J_{\varphi_t}(\boldsymbol{\theta}_t), \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t \rangle - \frac{S_J}{2} \mathbb{E} \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|_2^2 \\
&\geq \eta_t^{\theta} \langle \nabla_{\boldsymbol{\theta}} J_{\varphi_t}(\boldsymbol{\theta}_t), \bar{h}_{\boldsymbol{\theta}}(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \hat{\nu}_t) \rangle - \frac{S_J \eta_t^{\theta^2}}{2} \mathbb{E} \|h_{\boldsymbol{\theta}}(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \xi_t)\|_2^2 \\
&\geq \eta_t^{\theta} \langle \nabla_{\boldsymbol{\theta}} J_{\varphi_t}(\boldsymbol{\theta}_t), (1 - \gamma) \nabla_{\boldsymbol{\theta}} J_{\varphi_t}(\boldsymbol{\theta}_t) \rangle \\
&\quad + \eta_t^{\theta} \langle \nabla_{\boldsymbol{\theta}} J_{\varphi_t}(\boldsymbol{\theta}_t), \bar{h}_{\boldsymbol{\theta}}(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t^*, \boldsymbol{\varphi}_t, \nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}}) - (1 - \gamma) \nabla_{\boldsymbol{\theta}} J_{\varphi_t}(\boldsymbol{\theta}_t) \rangle \\
&\quad + \eta_t^{\theta} \langle \nabla_{\boldsymbol{\theta}} J_{\varphi_t}(\boldsymbol{\theta}_t), \bar{h}_{\boldsymbol{\theta}}(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}}) - \bar{h}_{\boldsymbol{\theta}}(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t^*, \boldsymbol{\varphi}_t, \nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}}) \rangle \\
&\quad + \eta_t^{\theta} \langle \nabla_{\boldsymbol{\theta}} J_{\varphi_t}(\boldsymbol{\theta}_t), \bar{h}_{\boldsymbol{\theta}}(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \hat{\nu}_t) - \bar{h}_{\boldsymbol{\theta}}(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}}) \rangle \\
&\quad - \frac{L^2 C_{\delta}^2 S_J \eta_t^{\theta^2}}{2} \\
&\geq \eta_t^{\theta} (1 - \gamma) \|\nabla_{\boldsymbol{\theta}} J_{\varphi_t}(\boldsymbol{\theta}_t)\|_2^2 \\
&\quad - \eta_t^{\theta} \|\nabla_{\boldsymbol{\theta}} J_{\varphi_t}(\boldsymbol{\theta}_t)\|_2 \underbrace{\|\bar{h}_{\boldsymbol{\theta}}(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t^*, \boldsymbol{\varphi}_t, \nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}}) - (1 - \gamma) \nabla_{\boldsymbol{\theta}} J_{\varphi_t}(\boldsymbol{\theta}_t)\|_2}_{I_1} \\
&\quad - \eta_t^{\theta} \|\nabla_{\boldsymbol{\theta}} J_{\varphi_t}(\boldsymbol{\theta}_t)\|_2 \underbrace{\|\bar{h}_{\boldsymbol{\theta}}(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}}) - \bar{h}_{\boldsymbol{\theta}}(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t^*, \boldsymbol{\varphi}_t, \nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}})\|_2}_{I_2} \\
&\quad - \eta_t^{\theta} \|\nabla_{\boldsymbol{\theta}} J_{\varphi_t}(\boldsymbol{\theta}_t)\|_2 \underbrace{\|\bar{h}_{\boldsymbol{\theta}}(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \hat{\nu}_t) - \bar{h}_{\boldsymbol{\theta}}(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}})\|_2}_{I_3} \\
&\quad - \frac{L^2 C_{\delta}^2 S_J \eta_t^{\theta^2}}{2} \tag{16}
\end{aligned}$$

I_1 is associated with the approximation error. Using Proposition 4.3, we have

$$\begin{aligned}
I_1 &= \left\| \mathbb{E}_{\nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}}, \pi_{\boldsymbol{\theta}_t}, \mathcal{P}} \left[\left(\left(\tilde{r}_{\boldsymbol{\varphi}_t, \boldsymbol{\theta}_t}(s, a) + \gamma \widehat{V}_{\boldsymbol{\omega}_t^*}(s') - \widehat{V}_{\boldsymbol{\omega}_t^*}(s) \right) \right. \right. \right. \\
&\quad \left. \left. \left. - \left(\tilde{r}_{\boldsymbol{\varphi}_t, \boldsymbol{\theta}_t}(s, a) + \gamma \tilde{V}_{\boldsymbol{\varphi}_t}^{\pi_{\boldsymbol{\theta}_t}}(s') - \tilde{V}_{\boldsymbol{\varphi}_t}^{\pi_{\boldsymbol{\theta}_t}}(s) \right) \right) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a|s) \right] \right\|_2 \\
&= \left\| \mathbb{E}_{\nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}}, \pi_{\boldsymbol{\theta}_t}, \mathcal{P}} \left[\left(\gamma (\widehat{V}_{\boldsymbol{\omega}_t^*}(s') - \tilde{V}_{\boldsymbol{\varphi}_t}^{\pi_{\boldsymbol{\theta}_t}}(s')) - (\widehat{V}_{\boldsymbol{\omega}_t^*}(s) - \tilde{V}_{\boldsymbol{\varphi}_t}^{\pi_{\boldsymbol{\theta}_t}}(s)) \right) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a|s) \right] \right\|_2 \\
&\leq L \sqrt{\mathbb{E}_{\nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}}, \pi_{\boldsymbol{\theta}_t}, \mathcal{P}} \left[\left(\gamma (\widehat{V}_{\boldsymbol{\omega}_t^*}(s') - \tilde{V}_{\boldsymbol{\varphi}_t}^{\pi_{\boldsymbol{\theta}_t}}(s')) - (\widehat{V}_{\boldsymbol{\omega}_t^*}(s) - \tilde{V}_{\boldsymbol{\varphi}_t}^{\pi_{\boldsymbol{\theta}_t}}(s)) \right)^2 \right]} \\
&\leq 2\sqrt{2}L\epsilon.
\end{aligned}$$

I_2 is associated with the critic error. We have

$$\begin{aligned}
I_2 &= \left\| \mathbb{E}_{\nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}}, \pi_{\boldsymbol{\theta}_t}, \mathcal{P}} \left[\left(\left(\tilde{r}_{\boldsymbol{\varphi}_t, \boldsymbol{\theta}_t}(s, a) + \gamma \widehat{V}_{\boldsymbol{\omega}_t}(s') - \widehat{V}_{\boldsymbol{\omega}_t}(s) \right) \right. \right. \right. \\
&\quad \left. \left. \left. - \left(\tilde{r}_{\boldsymbol{\varphi}_t, \boldsymbol{\theta}_t}(s, a) + \gamma \widehat{V}_{\boldsymbol{\omega}_t^*}(s') - \widehat{V}_{\boldsymbol{\omega}_t^*}(s) \right) \right) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a|s) \right] \right\|_2 \\
&= \left\| \mathbb{E}_{\nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}}, \pi_{\boldsymbol{\theta}_t}, \mathcal{P}} \left[(\gamma \phi(s') - \phi(s))^{\top} (\boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}_t}(a|s) \right] \right\|_2 \\
&\leq L \sqrt{\mathbb{E}_{\nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}}, \pi_{\boldsymbol{\theta}_t}, \mathcal{P}} \left[((\gamma \phi(s') - \phi(s))^{\top} (\boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*))^2 \right]} \\
&\leq 2L \|\boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*\|_2.
\end{aligned}$$

I_3 is associated with the Markovian noise. Using Proposition 4.9, we have

$$\begin{aligned}
I_3 &= \left\| \int_{\mathcal{S}} ds (\hat{\nu}_t(s) - \nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}}(s)) \mathbb{E}_{a, s'} \left[\hat{\delta}(s, a, s') \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}_t}(a|s) \right] \right\|_2 \\
&\leq LC_{\delta} \|\hat{\nu}_t - \nu_{\rho}^{\pi_{\boldsymbol{\theta}_t}}\|_1
\end{aligned}$$

$$\leq L^2 C_\delta^2 L_\nu \sum_{k=0}^{t-1} \gamma^{t-1-k} \eta_k^\theta + LC_\delta \gamma^t \|\rho - \nu_\rho^{\pi_{\theta_0}}\|_1.$$

Combining the above, we can derive from (16) that

$$\begin{aligned} (1-\gamma) \|\nabla_{\theta} J_{\varphi_t}(\theta_t)\|_2^2 &\leq \frac{1}{\eta_t^\theta} (\mathbb{E}[J_{\varphi_{t+1}}(\theta_{t+1})] - J_{\varphi_t}(\theta_t)) + \frac{D_J \mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2}{\eta_t^\theta} \\ &\quad + \|\nabla_{\theta} J_{\varphi_t}(\theta_t)\|_2 (I_1 + I_2 + I_3) + \frac{L^2 C_\delta^2 S_J \eta_t^\theta}{2} \\ &\leq \frac{1}{\eta_t^\theta} (\mathbb{E}[J_{\varphi_{t+1}}(\theta_{t+1})] - J_{\varphi_t}(\theta_t)) + \frac{D_J \mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2}{\eta_t^\theta} \\ &\quad + 2\sqrt{2} L L_J \epsilon + 2L \|\omega_t - \omega_t^*\|_2 \|\nabla_{\theta} J_{\varphi_t}(\theta_t)\|_2 \\ &\quad + L^2 C_\delta^2 L_J L_\nu \sum_{k=0}^{t-1} \gamma^{t-1-k} \eta_k^\theta + LC_\delta L_J \gamma^t \|\rho - \nu_\rho^{\pi_{\theta_0}}\|_1 \\ &\quad + \frac{L^2 C_\delta^2 S_J}{2} \eta_t^\theta. \end{aligned}$$

