# On the Convergence of Continuous Single-timescale Actor-critic

Xuyang Chen<sup>1</sup> Lin Zhao<sup>1</sup>

# Abstract

Actor-critic algorithms have been instrumental in boosting the performance of numerous challenging applications involving continuous control, such as highly robust and agile robot motion control. However, their theoretical understanding remains largely underdeveloped. Existing analyses mostly focus on finite state-action spaces and on simplified variants of actor-critic, such as doubleloop updates with i.i.d. sampling, which are often impractical for real-world applications. We consider the canonical and widely adopted singletimescale updates with Markovian sampling in continuous state-action space. Specifically, we establish finite-time convergence by introducing a novel Lyapunov analysis framework, which provides a unified convergence characterization of both the actor and the critic. Our approach is less conservative than previous methods and offers new insights into the coupled dynamics of actor-critic updates.

# 1. Introduction

Actor-critic methods have achieved substantial success in many challenging applications (Mnih et al., 2016; Silver et al., 2017; Vinyals et al., 2019; Lazaridis et al., 2020). In particular, it becomes instrumental in enabling highly robust and agile robot motion control involving continuous state-action spaces, such as quadruped locomotion control (Hoeller et al., 2024), humanoid whole-body control (Radosavovic et al., 2024), drone racing (Kaufmann et al., 2023), etc.

Despite substantial empirical success, the theoretical analysis of actor-critic is significantly behind. Most prior theoretical studies of actor-critic methods consider only finite state-action spaces and focus on their impractical variants to simplify the analysis, including the double-loop updates and the two-timescale updates. The double-loop updates perform multiple critic updates for a fixed actor (Yang et al., 2019; Kumar et al., 2023; Agarwal et al., 2021; Xu et al., 2020b). This facilitates more accurate value function estimation, which in turn enables a more precise policy gradient estimation for the fixed actor. It allows a simple decoupled analysis of the actor and the critic. However, such an implementation is impractical due to the high sampling complexity. Another variant is the two-timescale actor-critic method (Wu et al., 2020; Xu et al., 2020c; Chen et al., 2023; Shen et al., 2023; Hong et al., 2023), which assigns a smaller step size for the actor than that of the critic, with their ratio converging to zero as the number of iterations approaches infinity (i.e.,  $\lim_{t\to\infty} \alpha_t / \beta_t = 0$ ). It allows an asymptotically decoupling of the actor and the critic in the convergence analysis, similar to performing multiple critic updates at a fixed actor. However, such artificial slowing down of the critic update is often not desired in practice.

The canonical and more practical implementation of actorcritic is the single-timescale update, where the actor and the critic are updated simultaneously with proportional step sizes at each iteration (i.e.,  $\alpha_t/\beta_t = c$ ). However, analyzing its convergence is significantly more challenging than for the aforementioned simplified variants, as the actor and critic updates are strongly coupled. The aforementioned decoupled analysis is over-conservative and cannot establish the convergence of the single-timescale actor-critic. Recent efforts to study the convergence of the single-timescale actor-critic algorithm include Chen et al. (2021), Olshevsky & Gharesifard (2023), and Chen & Zhao (2024). However, these works are limited to finite action spaces with i.i.d. sampling and do not extend to the more practical yet complex setting of Markovian sampling in continuous state-action spaces under the single-timescale update scheme (See the comparison in Table 1). In particular, Chen et al. (2021) and Olshevsky & Gharesifard (2023) assume i.i.d. sampling directly from the *stationary distribution* for the critic and from the discounted state visitation distribution for the actor. However, both of these distributions are unknown for real-time online learning and hence are impractical. Additionally, Chen & Zhao (2024) considers the simpler undiscounted time-average reward setting rather than the widely adopted discounted reward setting. The key difference is

<sup>&</sup>lt;sup>1</sup>Department of Electrical and Computer Engineering, National University of Singapore, Singapore. Correspondence to: Lin Zhao <elezhli@nus.edu.sg>.

Proceedings of the  $42^{nd}$  International Conference on Machine Learning, Vancouver, Canada. PMLR 267, 2025. Copyright 2025 by the author(s).