Summing over iterations, we have

$$\begin{aligned} (1-\gamma) \sum_{t=T/2}^{T-1} \mathbb{E} \|\nabla_{\theta} J_{\varphi_t}(\theta_t)\|_2^2 &\leq \underbrace{\sum_{t=T/2}^{T-1} \frac{1}{\eta_t^\theta} (\mathbb{E}[J_{\varphi_{t+1}}(\theta_{t+1})] - \mathbb{E}[J_{\varphi_t}(\theta_t)])}_{S_1} + D_J \underbrace{\sum_{t=T/2}^{T-1} \frac{\mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2}{\eta_t^\theta}}_{S_2} \\ &\quad + \underbrace{\sqrt{2} T L L_J \epsilon + 2L \sum_{t=T/2}^{T-1} \|\omega_t - \omega_t^*\|_2 \|\nabla_{\theta} J_{\varphi_t}(\theta_t)\|_2}_{S_3} \\ &\quad + L^2 C_\delta^2 L_J L_\nu \underbrace{\sum_{t=T/2}^{T-1} \sum_{k=0}^{t-1} \gamma^{t-1-k} \eta_k^\theta}_{S_4} + LC_\delta L_J \|\rho - \nu_\rho^{\pi_{\theta_0}}\|_1 \underbrace{\sum_{j=T/2}^{T-1} \gamma^j}_{S_5} \\ &\quad + \frac{L^2 C_\delta^2 S_J}{2} \underbrace{\sum_{t=T/2}^{T-1} \eta_t^\theta}_{S_6}. \end{aligned} \tag{17}$$

For S_1 , by applying the telescoping skill, we have

$$\begin{aligned} S_1 &= \sum_{t=T/2+1}^{T-1} \left(\frac{1}{\eta_{t-1}^\theta} - \frac{1}{\eta_t^\theta} \right) \mathbb{E}[J_{\varphi_t}(\theta_t)] + \frac{\mathbb{E}[J_{\varphi_T}(\theta_T)]}{\eta_{T-1}^\theta} - \frac{\mathbb{E}[J_{\varphi_{T/2}}(\theta_{T/2})]}{\eta_{T/2}^\theta} \\ &\leq \sum_{t=T/2+1}^{T-1} \left(\frac{1}{\eta_t^\theta} - \frac{1}{\eta_{t-1}^\theta} \right) C_J + \frac{C_J}{\eta_{T-1}^\theta} + \frac{C_J}{\eta_{T/2}^\theta} \\ &= \frac{2C_J}{\eta_{T-1}^\theta} \\ &= O(\sqrt{T}). \end{aligned}$$

For S_2 , by applying the Cauchy-Schwartz inequality, we have

$$S_2 \leq \sqrt{\sum_{t=T/2}^{T-1} \mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2^2} \sqrt{\sum_{t=T/2}^{T-1} \frac{1}{\eta_t^\theta}}$$

$$\begin{aligned}
&= \sqrt{\frac{TF_T}{2} \sum_{t=T/2}^{T-1} \frac{1}{\eta_t^{\theta^2}}} \\
&= O\left(\sqrt{TF_T}\right),
\end{aligned}$$

where the last equality is due to the fact that

$$\sum_{t=T/2}^{T-1} \frac{1}{\eta_t^{\theta^2}} = \frac{1}{c_{\theta}^2} \sum_{t=T/2+1}^T \frac{1}{t} = \frac{1}{c_{\theta}^2} (H_T - H_{T/2}) \sim \frac{\ln 2}{c_{\theta}^2}.$$

For S_3 , by applying the Cauchy-Schwartz inequality, we have

$$\begin{aligned}
S_3 &\leq \sqrt{\sum_{t=T/2}^{T-1} \mathbb{E} \|\omega_t - \omega_t^*\|_2^2} \sqrt{\sum_{t=T/2}^{T-1} \|\nabla_{\theta} J_{\varphi_t}(\theta_t)\|_2^2} \\
&= \frac{T}{2} \sqrt{\frac{1}{T/2} \sum_{t=T/2}^{T-1} \mathbb{E} \|\omega_t - \omega_t^*\|_2^2} \sqrt{\frac{1}{T/2} \sum_{t=T/2}^{T-1} \|\nabla_{\theta} J_{\varphi_t}(\theta_t)\|_2^2} \\
&= \frac{T}{2} \sqrt{G_T W_T}.
\end{aligned}$$

For S_4 , S_5 and S_6 , we have

$$\begin{aligned}
S_4 &\leq \sum_{t=0}^{T-1} \sum_{k=0}^{t-1} \gamma^{t-1-k} \eta_k^{\theta} = \sum_{t=0}^{T-1} \eta_t^{\theta} \sum_{j=0}^{T-t-1} \gamma^j \leq \sum_{t=0}^{T-1} \frac{\eta_t^{\theta}}{1-\gamma} \\
&= \frac{1}{c_{\theta}(1-\gamma)} \sum_{t=1}^T \frac{1}{\sqrt{t}} = O\left(\sqrt{T}\right), \\
S_5 &\leq \frac{1}{\gamma^{T/2}(1-\gamma)}, \\
S_6 &= \frac{1}{c_{\theta}} \sum_{t=T/2+1}^T \frac{1}{\sqrt{t}} = O\left(\sqrt{T}\right).
\end{aligned}$$

Plug S_1 , S_2 , S_3 , S_4 , S_5 and S_6 into (17) and divide both sides by $(1-\gamma)\frac{T}{2}$, we obtain

$$G_T \leq \frac{2L}{1-\gamma} \sqrt{G_T W_T} + O\left(\sqrt{\frac{F_T}{T}}\right) + O\left(\frac{1}{\sqrt{T}}\right) + O(\epsilon),$$

thus completes the proof. \square

1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

1134 C.2 STEP 2: BOUNDING THE CRITIC ERROR
1135

1136 **Theorem C.2 (Critic Update)** *Tate* $\eta_t^\theta = \frac{c_\theta}{\sqrt{t}}$, $\eta_t^\omega = \frac{c_\omega}{\sqrt{t}}$ where c_θ and c_ω are constants such that
1137 $\frac{c_\theta}{c_\omega} \leq \frac{\lambda}{LS_\omega}$, then
1138

$$1139 W_T \leq 2(1 - \gamma)L_\omega \frac{c_\theta}{c_\omega} \sqrt{W_T G_T} + O\left(F_T \sqrt{T}\right) + O\left(\sqrt{\frac{F_T}{T}}\right) + O\left(\frac{1}{\sqrt{T}}\right) + O(\epsilon).
1140$$

1141 *Proof* For simplicity, we denote $\xi = (s, a, s')$ and

$$1142 h_\omega(\boldsymbol{\theta}, \boldsymbol{\omega}, \boldsymbol{\varphi}, \xi) = (\tilde{r}_{\boldsymbol{\varphi}, \boldsymbol{\theta}}(s, a) + (\gamma\phi(s') - \phi(s))^\top \boldsymbol{\omega}) \phi(s),
1143 \bar{h}_\omega(\boldsymbol{\theta}, \boldsymbol{\omega}, \boldsymbol{\varphi}, \nu) = \mathbb{E}_{\nu, \pi_{\boldsymbol{\theta}, \mathcal{P}}} [h_\omega(\boldsymbol{\theta}, \boldsymbol{\omega}, \boldsymbol{\varphi}, \xi)].
1144$$

1145 According to the critic update rule (9), we have
1146

$$1147 \|\boldsymbol{\omega}_{t+1} - \boldsymbol{\omega}_{t+1}^*\|_2^2 = \|\Pi_{C_\omega}(\boldsymbol{\omega}_t + \eta_t^\omega h_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \xi_t)) - \Pi_{C_\omega}(\boldsymbol{\omega}_{t+1}^*)\|_2^2
1148 \leq \|\boldsymbol{\omega}_t + \eta_t^\omega h_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \xi_t) - \boldsymbol{\omega}_{t+1}^*\|_2^2
1149 = \|(\boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*) + \eta_t^\omega h_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \xi_t) + (\boldsymbol{\omega}_t^* - \boldsymbol{\omega}_{t+1}^*)\|_2^2
1150 = \|\boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*\|_2^2 + \|\eta_t^\omega h_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \xi_t) + (\boldsymbol{\omega}_t^* - \boldsymbol{\omega}_{t+1}^*)\|_2^2
1151 + 2\langle \boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*, \eta_t^\omega h_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \xi_t) + (\boldsymbol{\omega}_t^* - \boldsymbol{\omega}_{t+1}^*) \rangle
1152 \leq \|\boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*\|_2^2 + 2\eta_t^{\omega^2} \|h_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \xi_t)\|_2^2 + 2\|\boldsymbol{\omega}_t^* - \boldsymbol{\omega}_{t+1}^*\|_2^2
1153 + 2\eta_t^\omega \langle \boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*, h_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \xi_t) \rangle + 2\langle \boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*, \boldsymbol{\omega}_t^* - \boldsymbol{\omega}_{t+1}^* \rangle. \quad (18)
1154
1155
1156
1157$$

1158 To capture the evolution of the critic error $\|\boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*\|_2^2$, we need to bound the other four terms on
1159 the right-hand side of (18). By taking the expectation, the two quadratic terms can be bounded by

$$1160 \mathbb{E} \|h_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \xi_t)\|_2^2 = \mathbb{E}_{\hat{\nu}_t, \pi_{\boldsymbol{\theta}_t, \mathcal{P}}} \left\| \hat{\delta}(s_t, a_t, s'_t) \phi(s_t) \right\|_2^2 \leq C_\delta^2 \quad (19)
1161$$

1162 and

$$1163 \mathbb{E} \|\boldsymbol{\omega}_t^* - \boldsymbol{\omega}_{t+1}^*\|_2^2 \leq \mathbb{E} \left[(L_\omega \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|_2 + D_\omega \|\boldsymbol{\varphi}_{t+1} - \boldsymbol{\varphi}_t\|_2)^2 \right]
1164 \leq 2L_\omega^2 \eta_t^{\theta^2} \mathbb{E} \|h_\theta(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \xi_t)\|_2^2 + 2D_\omega^2 \mathbb{E} \|\boldsymbol{\varphi}_{t+1} - \boldsymbol{\varphi}_t\|_2^2
1165 \leq 2L^2 C_\delta^2 L_\omega^2 \eta_t^{\theta^2} + 2D_\omega^2 \mathbb{E} \|\boldsymbol{\varphi}_{t+1} - \boldsymbol{\varphi}_t\|_2^2. \quad (20)
1166
1167$$