References	State Space	Action Space	Sampling for Critic	Sampling for Actor	Complexity
Chen et al. (2021)	Infinite	Finite	i.i.d. from	i.i.d. from state	$\mathcal{O}(\epsilon^{-2})$
			stationary distribution	visitation distribution	
Olshevsky & Gharesifard (2023)	Finite	Finite	i.i.d. from	i.i.d. from state	${\cal O}(\epsilon^{-2})$
			stationary distribution	visitation distribution	
This paper	Infinite	Infinite	Markovian	Markovian	$ ilde{\mathcal{O}}(\epsilon^{-2})$

*Table 1.* Comparison of existing works on single-timescale actor-critic methods in discounted reward setting with linear function approximation. Our work is the first to address the continuous state-action spaces and Markovian sampling.

that, in the former, the policy gradient only requires stationary distribution, whereas in the latter, it depends on the visitation distribution. How to accurately approximate the visitation distribution in the single-timescale update scheme with Markovian sampling remains an open question. To tackle the aforementioned challenges, specifically,

- 1. We introduce a new operator-based analysis to handle the intricacies arising from the uncountable continuous space. In particular, it enables us to generalize many important bounds to the continuous space successfully (see Appendix B).
- 2. For the Markovian samples used to update the actor, we prove that the resulted state distribution converges to the *discounted state visitation distribution* (see Proposition 3.2). We further utilize it to accurately estimate the policy gradient in the analysis.
- 3. We propose a Lyapunov-based convergence analysis framework, where a novel Lyapunov function is constructed specifically for the single-timescale actorcritic algorithm. We also establish a variety of new properties (see, for example, Proposition 3.1, Proposition 4.4), which enable us to demonstrate finite-time convergence for both the actor and the critic simultaneously, with a less conservative analysis.

Moreover, we highlight that our analysis builds on the same set of common assumptions that are widely adopted in many literature (Wu et al., 2020; Chen et al., 2021; Olshevsky & Gharesifard, 2023; Chen & Zhao, 2024). In particular, Assumptions 4.1 and 4.2 are about the regularity of the problem of interest, and Assumption 4.3 can be easily satisfied by many common policy parametrizations. Overall, our work takes a significant step toward a more practical analysis of actor-critic.

#### 1.1. Related Work

In this section, we review the existing works on actor-critic methods.

Actor-Critic methods. The actor-critic algorithm, initially proposed by (Konda & Tsitsiklis, 1999), was later extended to the natural actor-critic variant by (Kakade, 2001). The asymptotic convergence of actor-critic algorithms has been well established under various settings, as demonstrated in works by Kakade (2001), Bhatnagar et al. (2009), and Zhang et al. (2020b). More recently, many studies have focused on the finite-time convergence of actor-critic methods. They primarily focus on two variants for the ease of analysis: (1) double loop update, and (2) two-timescale update. For the double-loop variants, Kumar et al. (2023) investigated the finite-time local convergence of several actor-critic variants with linear function approximation. Wang et al. (2019) explored the global convergence of actor-critic methods with both the actor and the critic parameterized by neural networks with single hidden layers. Moreover, Gaur et al. (2024) established the last iterate convergence for actor-critic with neural networks.

For the two-timescale variants, Wu et al. (2020) established finite-time local convergence in the undiscounted time-average reward setting. Xu et al. (2020c) analyzed both local and global convergence for two-timescale (natural) actor-critic under the discounted reward setting, respectively, with multiple samples used for critic updates. Shen et al. (2023) investigated finite-time convergence for asynchronous actor-critic, while Hong et al. (2023) introduced a two-timescale stochastic approximation algorithm for bilevel optimization and two-timescale actor-critic.