1168 The expectation of first inner-product term can be decomposed as the follows:

$$1169 \mathbb{E} \langle \boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*, h_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \xi_t) \rangle = \langle \boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*, \bar{h}_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \hat{\nu}_t) \rangle
1170 = \underbrace{\langle \boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*, \bar{h}_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t^*, \boldsymbol{\varphi}_t, \nu_\rho^{\pi_{\boldsymbol{\theta}_t}}) \rangle}_{J_1}
1171 + \underbrace{\langle \boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*, \bar{h}_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \nu_\rho^{\pi_{\boldsymbol{\theta}_t}}) - \bar{h}_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t^*, \boldsymbol{\varphi}_t, \nu_\rho^{\pi_{\boldsymbol{\theta}_t}}) \rangle}_{J_2}
1172 + \underbrace{\langle \boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*, \bar{h}_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \hat{\nu}_t) - \bar{h}_\omega(\boldsymbol{\theta}_t, \boldsymbol{\omega}_t, \boldsymbol{\varphi}_t, \nu_\rho^{\pi_{\boldsymbol{\theta}_t}}) \rangle}_{J_3}.
1173
1174
1175
1176
1177
1178$$

1179 According to the definition of $\boldsymbol{\omega}_t^*$, it should be a stationary point of the update rule, hence we have
1180

$$1181 J_1 = 0.$$

1182 J_2 is associated with the critic error. We have

$$1183 J_2 = \left\langle \boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*, \mathbb{E}_{\nu_\rho^{\pi_{\boldsymbol{\theta}_t}}, \pi_{\boldsymbol{\theta}_t}, \mathcal{P}} [(\gamma\phi(s') - \phi(s))^\top (\boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*) \phi(s)] \right\rangle
1184 = \langle \boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*, -A_{\boldsymbol{\theta}_t}(\boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*) \rangle
1185 \leq -\lambda \|\boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*\|_2^2.
1186
1187$$

J_3 is associated with the Markovian noise. Using Lemma 4.9, we have

$$\begin{aligned}
1188 \quad J_3 &\leq \|\omega_t - \omega_t^*\|_2 \|\bar{h}_\omega(\theta_t, \omega_t, \varphi_t, \hat{\nu}_t) - \bar{h}_\omega(\theta_t, \omega_t, \varphi_t, \nu_\rho^{\pi_{\theta_t}})\|_2 \\
1189 \quad &\leq 2C_\omega \left\| \int_{\mathcal{S}} ds (\hat{\nu}_t(s) - \nu_\rho^{\pi_{\theta_t}}(s)) \mathbb{E}_{a \sim \pi_{\theta_t}(\cdot|s), s' \sim \mathcal{P}(\cdot|s,a)} [\hat{\delta}(s, a, s') \phi(s)] \right\|_2 \\
1190 \quad &\leq 2C_\omega C_\delta \|\hat{\nu}_t - \nu_\rho^{\pi_{\theta_t}}\|_1 \\
1191 \quad &\leq 2LC_\omega C_\delta^2 L_\nu \sum_{k=0}^{t-1} \gamma^{t-1-k} \eta_k^\theta + 2C_\omega C_\delta \gamma^t \|\rho - \nu_\rho^{\pi_{\theta_0}}\|_1.
\end{aligned}$$

Hence, we have

$$\begin{aligned}
1192 \quad \mathbb{E} \langle \omega_t - \omega_t^*, h_\omega(\theta_t, \omega_t, \varphi_t, \xi_t) \rangle &\leq -\lambda \|\omega_t - \omega_t^*\|_2^2 + 2LC_\omega C_\delta^2 L_\nu \sum_{k=0}^{t-1} \gamma^{t-1-k} \eta_k^\theta + 2C_\omega C_\delta \gamma^t \|\rho - \nu_\rho^{\pi_{\theta_0}}\|_1. \\
1193 \quad & \\
1194 \quad & \\
1195 \quad & \\
1196 \quad & \\
1197 \quad & \\
1198 \quad & \\
1199 \quad & \\
1200 \quad & \tag{21}
\end{aligned}$$

The analysis of the last term in (18) is similar to the analysis of the actor error. We first leverage the Lipschitz continuity of $\omega^*(\varphi, \theta)$ with respect to φ , then apply the Taylor expansion of $\omega^*(\varphi, \theta)$ with respect to θ around θ_t :

$$\begin{aligned}
1201 \quad \mathbb{E} \langle \omega_t - \omega_t^*, \omega_t^* - \omega_{t+1}^* \rangle &= \mathbb{E} \langle \omega_t - \omega_t^*, \omega^*(\varphi_t, \theta_t) - \omega^*(\varphi_{t+1}, \theta_t) \rangle \\
1202 \quad &\quad + \mathbb{E} \langle \omega_t - \omega_t^*, \omega^*(\varphi_{t+1}, \theta_t) - \omega^*(\varphi_{t+1}, \theta_{t+1}) \rangle \\
1203 \quad &= \mathbb{E} \langle \omega_t - \omega_t^*, \omega^*(\varphi_t, \theta_t) - \omega^*(\varphi_{t+1}, \theta_t) \rangle \\
1204 \quad &\quad + \mathbb{E} \langle \omega_t - \omega_t^*, \omega^*(\varphi_t, \theta_t) - \omega^*(\varphi_{t+1}, \theta_t) - \nabla_\theta \omega^*(\varphi_{t+1}, \theta_t)^\top (\theta_t - \theta_{t+1}) \rangle \\
1205 \quad &\quad + \mathbb{E} \langle \omega_t - \omega_t^*, \nabla_\theta \omega^*(\varphi_{t+1}, \theta_t)^\top (\theta_t - \theta_{t+1}) \rangle \\
1206 \quad &\leq D_\omega \|\omega_t - \omega_t^*\|_2 \mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2 + \frac{S_\omega}{2} \|\omega_t - \omega_t^*\|_2 \mathbb{E} \|\theta_{t+1} - \theta_t\|_2^2 \\
1207 \quad &\quad + \mathbb{E} \langle \omega_t - \omega_t^*, \nabla_\theta \omega^*(\varphi_{t+1}, \theta_t)^\top (\theta_t - \theta_{t+1}) \rangle \\
1208 \quad &\leq L^2 C_\delta^2 C_\omega S_\omega \eta_t^{\theta^2} + 2C_\omega D_\omega \mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2 \\
1209 \quad &\quad + \underbrace{\mathbb{E} \langle \omega_t - \omega_t^*, \nabla_\theta \omega^*(\varphi_{t+1}, \theta_t)^\top (\theta_t - \theta_{t+1}) \rangle}_{I},
\end{aligned}$$

where the last inequality uses the facts that

$$\begin{aligned}
1210 \quad \mathbb{E} \|\theta_t - \theta_{t+1}\|_2^2 &\leq L^2 C_\delta^2 \eta_t^{\theta^2} \quad \text{and} \quad \|\omega_t - \omega_t^*\|_2^2 \leq 2C_\omega.
\end{aligned}$$

The remaining inner product term I can then be decomposed into terms related to I_1 , I_2 and I_3 .

$$\begin{aligned}
1211 \quad I &= -\eta_t^\theta \langle \omega_t - \omega_t^*, \nabla_\theta \omega^*(\varphi_{t+1}, \theta_t)^\top \bar{h}_\theta(\theta_t, \omega_t, \varphi_t, \xi_t) \rangle \\
1212 \quad &= -(1-\gamma) \eta_t^\theta \langle \omega_t - \omega_t^*, \nabla_\theta \omega^*(\varphi_{t+1}, \theta_t)^\top \nabla_\theta J_{\varphi_t}(\theta_t) \rangle \\
1213 \quad &\quad - \eta_t^\theta \langle \omega_t - \omega_t^*, \nabla_\theta \omega^*(\varphi_{t+1}, \theta_t)^\top (\bar{h}_\theta(\theta_t, \omega_t^*, \varphi_t, \nu_\rho^{\pi_{\theta_t}}) - (1-\gamma) \nabla_\theta J_{\varphi_t}(\theta_t)) \rangle \\
1214 \quad &\quad - \eta_t^\theta \langle \omega_t - \omega_t^*, \nabla_\theta \omega^*(\varphi_{t+1}, \theta_t)^\top (\bar{h}_\theta(\theta_t, \omega_t, \varphi_t, \nu_\rho^{\pi_{\theta_t}}) - \bar{h}_\theta(\theta_t, \omega_t^*, \varphi_t, \nu_\rho^{\pi_{\theta_t}})) \rangle \\
1215 \quad &\quad - \eta_t^\theta \langle \omega_t - \omega_t^*, \nabla_\theta \omega^*(\varphi_{t+1}, \theta_t)^\top (\bar{h}_\theta(\theta_t, \omega_t, \varphi_t, \hat{\nu}_t) - \bar{h}_\theta(\theta_t, \omega_t, \varphi_t, \nu_\rho^{\pi_{\theta_t}})) \rangle \\
1216 \quad &\leq L^2 C_\delta^2 C_\omega S_\omega \eta_t^{\theta^2} + 2C_\omega D_\omega \mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2 \\
1217 \quad &\quad + (1-\gamma) S_\omega \eta_t^\theta \|\omega_t - \omega_t^*\|_2 \|\nabla_\theta J_{\varphi_t}(\theta_t)\|_2 \\
1218 \quad &\quad + S_\omega \eta_t^\theta \|\omega_t - \omega_t^*\|_2 \underbrace{\|\bar{h}_\theta(\theta_t, \omega_t^*, \varphi_t, \nu_\rho^{\pi_{\theta_t}}) - (1-\gamma) \nabla_\theta J_{\varphi_t}(\theta_t)\|}_{I_1} \\
1219 \quad &\quad + S_\omega \eta_t^\theta \|\omega_t - \omega_t^*\|_2 \underbrace{\|\bar{h}_\theta(\theta_t, \omega_t, \varphi_t, \nu_\rho^{\pi_{\theta_t}}) - \bar{h}_\theta(\theta_t, \omega_t^*, \varphi_t, \nu_\rho^{\pi_{\theta_t}})\|}_{I_2} \\
1220 \quad &\quad + S_\omega \eta_t^\theta \|\omega_t - \omega_t^*\|_2 \underbrace{\|\bar{h}_\theta(\theta_t, \omega_t, \varphi_t, \hat{\nu}_t) - \bar{h}_\theta(\theta_t, \omega_t, \varphi_t, \nu_\rho^{\pi_{\theta_t}})\|}_{I_3},
\end{aligned}$$

where

$$\begin{aligned}
1242 \quad I_1 &\leq 2\sqrt{2}L\epsilon, \quad I_2 \leq 2L\|\omega_t - \omega_t^*\|_2, \\
1243 \\
1244 \quad I_3 &\leq L^2 C_\delta^2 L_\nu \sum_{k=0}^{t-1} \gamma^{t-1-k} \eta_k^\theta + LC_\delta \gamma^t \|\rho - \nu_\rho^{\pi_{\theta_0}}\|_1. \\
1245 \\
1246
\end{aligned}$$

1247 Hence, we have

$$\begin{aligned}
1248 \quad \langle \omega_t - \omega_t^*, \omega_t^* - \omega_{t+1}^* \rangle &\leq L^2 C_\delta^2 C_\omega S_\omega \eta_t^{\theta^2} + 2C_\omega D_\omega \|\varphi_{t+1} - \varphi_t\|_2 \\
1249 &\quad + (1 - \gamma) S_\omega \eta_t^\theta \|\omega_t - \omega_t^*\|_2 \|\nabla_{\theta} J_{\varphi_t}(\theta_t)\|_2 \\
1250 &\quad + 4\sqrt{2}LC_\omega S_\omega \epsilon \eta_t^\theta + 2LS_\omega \eta_t^\theta \|\omega_t - \omega_t^*\|_2^2 \\
1251 &\quad + 2L^2 C_\omega C_\delta^2 L_\nu S_\omega \eta_t^\theta \sum_{k=0}^{t-1} \gamma^{t-1-k} \eta_k^\theta + 2LC_\omega C_\delta S_\omega \eta_t^\theta \gamma^t \|\rho - \nu_\rho^{\pi_{\theta_0}}\|_1. \\
1252 \\
1253 \\
1254 \\
1255 \end{aligned} \tag{22}$$

1256 Plugging (19), (20), (21) and (22) into (18) gives

$$\begin{aligned}
1257 \quad \mathbb{E} \|\omega_{t+1} - \omega_{t+1}^*\|_2^2 &\leq (1 - 2\lambda \eta_t^\omega + 4LS_\omega \eta_t^\theta) \|\omega_t - \omega_t^*\|_2^2 \\
1258 &\quad + 2(1 - \gamma) S_\omega \eta_t^\theta \|\omega_t - \omega_t^*\|_2 \|\nabla_{\theta} J_t(\theta_t)\|_2 \\
1259 &\quad + 2D_\omega^2 \mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2^2 + 4C_\omega D_\omega \mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2 \\
1260 &\quad + 4LC_\omega C_\delta^2 L_\nu (LS_\omega \eta_t^\theta + \eta_t^\omega) \sum_{k=0}^{t-1} \gamma^{t-1-k} \eta_k^\theta \\
1261 &\quad + 4C_\omega C_\delta (LS_\omega \eta_t^\theta + \eta_t^\omega) \gamma^t \|\rho - \nu_\rho^{\pi_{\theta_0}}\|_1 \\
1262 &\quad + 4L^2 C_\delta^2 L_\omega^2 \eta_t^{\theta^2} + 8\sqrt{2}LC_\omega S_\omega \epsilon \eta_t^\theta. \\
1263 \\
1264 \\
1265 \\
1266 \\
1267
\end{aligned}$$

1268 Note that $\frac{\eta_t^\theta}{\eta_t^\omega} = \frac{c_\theta}{c_\omega} \leq \frac{\lambda}{LS_\omega}$, thus

$$\begin{aligned}
1270 \quad \lambda \sum_{t=T/2}^{T-1} \mathbb{E} \|\omega_t - \omega_t^*\|_2^2 &\leq \underbrace{\sum_{t=T/2}^{T-1} \frac{1}{\eta_t^\omega} \left(\mathbb{E} \|\omega_t - \omega_t^*\|_2^2 - \mathbb{E} \|\omega_{t+1} - \omega_{t+1}^*\|_2^2 \right)}_{S_1} \\
1271 &\quad + 2(1 - \gamma) L_\omega \underbrace{\frac{c_\theta}{c_\omega} \sum_{t=T/2}^{T-1} \|\omega_t - \omega_t^*\|_2 \|\nabla_{\theta} J_t(\theta_t)\|_2}_{S_2} \\
1272 &\quad + 2D_\omega^2 \underbrace{\sum_{t=T/2}^{T-1} \frac{\mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2^2}{\eta_t^\omega}}_{S_3} + 2C_\omega D_\omega \underbrace{\sum_{t=T/2}^{T-1} \frac{\mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2}{\eta_t^\omega}}_{S_4} \\
1273 &\quad + (1 + \lambda) LC_\omega C_\delta^2 L_\nu \underbrace{\sum_{t=T/2}^{T-1} \sum_{k=0}^{t-1} \gamma^{t-1-k} \eta_k^\theta}_{S_5} \\
1274 &\quad + (1 + \lambda) C_\omega C_\delta L_J \|\rho - \nu_\rho^{\pi_{\theta_0}}\|_1 \underbrace{\sum_{j=T/2}^{T-1} \gamma^j}_{S_6} \\
1275 &\quad + 4L^2 C_\delta^2 L_\omega^2 \underbrace{\frac{c_\theta}{c_\omega} \sum_{t=T/2}^{T-1} \eta_t^\theta}_{S_7} + 8\sqrt{2}LC_\omega S_\omega \frac{c_\theta}{c_\omega} \epsilon \\
1276 \\
1277 \\
1278 \\
1279 \\
1280 \\
1281 \\
1282 \\
1283 \\
1284 \\
1285 \\
1286 \\
1287 \\
1288 \\
1289 \\
1290 \\
1291 \\
1292 \\
1293 \\
1294 \\
1295 \end{aligned} \tag{23}$$

For S_1 , by applying the telescoping skill, we have

$$\begin{aligned}
S_1 &= \sum_{t=T/2+1}^{T-1} \left(\frac{1}{\eta_t^\omega} - \frac{1}{\eta_{t-1}^\omega} \right) \mathbb{E} \|\omega_t - \omega_t^*\|_2^2 - \frac{\mathbb{E} \|\omega_T - \omega_T^*\|_2^2}{\eta_{T-1}^\omega} + \frac{\mathbb{E} \|\omega_{T/2} - \omega_{T/2}^*\|_2^2}{\eta_{T/2}^\omega} \\
&\leq \sum_{t=T/2+1}^{T-1} \left(\frac{1}{\eta_t^\omega} - \frac{1}{\eta_{t-1}^\omega} \right) 2C_\omega + \frac{2C_\omega}{\eta_{T-1}^\omega} + \frac{2C_\omega}{\eta_{T/2}^\omega} \\
&= \frac{4C_\omega}{\eta_{T-1}^\omega} \\
&= O(\sqrt{T}).
\end{aligned}$$

For S_2 , by applying the Cauchy-Schwartz inequality, we have

$$\begin{aligned}
S_2 &\leq \sqrt{\sum_{t=T/2}^{T-1} \mathbb{E} \|\omega_t - \omega_t^*\|_2^2} \sqrt{\sum_{t=T/2}^{T-1} \|\nabla_{\theta} J_{\varphi_t}(\theta_t)\|_2^2} \\
&= \frac{T}{2} \sqrt{\frac{1}{T/2} \sum_{t=T/2}^{T-1} \mathbb{E} \|\omega_t - \omega_t^*\|_2^2} \sqrt{\frac{1}{T/2} \sum_{t=T/2}^{T-1} \|\nabla_{\theta} J_{\varphi_t}(\theta_t)\|_2^2} \\
&= \frac{T}{2} \sqrt{W_T G_T}.
\end{aligned}$$

For S_4 , note that η_t^ω is decreasing as t grows, we have

$$S_4 = \sum_{t=T/2}^{T-1} \frac{\mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2^2}{\eta_t^\omega} \leq \frac{1}{\eta_{T-1}^\omega} \sum_{t=T/2}^{T-1} \mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2^2 = O(F_T T \sqrt{T})$$

For S_5 , by applying the Cauchy-Schwartz inequality, we have

$$\begin{aligned}
S_5 &\leq \sqrt{\sum_{t=T/2}^{T-1} \mathbb{E} \|\varphi_{t+1} - \varphi_t\|_2^2} \sqrt{\sum_{t=T/2}^{T-1} \frac{1}{\eta_t^{\omega^2}}} \\
&= \sqrt{\frac{TF_T}{2} \sum_{t=T/2}^{T-1} \frac{1}{\eta_t^{\omega^2}}} \\
&= O(\sqrt{TF_T}),
\end{aligned}$$

where the last equality is due to the fact that

$$\sum_{t=T/2}^{T-1} \frac{1}{\eta_t^{\omega^2}} = \frac{1}{c_\omega^2} \sum_{t=T/2+1}^T \frac{1}{t} = \frac{1}{c_\omega^2} (H_T - H_{T/2}) \sim \frac{\ln 2}{c_\omega^2}.$$

For S_5 , S_6 and S_7 , we have

$$\begin{aligned}
S_5 &\leq \sum_{t=0}^{T-1} \sum_{k=0}^{t-1} \gamma^{t-1-k} \eta_k^\theta = \sum_{t=0}^{T-1} \eta_t^\theta \sum_{j=0}^{T-t-1} \gamma^j \leq \sum_{t=0}^{T-1} \frac{\eta_t^\theta}{1-\gamma} \\
&= \frac{1}{c_\theta(1-\gamma)} \sum_{t=1}^T \frac{1}{\sqrt{t}} = O(\sqrt{T}), \\
S_6 &\leq \frac{1}{\gamma^{T/2}(1-\gamma)}, \\
S_7 &= \frac{1}{c_\theta} \sum_{t=T/2+1}^T \frac{1}{\sqrt{t}} = O(\sqrt{T}).
\end{aligned}$$

1350 Plug $S_1, S_2, S_3, S_4, S_5, S_6$ and S_7 into (23) and divide both sides by $\frac{\lambda T}{2}$, we obtain
 1351

$$1352 \quad W_T \leq 2(1-\gamma)L\omega \frac{c\theta}{c_\omega} \sqrt{W_T G_T} + O\left(F_T \sqrt{T}\right) + O\left(\sqrt{\frac{F_T}{T}}\right) + O\left(\frac{1}{\sqrt{T}}\right) + O(\epsilon),$$