There are only a few works considering the canonical and most widely adopted single-timescale variant. Fu et al. (2020) explored the least-squares temporal difference (LSTD) update for the critic, achieving the optimal policy within the energy-based policy class for both linear function approximation and neural network approximation. Zhou & Lu (2023) and Chen et al. (2024) established the global convergence of actor-critic methods for solving linear quadratic regulator. Recently, Chen et al. (2021); Olshevsky & Gharesifard (2023); Chen & Zhao (2024) investigated single-timescale actor-critic methods in general Markov Decision Processes (MDPs) with linear function approximation, aligning with the focus of this work. Specifically, Chen et al. (2021) and Olshevsky & Gharesifard (2023) addressed the commonly used discounted reward setting, while Chen & Zhao (2024) improved upon (Wu et al., 2020) by advancing from the two-timescale to the single-timescale approach under the undiscounted time-average reward setting. A detailed review and comparison of these results can be found in Table 1 and the introduction.

**Notation.** We use san-serif letters to denote scalars and use lower and upper case bold letters to denote vectors and matrices respectively. We also use  $\|\boldsymbol{\omega}\|$  to denote the  $\ell_2$ -norm of a vector  $\boldsymbol{\omega}$  and  $\|\boldsymbol{A}\|$  to denote the spectral norm of a matrix  $\boldsymbol{A}$ . Without further specification, we write  $x_n = \mathcal{O}(y_n)$ if there exists an absolute positive constant C such that  $x_n \leq Cy_n$ , for two sequences  $\{x_n\}$  and  $\{y_n\}$ . We use  $\tilde{\mathcal{O}}(\cdot)$  to hide logarithm factors. The total variation distance of two probability measure  $\mu$  and  $\nu$  is defined by  $d_{TV}(\mu,\nu) := 1/2 \int_{\mathcal{X}} |\mu(dx) - \nu(dx)|$ .

## 2. Preliminaries

**Markov Decision Process.** In this paper, we consider a discrete-time Markov Decision Process (MDP) defined by a tuple  $\mathcal{M} = \{S, \mathcal{A}, P, r, \gamma\}$ , where S is the state space and  $\mathcal{A}$  is the action space. The spaces S and  $\mathcal{A}$  are allowed to be either finite sets or real vector spaces, i.e.,  $S \subset \mathbb{R}^{d_s}$  and  $\mathcal{A} \subset \mathbb{R}^{d_a}$ . The transition kernel is denoted by  $P(s_{t+1} | s_t, a_t) \in \mathbb{R}_{\geq 0}$ , the reward function is  $r : S \times \mathcal{A} \rightarrow [-\bar{r}, \bar{r}]$ , and  $\gamma \in (0, 1)$  is the discounted factor. We also assume that the initial state is sampled from a fixed initial distribution  $\eta$ .

A policy  $\pi_{\theta}$  parameterized by  $\theta \in \mathcal{X}_{\Theta}$  maps a given state to a probability distribution over the action space, i.e.,  $a_t \sim \pi_{\theta}(\cdot | s_t)$ . We denote the stationary distribution induced by the policy  $\pi_{\theta}$  and the transition kernel P by  $\mu_{\theta}$ . The value function of a state s under a policy  $\pi_{\theta}$  is the expected cumulative return when starting in s and following  $\pi_{\theta}$  thereafter. It is defined as

$$V_{\boldsymbol{\theta}}(s) = \mathbb{E}_{a_t \sim \pi_{\boldsymbol{\theta}}(\cdot \mid s_t)} \bigg[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \bigg], \quad (1)$$

where the expectation takes over the randomness of the policy  $\pi_{\theta}$  and the transition function *P*. The corresponding action-value function is the expected cumulative return when starting from state *s*, taking action *a*, and following  $\pi_{\theta}$  thereafter, which is defined as

$$Q_{\boldsymbol{\theta}}(s,a) = \mathbb{E}\bigg[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \, \big| \, s_0 = s, a_0 = a\bigg], \quad (2)$$

where we simplified the expectation notation when there is no confusion. The reinforcement learning (RL) tasks typically aim to find a policy  $\pi_{\theta}$  that maximizes the following objective function:

$$J(\boldsymbol{\theta}) = \int_{\mathcal{S}} \eta(s) V_{\boldsymbol{\theta}}(s) \, ds, \tag{3}$$

where  $\eta(s)$  is a fixed initial distribution.