1353 thus completes the proof.
 1354

□

1355 C.3 STEP 3: SOLVING THE SYSTEM OF INEQUALITIES

1356 **Proof of Theorem 4.7** According to Theorem C.1 and Theorem C.2, we have
 1357

$$1358 \quad (1-\gamma)G_T \leq 2L\sqrt{G_T W_T} + O\left(\sqrt{\frac{F_T}{T}}\right) + O\left(\frac{1}{\sqrt{T}}\right) + O(\epsilon), \quad (24)$$

$$1359 \quad \frac{1}{1-\gamma}W_T \leq 2L\omega \frac{c\theta}{c_\omega} \sqrt{G_T W_T} + O\left(\sqrt{T}F_T\right) + O\left(\frac{F_T}{T}\right) + O\left(\frac{1}{\sqrt{T}}\right) + O(\epsilon). \quad (25)$$

1360 Note that

$$1361 \quad 2\sqrt{G_T W_T} = 2\sqrt{\frac{1-\gamma}{2L}G_T \cdot \frac{2L}{1-\gamma}W_T} \leq \frac{1-\gamma}{2L}G_T + \frac{2L}{1-\gamma}W_T. \quad (26)$$

1362 Plug (26) into (24), we have

$$1363 \quad \frac{1-\gamma}{2L}G_T \leq \frac{2L}{1-\gamma}W_T + O\left(\sqrt{\frac{F_T}{T}}\right) + O\left(\frac{1}{\sqrt{T}}\right) + O(\epsilon). \quad (27)$$

1364 Combining (26) and (27), we have

$$1365 \quad 2\sqrt{G_T W_T} \leq \frac{4L}{1-\gamma}W_T + O\left(\sqrt{\frac{F_T}{T}}\right) + O\left(\frac{1}{\sqrt{T}}\right) + O(\epsilon). \quad (28)$$

1366 Plug (28) into (25), we have

$$1367 \quad \frac{1-8LL\omega \frac{c\theta}{c_\omega}}{1-\gamma}W_T \leq O\left(\sqrt{T}F_T\right) + O\left(\frac{F_T}{T}\right) + O\left(\frac{1}{\sqrt{T}}\right) + O(\epsilon).$$

1368 Therefore, when $\frac{c\theta}{c_\omega} \leq \frac{1}{16LL\omega}$,

$$1369 \quad W_T = O\left(\frac{1}{\sqrt{T}}\right) + O\left(\sqrt{T}F_T\right) + O\left(\frac{F_T}{T}\right) + O(\epsilon).$$

1370 Combined with (27), we have

$$1371 \quad G_T = O\left(\frac{1}{\sqrt{T}}\right) + O\left(\sqrt{T}F_T\right) + O\left(\frac{F_T}{T}\right) + O(\epsilon).$$

1404 **D PROOF OF PROPOSITIONS, PRELIMINARY LEMMAS AND COROLLARIES**

1405 **Proof of Proposition 4.2** For any vector x ,

1406

1407

1408
$$x^\top A_\theta x = x^\top \mathbb{E}_{\nu_\rho^{\pi_\theta}, \pi_\theta, \mathcal{P}} \left[\phi(s) (\phi(s) - \gamma \phi(s'))^\top \right] x$$

1409
$$= x^\top \left(\mathbb{E}_s [\phi(s) \phi(s)^\top] - \gamma \mathbb{E}_{s, s'} [\phi(s) \phi(s')^\top] \right) x$$

1410

1411
$$= \mathbb{E}_s [x^\top \phi(s) \phi(s)^\top x] - \gamma \mathbb{E}_{s, s'} [x^\top \phi(s) \phi(s')^\top x]$$

1412

1413 According to the Cauchy-Schwartz inequality,

1414
$$\mathbb{E}_{s, s'} [x^\top \phi(s) \phi(s')^\top x] \leq \sqrt{\mathbb{E}_s [x^\top \phi(s) \phi(s)^\top x]} \sqrt{\mathbb{E}_{s'} [x^\top \phi(s') \phi(s')^\top x]}.$$

1415

1416 Note that

1417

1418
$$\Pr(s' = x) = \frac{\nu_\rho^{\pi_\theta}(x) - (1 - \gamma)\rho(x)}{\gamma} \leq \frac{\nu_\rho^{\pi_\theta}(x)}{\gamma},$$

1419

1420 so

1421
$$\mathbb{E}_{s'} [x^\top \phi(s') \phi(s')^\top x] \leq \frac{1}{\gamma} \mathbb{E}_s [x^\top \phi(s) \phi(s)^\top x],$$

1422

1423 and hence

1424

1425
$$x^\top A_\theta x \geq \mathbb{E}_s [x^\top \phi(s) \phi(s)^\top x] - \gamma \sqrt{\mathbb{E}_s [x^\top \phi(s) \phi(s)^\top x]} \sqrt{\frac{1}{\gamma} \mathbb{E}_s [x^\top \phi(s) \phi(s)^\top x]}$$

1426

1427
$$= (1 - \sqrt{\gamma}) x^\top \Sigma_\theta x$$

1428
$$\geq (1 - \sqrt{\gamma}) \lambda_\Sigma \|x\|_2^2$$

1429

1430 Therefore, A_θ is positive definite with singular values lower-bounded by $\lambda = (1 - \sqrt{\gamma}) \lambda_\Sigma$.

1431

1432 **Proof of Proposition 4.3**

1433

1434
$$LHS = \sqrt{\mathbb{E}_{\nu_\rho^{\pi_\theta}, \pi_\theta, \mathcal{P}} \left[\left(\left(\gamma \widehat{V}_{\omega^*}(s') - \widehat{V}_{\omega^*}(s) \right) - \left(\gamma \widetilde{V}^{\pi_\theta}(s') - \widetilde{V}^{\pi_\theta}(s) \right) \right)^2 \right]}$$

1435

1436
$$\leq \sqrt{\mathbb{E}_{\nu_\rho^{\pi_\theta}, \pi_\theta, \mathcal{P}} \left[2 \left(\gamma \left(\widehat{V}_{\omega^*}(s') - \widetilde{V}^{\pi_\theta}(s') \right) \right)^2 + 2 \left(\widehat{V}_{\omega^*}(s) - \widetilde{V}^{\pi_\theta}(s) \right)^2 \right]}$$

1437

1438
$$\leq \sqrt{2 \mathbb{E}_s \left[\left(\widehat{V}_{\omega^*}(s) - \widetilde{V}^{\pi_\theta}(s) \right)^2 \right] + 2\gamma^2 \mathbb{E}_{s'} \left[\left(\widehat{V}_{\omega^*}(s') - \widetilde{V}^{\pi_\theta}(s') \right)^2 \right]}$$

1439

1440
$$\leq \sqrt{2} \left(\underbrace{\sqrt{\mathbb{E}_s \left[\left(\widehat{V}_{\omega^*}(s) - \widetilde{V}^{\pi_\theta}(s) \right)^2 \right]}}_{I_1} + \gamma \underbrace{\sqrt{\mathbb{E}_{s'} \left[\left(\widehat{V}_{\omega^*}(s') - \widetilde{V}^{\pi_\theta}(s') \right)^2 \right]}}_{I_2} \right)$$

1441

1442

1443

1444

1445

1446

1447

1448 According to the definition of ϵ (13), $I_1 \leq \epsilon$.

1449 For I_2 , note that

1450

1451
$$\Pr(s' = x) = \frac{\nu_\rho^{\pi_\theta}(x) - (1 - \gamma)\rho(x)}{\gamma} \leq \frac{\nu_\rho^{\pi_\theta}(x)}{\gamma},$$

1452

1453 so

1454
$$I_2 \leq \gamma \sqrt{\frac{1}{\gamma} \mathbb{E}_s \left[\left(\widehat{V}_{\omega^*}(s) - \widetilde{V}^{\pi_\theta}(s) \right)^2 \right]} \leq \sqrt{\gamma} \epsilon \leq \epsilon.$$

1455

1456

1457

Therefore, $LHS \leq 2\sqrt{2}\epsilon$.

Proof of Proposition 4.5 For any $\theta_1, \theta_2 \in \Omega(\theta)$, let

$$f(a) = \begin{cases} 1 & , \pi_{\theta_1}(a|s) \geq \pi_{\theta_2}(a|s) \\ -1 & , \text{otherwise} \end{cases},$$

then

$$\|\pi_{\theta_1}(\cdot|s) - \pi_{\theta_2}(\cdot|s)\|_1 = \mathbb{E}_{a \sim \pi_{\theta_1}(\cdot|s)}[f(a)] - \mathbb{E}_{a \sim \pi_{\theta_2}(\cdot|s)}[f(a)].$$

Note that

$$\begin{aligned} \|\nabla_{\theta} \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[f(a)]\|_2 &= \left\| \nabla_{\theta} \int_{\mathcal{A}} \pi_{\theta}(a|s) f(a) da \right\|_2 \\ &= \left\| \int_{\mathcal{A}} \nabla_{\theta} \pi_{\theta}(a|s) f(a) da \right\|_2 \\ &= \left\| \int_{\mathcal{A}} \pi_{\theta}(a|s) \nabla_{\theta} \log \pi_{\theta}(a|s) f(a) da \right\|_2 \\ &= \left\| \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[\nabla_{\theta} \log \pi_{\theta}(a|s) f(a)] \right\|_2 \\ &\leq \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[\|\nabla_{\theta} \log \pi_{\theta}(a|s)\|_2 |f(a)|] \\ &\leq L. \end{aligned}$$

Therefore,

$$\|\pi_{\theta_1}(\cdot|s) - \pi_{\theta_2}(\cdot|s)\|_1 \leq L \|\theta_1 - \theta_2\|.$$

Proof of Corollary 4.8 Assume that $\mathbb{E}\|h_{\varphi}(t)\|_2^2 \leq C_{\varphi}^2$ and $\eta_t^{\varphi} = \frac{c_{\varphi}}{\sqrt{t}}$, we have

$$\begin{aligned} F_T &= \frac{1}{T/2} \sum_{t=T/2}^{T-1} \mathbb{E}\|\varphi_{t+1} - \varphi_t\|_2^2 \\ &= \frac{1}{T/2} \sum_{t=T/2}^{T-1} \eta_t^{\varphi^2} \mathbb{E}\|h_{\varphi}(t)\|_2^2 \\ &\leq \frac{C_{\varphi}^2}{T/2} \sum_{t=T/2}^{T-1} \frac{c_{\varphi}^2}{t} \\ &= O(1/T) \end{aligned}$$

Hence, the terms $O(F_T \sqrt{T})$ and $O(\sqrt{F_T/T})$ are both dominated by $O(1/\sqrt{T})$, leading to an overall $O(1/\sqrt{T}) + O(\epsilon)$ bound.

Proof of Proposition 4.9 We abuse the notation $\widehat{\mathcal{P}}_{\theta} : \Delta(S) \rightarrow \Delta(S)$ to denote an operator that acts on a state distribution ν , defined by

$$\begin{aligned} (\widehat{\mathcal{P}}_{\theta} \nu)(s') &= \int_S ds \int_{\mathcal{A}} da \nu(s) \pi_{\theta}(a|s) \widehat{\mathcal{P}}(s'|s, a) \\ &= \gamma \int_S ds \int_{\mathcal{A}} da \nu(s) \pi_{\theta}(a|s) \mathcal{P}(s'|s, a) + (1 - \gamma) \rho(s') \end{aligned}$$

Then, $\widehat{\mathcal{P}}_{\theta}$ is a contraction mapping and $\nu_{\rho}^{\pi_{\theta}}$ is the unique fix point of it. Formally, $\forall \nu_1, \nu_2 \in \Delta(S)$, we have

$$\begin{aligned} \|\widehat{\mathcal{P}}_{\theta} \nu_1 - \widehat{\mathcal{P}}_{\theta} \nu_2\|_1 &= \int_S ds' \left| (\widehat{\mathcal{P}}_{\theta} \nu_1)(s') - (\widehat{\mathcal{P}}_{\theta} \nu_2)(s') \right| \\ &= \gamma \int_S ds' \left| \int_S ds \int_{\mathcal{A}} da (\nu_1(s) - \nu_2(s)) \pi_{\theta}(a|s) \mathcal{P}(s'|s, a) \right| \\ &\leq \gamma \int_S ds |\nu_1(s) - \nu_2(s)| \int_S ds' \int_{\mathcal{A}} da \pi_{\theta}(a|s) \mathcal{P}(s'|s, a) \end{aligned}$$

$$= \gamma \|\nu_1 - \nu_2\|_1,$$

and

$$(\widehat{\mathcal{P}}_{\theta} \nu_{\rho}^{\pi_{\theta}})(s) = \nu_{\rho}^{\pi_{\theta}}(s), \forall s \in \mathcal{S}.$$

Therefore,

$$\begin{aligned} \mathbb{E} \|\hat{\nu}_t - \nu_{\rho}^{\pi_{\theta_t}}\|_1 &\leq \mathbb{E} \|\hat{\nu}_t - \nu_{\rho}^{\pi_{\theta_{t-1}}}\|_1 + \mathbb{E} \|\nu_{\rho}^{\pi_{\theta_{t-1}}} - \nu_{\rho}^{\pi_{\theta_t}}\|_1 \\ &\leq \mathbb{E} \left\| \widehat{\mathcal{P}}_{\theta_{t-1}} \hat{\nu}_{t-1} - \widehat{\mathcal{P}}_{\theta_{t-1}} \nu_{\rho}^{\pi_{\theta_{t-1}}} \right\|_1 + L_{\nu} \mathbb{E} \|\theta_{t-1} - \theta_t\|_2 \\ &\leq \gamma \mathbb{E} \|\hat{\nu}_{t-1} - \nu_{\rho}^{\pi_{\theta_{t-1}}}\|_1 + LC_{\delta} L_{\nu} \eta_{t-1}^{\theta} \\ &\leq LC_{\delta} L_{\nu} \sum_{k=0}^{t-1} \gamma^{t-1-k} \eta_k^{\theta} + \gamma^t \|\rho - \nu_{\rho}^{\pi_{\theta_0}}\|_1. \end{aligned}$$

Proof of Lemma B.1

$$\begin{aligned} \left| \widetilde{V}_{\varphi}^{\pi_{\theta}}(s) \right| &= \left| \frac{1}{1-\gamma} \mathbb{E}_{s \sim \nu_{\rho}^{\pi_{\theta}}(\cdot)} \left[\mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [\tilde{r}_{\varphi, \theta}(s, a)] \right] \right| \\ &\leq \frac{1}{1-\gamma} \mathbb{E}_{s \sim \nu_{\rho}^{\pi_{\theta}}(\cdot)} \left[\left| \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [\tilde{r}_{\varphi, \theta}(s, a)] \right| \right] \\ &\leq \frac{C}{1-\gamma} \end{aligned}$$

Hence, $C_J = O((1-\gamma)^{-1})$. Note that by letting $\rho(s) = \mathbb{I}[s = s_0]$, we have $\left| \widetilde{V}_{\varphi}^{\pi_{\theta}}(s_0) \right| \leq C_J$ for any $s_0 \in \mathcal{S}$.

$$\begin{aligned} \|\nabla_{\theta} J_{\varphi}(\theta)\|_2 &= \left\| \frac{1}{1-\gamma} \mathbb{E}_{s \sim \nu_{\rho}^{\pi_{\theta}}(\cdot)} \left[\mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} \left[\widetilde{Q}_{\varphi}^{\pi_{\theta}}(s, a) \nabla_{\theta} \log \pi_{\theta}(a|s) \right] \right] \right\|_2 \\ &\leq \frac{L}{1-\gamma} \mathbb{E}_{s \sim \nu_{\rho}^{\pi_{\theta}}(\cdot), a \sim \pi_{\theta}(\cdot|s)} \left| \tilde{r}(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s, a)} \left[\widetilde{V}_{\varphi}^{\pi_{\theta}}(s') \right] \right| \\ &\leq \frac{CL}{(1-\gamma)^2} \end{aligned}$$

Hence, $L_J = O((1-\gamma)^{-2})$. Similarly, by letting $\rho(s) = \mathbb{I}[s = s_0]$, we have $\left| \nabla_{\theta} \widetilde{V}_{\varphi}^{\pi_{\theta}}(s_0) \right| \leq L_J$ for any $s_0 \in \mathcal{S}$.

$$\begin{aligned} \nabla_{\theta}^2 J_{\varphi}(\theta) &= \frac{1}{1-\gamma} \mathbb{E}_{\nu_{\rho}^{\pi_{\theta}}, \pi_{\theta}} \left[\widetilde{Q}^{\pi_{\theta}}(s, a) (\nabla_{\theta} \log \pi_{\theta}(a|s) \nabla_{\theta} \log \pi_{\theta}(a|s)^{\top} + \nabla_{\theta}^2 \log \pi_{\theta}(a|s)) \right] \\ &\quad + \frac{\gamma}{1-\gamma} \mathbb{E}_{\nu_{\rho}^{\pi_{\theta}}, \pi_{\theta}, \mathcal{P}} \left[\nabla_{\theta} \log \pi_{\theta}(a|s) \nabla_{\theta} \widetilde{V}_{\varphi}^{\pi_{\theta}}(s')^{\top} + \nabla_{\theta} \widetilde{V}_{\varphi}^{\pi_{\theta}}(s') \nabla_{\theta} \log \pi_{\theta}(a|s)^{\top} \right] \\ \|\nabla_{\theta}^2 J_{\varphi}(\theta)\|_2 &\leq \frac{C_J(L^2 + S)}{1-\gamma} + \frac{2\gamma LL_J}{1-\gamma} = O((1-\gamma)^{-3}). \end{aligned}$$

Hence, $S_J = O((1-\gamma)^{-3})$.

$$\begin{aligned} \|\nabla_{\varphi} J_{\varphi}(\theta)\|_2 &= \left\| \nabla_{\varphi} \left(\frac{1}{1-\gamma} \mathbb{E}_{\nu_{\rho}^{\pi_{\theta}}, \pi_{\theta}} [\tilde{r}_{\varphi, \theta}(s, a)] \right) \right\|_2 \\ &= \left\| \frac{1}{1-\gamma} \mathbb{E}_{s \sim \nu_{\rho}^{\pi_{\theta}}(\cdot)} \left[\nabla_{\varphi} \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [\tilde{r}_{\varphi, \theta}(s, a)] \right] \right\|_2 \\ &\leq \frac{D}{1-\gamma} \end{aligned}$$

Hence, $D_J = O((1-\gamma)^{-1})$.

1566 **Proof of Corollary B.2** For any state $\theta_1, \theta_2 \in \Omega(\theta)$, consider the MDP $(\mathcal{S}, \mathcal{A}, \mathcal{P}, r', \gamma)$ where for
 1567 any $s \in \mathcal{S}, a \in \mathcal{A}$,

$$1568 \quad r'(s, a) = f(s) := \begin{cases} 1 & , \nu_\rho^{\pi_{\theta_1}}(s) > \nu_\rho^{\pi_{\theta_2}}(s) \\ -1 & , \text{otherwise} \end{cases}.$$

1570 Assume the regularization factor $\alpha = 0$, then for this RL problem, $\tilde{r}(s, a) = r'(s, a)$, and

$$\begin{aligned} 1572 \quad \|\nu_\rho^{\pi_{\theta_1}} - \nu_\rho^{\pi_{\theta_2}}\|_1 &= \int_{\mathcal{S}} ds |\nu_\rho^{\pi_{\theta_1}}(s) - \nu_\rho^{\pi_{\theta_2}}(s)| \\ 1573 &= \int_{\mathcal{S}} ds (\nu_\rho^{\pi_{\theta_1}}(s) - \nu_\rho^{\pi_{\theta_2}}(s)) f(s) \\ 1574 &= \int_{\mathcal{S}} ds \nu_\rho^{\pi_{\theta_1}}(s) f(s) - \int_{\mathcal{S}} ds \nu_\rho^{\pi_{\theta_2}}(s) f(s) \\ 1575 &= \mathbb{E}_{s \sim \nu_\rho^{\pi_{\theta_1}}(\cdot)} \left[\mathbb{E}_{a \sim \pi_{\theta_1}(\cdot|s)} [\tilde{r}(s, a)] \right] - \mathbb{E}_{s \sim \nu_\rho^{\pi_{\theta_2}}(\cdot)} \left[\mathbb{E}_{a \sim \pi_{\theta_2}(\cdot|s)} [\tilde{r}(s, a)] \right] \\ 1576 &= (1 - \gamma)(J(\theta_1) - J(\theta_2)). \end{aligned}$$