We denote the density at state s' after transitioning for one time step from state s by  $P_{\theta}(s' \mid s)$ , which is defined as

$$P_{\boldsymbol{\theta}}(s' \mid s) = \int_{\mathcal{A}} P(s' \mid s, a) \pi_{\boldsymbol{\theta}}(a \mid s) \, da.$$

The corresponding state density after transitioning for t time steps can be acquired by recursively applying  $P_{\theta}(s' | s)$ , i.e.,

$$P_{\boldsymbol{\theta}}^{t}(s' \mid s) = \int_{\mathcal{S}} P_{\boldsymbol{\theta}}(s' \mid x) P_{\boldsymbol{\theta}}^{t-1}(x \mid s) \, dx, \, t > 1.$$

Consequently, we define the *discounted state visitation distribution* under policy  $\pi_{\theta}$  as

$$\nu_{\boldsymbol{\theta}}(s') = (1 - \gamma) \int_{\mathcal{S}} \sum_{t=0}^{\infty} \gamma^t \eta(s) P_{\boldsymbol{\theta}}^t(s' \,|\, s) \, ds. \tag{4}$$

It is worth noting that previous works (Chen et al., 2021; Olshevsky & Gharesifard, 2023) rely on sampling from this distribution, which is infeasible. In this work, we propose a practical sampling scheme to circumvent this impediment.

With the discounted state visitation distribution, the objective function can be reformulated as (Sutton et al., 1999)

$$J(\boldsymbol{\theta}) = \frac{1}{1-\gamma} \int_{\mathcal{S}} \nu_{\boldsymbol{\theta}}(s) \int_{\mathcal{A}} \pi_{\boldsymbol{\theta}}(a \mid s) r(s, a) \, dads$$
$$= \frac{1}{1-\gamma} \mathbb{E}_{s \sim \nu_{\boldsymbol{\theta}}, a \sim \pi_{\boldsymbol{\theta}}} \Big[ r(s, a) \Big].$$

**Policy Gradient Theorem.** Policy gradient algorithms are among the most widely used approaches in continuousaction reinforcement learning. Their core concept involves adjusting the policy parameter  $\theta$  in the direction of the performance gradient  $\nabla_{\theta} J(\theta)$ . These algorithms are built upon the foundational result known as the policy gradient theorem (Sutton et al., 1999):

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{1-\gamma} \int_{\mathcal{S}} \nu_{\boldsymbol{\theta}}(s) \int_{\mathcal{A}} Q_{\boldsymbol{\theta}}(s, a) \nabla_{\boldsymbol{\theta}} \pi_{\boldsymbol{\theta}}(a \mid s) dads$$
$$= \frac{1}{1-\gamma} \mathbb{E}_{s \sim \nu_{\boldsymbol{\theta}}, a \sim \pi_{\boldsymbol{\theta}}} \Big[ Q_{\boldsymbol{\theta}}(s, a) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a \mid s) \Big].$$
(5)

Computing this gradient necessitates the Q-value associated with the current policy  $\pi_{\theta}$ . The REINFORCE algorithm (Williams, 1992), an episodic Monte Carlo-based method, approximates the true Q-value by utilizing the cumulative rewards gathered along the sampled trajectory.

Note that for any function  $b : S \to \mathbb{R}$  that is independent of the action, we have

$$\int_{\mathcal{A}} b(s) \nabla \pi_{\boldsymbol{\theta}}(a|s) = b(s) \nabla \int_{\mathcal{A}} \pi_{\boldsymbol{\theta}}(a|s) = b(s) \nabla 1 = 0.$$

Therefore, the policy gradient theorem can be written equivalently as:

$$\nabla J(\boldsymbol{\theta}) = \frac{1}{1-\gamma} \mathbb{E}_{s,a} \big[ (Q_{\boldsymbol{\theta}}(s,a) - b(s)) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a|s) \big],$$

where b(s) is called the baseline function. A popular choice of baseline is the state-value function, which leads to the following advantage-based policy gradient

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim \nu_{\boldsymbol{\theta}}, a \sim \pi_{\boldsymbol{\theta}}} \big[ G_{\boldsymbol{\theta}}(s, a) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a|s) \big],$$

where  $G_{\theta}(s, a) = Q_{\theta}(s, a) - V_{\theta}(s)$  is known as the advantage function. This is the "REINFORCE with a baseline".

The baseline function can help reduce variance. However, like all Monte Carlo-based methods, it can still suffer from high variance and thus learns slowly. An alternative approach involves introducing an additional trainable model to approximate the value function, a method typically known as actor-critic methods.

## 3. Actor-Critic Methods

In this work, we analyze the classic single-sample singletimescale actor-critic method, where the critic employs bootstrapping by using a single sampled reward to update its value estimate at each iteration. We consider the following linear function approximation of the state-value function:

$$\widehat{V}_{\boldsymbol{\theta}}(s;\boldsymbol{\omega}) = \boldsymbol{\phi}(s)^{\top}\boldsymbol{\omega},$$

where  $\phi(\cdot) : S \to \mathbb{R}^d$  is a known feature mapping, which satisfies  $\|\phi(\cdot)\| \leq 1$ . To align  $\hat{V}_{\theta}(s; \omega)$  with its true value  $V_{\theta}(s)$ , the semi-gradient TD(0) update is employed to estimate the linear coefficient  $\omega$  (hereafter referred to as the critic):

$$\boldsymbol{\omega}_{t+1} = \boldsymbol{\omega}_t + \beta \left( \boldsymbol{r}_t + \gamma \boldsymbol{\phi}(\boldsymbol{s}_{t+1})^\top \boldsymbol{\omega}_t - \boldsymbol{\phi}(\boldsymbol{s}_t)^\top \boldsymbol{\omega}_t \right) \boldsymbol{\phi}(\boldsymbol{s}_t),$$

where  $\beta$  is the step size of the critic  $\omega$  and  $r_t := r(s_t, a_t)$ . Denote the transition tuple as O := (s, a, s') and we define the following temporal difference error

$$\delta(O, \boldsymbol{\omega}) = r(s, a) + \gamma \boldsymbol{\phi}(s')^{\top} \boldsymbol{\omega} - \boldsymbol{\phi}(s)^{\top} \boldsymbol{\omega}_{s}$$

and the update rule for the critic is then given by

$$\boldsymbol{\omega}_{t+1} = \boldsymbol{\omega}_t + \beta \delta(O_t, \boldsymbol{\omega}_t) \boldsymbol{\phi}(s_t), \tag{6}$$

where  $O_t = (s_t, a_t, s_{t+1})$  denotes the *t*-th transition tuple for the critic, generated via Markovian sampling under the policy  $\pi_{\theta}$  and transition kernel *P*, such that

$$O_t = \left(s_t, a_t \sim \pi_{\boldsymbol{\theta}_t}(\cdot \mid s_t), s_{t+1} \sim P(\cdot \mid s_t, a_t)\right).$$
(7)

Since  $\delta$  is an approximation of the advantage function, similar to REINFORCE with a baseline, the corresponding update rule for the actor can be written as:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha \delta(\widehat{O}_t, \boldsymbol{\omega}_t) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}_t}(\hat{a}_t \,|\, \hat{s}_t), \quad (8)$$

where  $\alpha$  is the step size of the actor and  $\hat{O}_t = (\hat{s}_t, \hat{a}_t, \hat{s}_{t+1})$ denotes the *t*-th transition tuple for the actor. Specifically,  $\hat{O}_t$  is also generated via the following Markovian sampling (Konda & Tsitsiklis, 2003; Shen et al., 2023)

$$\widehat{O}_{t} = \left(\widehat{s}_{t}, \widehat{a}_{t} \sim \pi_{\boldsymbol{\theta}_{t}}(\cdot \mid \widehat{s}_{t}), \widehat{s}_{t+1} \sim \widehat{P}(\cdot \mid \widehat{s}_{t}, \widehat{a}_{t})\right), \quad (9)$$
where  $\widehat{P} = \gamma P + (1 - \gamma)\eta.$ 

Here the transition kernel  $\hat{P}$  is defined as with probability  $\gamma$ , the next state follows the original transition kernel P; Otherwise, with probability  $1 - \gamma$ , the next state is sampled from the initial distribution  $\eta$ . Note that the above Markovian sampling generally requires a simulator whose state can be arbitrarily reset. It has a few nice properties that will be discussed shortly, which facilitate an accurate estimation of the policy gradient.