1577 Then we can apply Lemma B.1 with $C = 1$ to obtain $L_\nu = (1 - \gamma)L_J$ and $S_\nu = (1 - \gamma)S_J$.

1583 **Proof of Lemma B.3**

$$\begin{aligned} 1584 \quad \|\nabla_{\theta} \mathbf{A}_{\theta}\|_2 &= \left\| \nabla_{\theta} \mathbb{E}_{\nu_\rho^{\pi_{\theta}}, \pi_{\theta}, \mathcal{P}} [\phi(s)(\phi(s) - \gamma\phi(s'))^{\top}] \right\|_2 \\ 1585 &= \left\| \int_{\mathcal{S}} ds \int_{\mathcal{A}} da \nabla_{\theta} (\nu_\rho^{\pi_{\theta}}(s) \pi_{\theta}(a|s)) \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s,a)} [\phi(s)(\phi(s) - \gamma\phi(s'))^{\top}] \right\|_2 \\ 1586 &\leq \left(\int_{\mathcal{S}} ds \int_{\mathcal{A}} da \|\nabla_{\theta} (\nu_\rho^{\pi_{\theta}}(s) \pi_{\theta}(a|s))\|_2 \right) \left(\max_s \|\mathbb{E}_{s' \sim \mathcal{P}(\cdot|s,a)} [\phi(s)(\phi(s) - \gamma\phi(s'))^{\top}]\|_2 \right) \\ 1587 &\leq (1 + \gamma) \int_{\mathcal{S}} ds \int_{\mathcal{A}} da \|\nabla_{\theta} (\nu_\rho^{\pi_{\theta}}(s) \pi_{\theta}(a|s))\|_2 \\ 1588 &= (1 + \gamma) \int_{\mathcal{S}} ds \int_{\mathcal{A}} da \|\nabla_{\theta} \nu_\rho^{\pi_{\theta}}(s) \pi_{\theta}(a|s) + \nu_\rho^{\pi_{\theta}}(s) \nabla_{\theta} \pi_{\theta}(a|s)\|_2 \\ 1589 &\leq (1 + \gamma) \left(\int_{\mathcal{S}} ds \|\nabla_{\theta} \nu_\rho^{\pi_{\theta}}(s)\|_2 \int_{\mathcal{A}} da \pi_{\theta}(a|s) + \int_{\mathcal{S}} ds \nu_\rho^{\pi_{\theta}}(s) \int_{\mathcal{A}} da \pi_{\theta}(a|s) \|\nabla_{\theta} \log \pi_{\theta}(a|s)\|_2 \right) \\ 1590 &\leq (1 + \gamma)(L_\mu + L) \\ 1591 &= O((1 - \gamma)^{-1}) \end{aligned}$$

1592 Hence, $L_A = O((1 - \gamma)^{-1})$.

$$\begin{aligned} 1600 \quad \|\nabla_{\theta}^2 \mathbf{A}_{\theta}\|_2 &= \left\| \nabla_{\theta}^2 \mathbb{E}_{\nu_\rho^{\pi_{\theta}}, \pi_{\theta}, \mathcal{P}} [\phi(s)(\phi(s) - \gamma\phi(s'))^{\top}] \right\|_2 \\ 1601 &= \left\| \int_{\mathcal{S}} ds \int_{\mathcal{A}} da \nabla_{\theta}^2 (\nu_\rho^{\pi_{\theta}}(s) \pi_{\theta}(a|s)) \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s,a)} [\phi(s)(\phi(s) - \gamma\phi(s'))^{\top}] \right\|_2 \\ 1602 &\leq \left(\int_{\mathcal{S}} ds \int_{\mathcal{A}} da \|\nabla_{\theta}^2 (\nu_\rho^{\pi_{\theta}}(s) \pi_{\theta}(a|s))\|_2 \right) \left(\max_s \|\mathbb{E}_{s' \sim \mathcal{P}(\cdot|s,a)} [\phi(s)(\phi(s) - \gamma\phi(s'))^{\top}]\|_2 \right) \\ 1603 &\leq (1 + \gamma) \int_{\mathcal{S}} ds \int_{\mathcal{A}} da \|\nabla_{\theta}^2 (\nu_\rho^{\pi_{\theta}}(s) \pi_{\theta}(a|s))\|_2 \\ 1604 &\leq (1 + \gamma) \left[\int_{\mathcal{S}} ds \|\nabla_{\theta}^2 \nu_\rho^{\pi_{\theta}}(s)\|_2 \int_{\mathcal{A}} da \pi_{\theta}(a|s) \right. \\ 1605 &\quad \left. + 2 \int_{\mathcal{S}} ds \|\nabla_{\theta} \nu_\rho^{\pi_{\theta}}(s)\|_2 \int_{\mathcal{A}} da \pi_{\theta}(a|s) \|\nabla_{\theta} \log \pi_{\theta}(a|s)\|_2 \right. \\ 1606 &\quad \left. + \int_{\mathcal{S}} ds \nu_\rho^{\pi_{\theta}}(s) \int_{\mathcal{A}} da \pi_{\theta}(a|s) \|\nabla_{\theta} \log \pi_{\theta}(a|s) \nabla_{\theta} \log \pi_{\theta}(a|s)^{\top} + \nabla_{\theta}^2 \log \pi_{\theta}(a|s)\|_2 \right] \\ 1607 &\leq (1 + \gamma) (S_\nu + 2LL_\nu + L^2 + S) \\ 1608 &= O((1 - \gamma)^{-2}) \end{aligned}$$

1618 Hence, $S_A = O((1 - \gamma)^{-2})$.

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

Proof of Lemma B.4

$$\begin{aligned}\|\mathbf{b}_{\varphi,\theta}\|_2 &= \left\| \mathbb{E}_{\nu_\rho^{\pi_\theta}, \pi_\theta} [\tilde{r}_{\varphi,\theta}(s, a)\phi(s)] \right\|_2 \\ &\leq \left| \mathbb{E}_{\nu_\rho^{\pi_\theta}, \pi_\theta} [\tilde{r}_{\varphi,\theta}(s, a)] \right| \\ &\leq C\end{aligned}$$

Hence, $C_b = O(1)$.

$$\begin{aligned}\|\nabla_{\theta} \mathbf{b}_{\varphi,\theta}\|_2 &= \left\| \nabla_{\theta} \mathbb{E}_{\nu_\rho^{\pi_\theta}, \pi_\theta} [\tilde{r}_{\varphi,\theta}(s, a)\phi(s)] \right\|_2 \\ &= \left\| \int_{\mathcal{S}} ds \nabla_{\theta} (\nu_\rho^{\pi_\theta}(s) \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [\tilde{r}_{\varphi,\theta}(s, a)]) \phi(s) \right\|_2 \\ &\leq \left\| \int_{\mathcal{S}} ds \nabla_{\theta} (\nu_\rho^{\pi_\theta}(s) \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [\tilde{r}_{\varphi,\theta}(s, a)]) \right\|_2 \left(\max_s \|\phi(s)\|_2 \right) \\ &\leq \int_{\mathcal{S}} ds \|\nabla_{\theta} (\nu_\rho^{\pi_\theta}(s) \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [\tilde{r}_{\varphi,\theta}(s, a)])\|_2 \\ &\leq \int_{\mathcal{S}} ds \|\nabla_{\theta} \nu_\rho^{\pi_\theta}(s)\|_2 \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [\tilde{r}_{\varphi,\theta}(s, a)] \\ &\quad + \int_{\mathcal{S}} ds \nu_\rho^{\pi_\theta}(s) \|\nabla_{\theta} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [\tilde{r}_{\varphi,\theta}(s, a)]\|_2 \\ &\leq CL_\nu + CL + 0 \\ &= O((1 - \gamma)^{-1})\end{aligned}$$

Hence, $L_b = O((1 - \gamma)^{-1})$.

$$\begin{aligned}\|\nabla_{\theta}^2 \mathbf{b}_{\varphi,\theta}\|_2 &= \left\| \nabla_{\theta}^2 \mathbb{E}_{\nu_\rho^{\pi_\theta}, \pi_\theta} [\tilde{r}_{\varphi,\theta}(s, a)\phi(s)] \right\|_2 \\ &= \left\| \int_{\mathcal{S}} ds \nabla_{\theta}^2 (\nu_\rho^{\pi_\theta}(s) \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [\tilde{r}_{\varphi,\theta}(s, a)]) \phi(s) \right\|_2 \\ &\leq \left\| \int_{\mathcal{S}} ds \nabla_{\theta}^2 (\nu_\rho^{\pi_\theta}(s) \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [\tilde{r}_{\varphi,\theta}(s, a)]) \right\|_2 \left(\max_s \|\phi(s)\|_2 \right) \\ &\leq \int_{\mathcal{S}} ds \|\nabla_{\theta}^2 (\nu_\rho^{\pi_\theta}(s) \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [\tilde{r}_{\varphi,\theta}(s, a)])\|_2 \\ &\leq \int_{\mathcal{S}} ds \|\nabla_{\theta}^2 \nu_\rho^{\pi_\theta}(s)\|_2 \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [\tilde{r}_{\varphi,\theta}(s, a)] \\ &\quad + \int_{\mathcal{S}} ds \|\nabla_{\theta} \nu_\rho^{\pi_\theta}(s)\|_2 \|\nabla_{\theta} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [\tilde{r}_{\varphi,\theta}(s, a)]\|_2 \\ &\quad + \int_{\mathcal{S}} ds \nu_\rho^{\pi_\theta}(s) \|\nabla_{\theta}^2 \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [\tilde{r}_{\varphi,\theta}(s, a)]\|_2 \\ &\leq CS_\nu + CLL_\nu + C(L^2 + S) \\ &= O((1 - \gamma)^{-2})\end{aligned}$$

Hence, $S_b = O((1 - \gamma)^{-2})$.

1674
 1675
 1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684
 1685
 1686
 1687
 1688
 1689
 1690
 1691
 1692
 1693
 1694
 1695
 1696
 1697
 1698
 1699
 1700
 1701
 1702
 1703
 1704
 1705
 1706
 1707
 1708
 1709
 1710
 1711
 1712
 1713
 1714
 1715
 1716
 1717
 1718
 1719
 1720
 1721
 1722
 1723
 1724
 1725
 1726
 1727

$$\begin{aligned}
 \|\nabla_{\varphi} \mathbf{b}_{\varphi, \theta}\|_2 &= \left\| \nabla_{\varphi} \mathbb{E}_{\nu_{\rho}^{\pi_{\theta}}, \pi_{\theta}} [\tilde{r}_{\varphi, \theta}(s, a) \phi(s)] \right\|_2 \\
 &= \left\| \int_{\mathcal{S}} ds \nabla_{\varphi} (\nu_{\rho}^{\pi_{\theta}}(s) \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [\tilde{r}_{\varphi, \theta}(s, a)]) \phi(s) \right\|_2 \\
 &\leq \left\| \int_{\mathcal{S}} ds \nabla_{\varphi} (\nu_{\rho}^{\pi_{\theta}}(s) \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [\tilde{r}_{\varphi, \theta}(s, a)]) \right\|_2 \left(\max_s \|\phi(s)\|_2 \right) \\
 &\leq \int_{\mathcal{S}} ds \|\nabla_{\varphi} (\nu_{\rho}^{\pi_{\theta}}(s) \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [\tilde{r}_{\varphi, \theta}(s, a)])\|_2 \\
 &= \int_{\mathcal{S}} ds \nu_{\rho}^{\pi_{\theta}}(s) \|\nabla_{\varphi} \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [\tilde{r}_{\varphi, \theta}(s, a)]\|_2 \\
 &\leq D
 \end{aligned}$$

Hence, $D_b = O(1)$.

Proof of Lemma B.5

$$\|\omega^*(\varphi, \theta)\|_2 = \|\mathbf{A}_{\theta}^{-1} \mathbf{b}_{\varphi, \theta}\|_2 \leq \|\mathbf{A}_{\theta}^{-1}\|_2 \|\mathbf{b}_{\varphi, \theta}\|_2 \leq \frac{C_b}{\lambda} = \frac{C}{\lambda}$$

Hence, $C_{\omega} = O(\lambda^{-1})$.

$$\begin{aligned}
 \|\nabla_{\theta} \omega^*(\varphi, \theta)\|_2 &= \|\nabla_{\theta} (\mathbf{A}_{\theta}^{-1} \mathbf{b}_{\varphi, \theta})\|_2 \\
 &= \|\nabla_{\theta} (\mathbf{A}_{\theta}^{-1}) \mathbf{b}_{\varphi, \theta} + \mathbf{A}_{\theta}^{-1} \nabla_{\theta} \mathbf{b}_{\varphi, \theta}\|_2 \\
 &= \|\mathbf{A}_{\theta}^{-1} \nabla_{\theta} \mathbf{A}_{\theta} \mathbf{A}_{\theta}^{-1} \mathbf{b}_{\varphi, \theta} + \mathbf{A}_{\theta}^{-1} \nabla_{\theta} \mathbf{b}_{\varphi, \theta}\|_2 \\
 &\leq \|\mathbf{A}_{\theta}^{-1}\|_2 \|\nabla_{\theta} \mathbf{A}_{\theta}\|_2 \|\mathbf{A}_{\theta}^{-1}\|_2 \|\mathbf{b}_{\varphi, \theta}\|_2 + \|\mathbf{A}_{\theta}^{-1}\|_2 \|\nabla_{\theta} \mathbf{b}_{\varphi, \theta}\|_2 \\
 &\leq L_A C_b \lambda^{-2} + L_b \lambda^{-1} \\
 &= O((1 - \gamma)^{-1} \lambda^{-2})
 \end{aligned}$$

Hence, $L_{\omega} = O((1 - \gamma)^{-1} \lambda^{-2})$.

$$\begin{aligned}
 \|\nabla_{\theta}^2 \omega^*(\varphi, \theta)\|_2 &= \|\nabla_{\theta}^2 (\mathbf{A}_{\theta}^{-1} \mathbf{b}_{\varphi, \theta})\|_2 \\
 &= \|\nabla_{\theta} (\mathbf{A}_{\theta}^{-1} \nabla_{\theta} \mathbf{A}_{\theta} \mathbf{A}_{\theta}^{-1} \mathbf{b}_{\varphi, \theta} + \mathbf{A}_{\theta}^{-1} \nabla_{\theta} \mathbf{b}_{\varphi, \theta})\|_2 \\
 &= \|\mathbf{A}_{\theta}^{-1} \nabla_{\theta}^2 \mathbf{A}_{\theta} \mathbf{A}_{\theta}^{-1} \mathbf{b}_{\varphi, \theta} + 2 \mathbf{A}_{\theta}^{-1} \nabla_{\theta} \mathbf{A}_{\theta} \mathbf{A}_{\theta}^{-1} \nabla_{\theta} \mathbf{A}_{\theta} \mathbf{A}_{\theta}^{-1} \mathbf{b}_{\varphi, \theta} \\
 &\quad + 2 \mathbf{A}_{\theta}^{-1} \nabla_{\theta} \mathbf{A}_{\theta} \mathbf{A}_{\theta}^{-1} \nabla_{\theta} \mathbf{b}_{\varphi, \theta} + \mathbf{A}_{\theta}^{-1} \nabla_{\theta}^2 \mathbf{b}_{\varphi, \theta}\|_2 \\
 &\leq \|\mathbf{A}_{\theta}^{-1}\|_2 \|\nabla_{\theta}^2 \mathbf{A}_{\theta}\|_2 \|\mathbf{A}_{\theta}^{-1}\|_2 \|\mathbf{b}_{\varphi, \theta}\|_2 \\
 &\quad + 2 \|\mathbf{A}_{\theta}^{-1}\|_2 \|\nabla_{\theta} \mathbf{A}_{\theta}\|_2 \|\mathbf{A}_{\theta}^{-1}\|_2 \|\nabla_{\theta} \mathbf{A}_{\theta}\|_2 \|\mathbf{A}_{\theta}^{-1}\|_2 \|\mathbf{b}_{\varphi, \theta}\|_2 \\
 &\quad + 2 \|\mathbf{A}_{\theta}^{-1}\|_2 \|\nabla_{\theta} \mathbf{A}_{\theta}\|_2 \|\mathbf{A}_{\theta}^{-1}\|_2 \|\nabla_{\theta} \mathbf{b}_{\varphi, \theta}\|_2 \\
 &\quad + \|\mathbf{A}_{\theta}^{-1}\|_2 \|\nabla_{\theta}^2 \mathbf{b}_{\varphi, \theta}\|_2 \\
 &\leq S_A C_b \lambda^{-2} + 2 L_A^2 C_b \lambda^{-3} + 2 L_A L_b \lambda^{-2} + S_b \lambda^{-1} \\
 &= O((1 - \gamma)^{-2} \lambda^{-3})
 \end{aligned}$$

Hence, $S_{\omega} = O((1 - \gamma)^{-2} \lambda^{-3})$.

$$\begin{aligned}
 \|\nabla_{\varphi} \omega^*(\varphi, \theta)\|_2 &= \|\nabla_{\varphi} (\mathbf{A}_{\theta}^{-1} \mathbf{b}_{\varphi, \theta})\|_2 \\
 &= \|\mathbf{A}_{\theta}^{-1} \nabla_{\varphi} \mathbf{b}_{\varphi, \theta}\|_2 \\
 &\leq \|\mathbf{A}_{\theta}^{-1}\|_2 \|\nabla_{\varphi} \mathbf{b}_{\varphi, \theta}\|_2 \\
 &\leq \frac{D_b}{\lambda} \\
 &= O(\lambda^{-1})
 \end{aligned}$$

Hence, $D_{\omega} = O(\lambda^{-1})$.

1728
 1729
 1730
 1731
 1732
 1733
 1734
 1735
 1736
 1737
 1738
 1739
 1740
 1741
 1742
 1743
 1744
 1745
 1746
 1747
 1748
 1749
 1750
 1751
 1752
 1753
 1754
 1755
 1756
 1757
 1758
 1759
 1760
 1761
 1762
 1763
 1764
 1765
 1766
 1767
 1768
 1769
 1770
 1771
 1772
 1773
 1774
 1775
 1776
 1777
 1778
 1779
 1780
 1781

Proof of Lemma B.6

$$\begin{aligned}
 \mathbb{E}_{\nu, \pi_{\theta}, \mathcal{P}} \|\hat{\delta}(s, a, s')\|_2^2 &= \int_{\mathcal{S}} ds \nu(s) \mathbb{E}_{\pi_{\theta}, \mathcal{P}} \left[(\tilde{r}(s, a) + (\phi(s') - \phi(s))^{\top} \boldsymbol{\omega})^2 \right] \\
 &\leq \int_{\mathcal{S}} ds \nu(s) \mathbb{E}_{\pi_{\theta}} \left[(\tilde{r}(s, a) + 2C_{\boldsymbol{\omega}})^2 \right] \\
 &\leq \int_{\mathcal{S}} ds \nu(s) \left(\mathbb{E}_{\pi_{\theta}} [\tilde{r}(s, a)^2] + 4C_{\boldsymbol{\omega}} \mathbb{E}_{\pi_{\theta}} [\tilde{r}(s, a)] + 4C_{\boldsymbol{\omega}}^2 \right) \\
 &\leq (C^2 + 4CC_{\boldsymbol{\omega}} + 4C_{\boldsymbol{\omega}}^2) \\
 &= (C + 2C_{\boldsymbol{\omega}})^2
 \end{aligned}$$

Hence, $C_{\delta} = C + 2C_{\boldsymbol{\omega}} = O(\lambda^{-1})$.

1782 E LLM USAGE
1783

1784 Large Language Models (LLMs) were used to aid in the writing and polishing of the manuscript.
1785 Specifically, we used an LLM to assist in refining the language, improving readability, and ensuring
1786 clarity in various sections of the paper. The model helped with tasks such as sentence rephrasing,
1787 grammar checking, and enhancing the overall flow of the text.

1788 It is important to note that the LLM was not involved in the ideation, research methodology, or
1789 experimental design. All research concepts, ideas, and analyses were developed and conducted by
1790 the authors. The contributions of the LLM were solely focused on improving the linguistic quality
1791 of the paper, with no involvement in the scientific content or data analysis.

1792 The authors take full responsibility for the content of the manuscript, including any text generated
1793 or polished by the LLM. We have ensured that the LLM-generated text adheres to ethical guidelines
1794 and does not contribute to plagiarism or scientific misconduct.
1795

1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835