Denote the class of real-valued functions on the state space S by  $\mathcal{F} := \{f \mid f : S \to \mathbb{R}\}$ . We define the operator  $\mathcal{P} : \mathcal{F} \to \mathcal{F}$  acts on a state distribution  $f \in \mathcal{F}$  by

$$(\mathcal{P}f)(s') = \int_{\mathcal{S}} \int_{\mathcal{A}} f(s)\pi_{\theta}(a \mid s) P(s' \mid s, a) \, dads.$$
(10)

We further define a reset operator  $\mathcal{E} : \mathcal{F} \to \mathcal{F}$  such that it reset any state distribution f to the initial distribution  $\eta$ :

$$(\mathcal{E}f)(s) = \eta(s).$$

Therefore, the operator  $\widehat{\mathcal{P}}$  that acts on a state distribution, describing how the distribution evolves after a single step of the Markov chain induced by the policy  $\pi_{\theta}$  and the transition kernel  $\widehat{P}$ , can be written compactly as:

$$\widehat{\mathcal{P}} = \gamma \mathcal{P} + (1 - \gamma) \mathcal{E}$$

We show in the following proposition that the discounted state visitation distribution  $\nu_{\theta}$  defined in Eq. (4) is the stationary distribution of the Markov chain induced by policy  $\pi_{\theta}$  and transition kernel  $\hat{P}$  by showing that  $\nu_{\theta}$  is the unique fixed point of the operator  $\hat{\mathcal{P}}$ .

**Proposition 3.1.**  $\nu_{\theta}(s)$  is the unique fixed point of the operator  $\widehat{\mathcal{P}}$ , that is,

$$(\widehat{\mathcal{P}}\nu_{\theta})(s) = \nu_{\theta}(s),$$

Algorithm 1 Continuous Single-sample Single-timescale Actor-Critic with Markovian Sampling

- Initialize: actor parameter θ<sub>0</sub>, critic parameter ω<sub>0</sub>, initial states s<sub>0</sub>, ŝ<sub>0</sub> ~ η, stepsizes α for actor, β for critic.
   for t = 0, 1, 2, · · · , T − 1 do
- 3: Markovian sampling:
- 4:  $O_t = (s_t, a_t \sim \pi_{\theta_t}(\cdot | s_t), s_{t+1} \sim P(\cdot | s_t, a_t)),$
- 5:  $\widehat{O}_t = (\widehat{s}_t, \widehat{a}_t \sim \pi_{\boldsymbol{\theta}_t}(\cdot \mid \widehat{s}_t), \widehat{s}_{t+1} \sim \widehat{P}(\cdot \mid \widehat{s}_t, \widehat{a}_t)).$
- 6: Critic and actor update:
- 7:  $\omega_{t+1} = proj_{\bar{\omega}} (\boldsymbol{\omega}_t + \beta \delta(O_t, \boldsymbol{\omega}_t) \boldsymbol{\phi}(s_t)),$
- 8:  $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha \delta(\widehat{O}_t, \boldsymbol{\omega}_t) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}_t}(\widehat{a}_t \mid \widehat{s}_t).$
- 9: end for

and therefore the stationary distribution of the following Markov chain induced by  $\pi_{\theta}$  and  $\hat{P}$ ,

$$\hat{s}_0 \xrightarrow{(\pi_{\theta}, \hat{P})} \hat{s}_1 \xrightarrow{(\pi_{\theta}, \hat{P})} \cdots \xrightarrow{(\pi_{\theta}, \hat{P})} \hat{s}_t \xrightarrow{(\pi_{\theta}, \hat{P})} \hat{s}_{t+1}.$$
(11)

The above proposition justifies the actor's sampling scheme in Eq. (9), as  $\hat{O}_t$  effectively approximates the discounted state visitation distribution, which is required by the policy gradient theorem (Eq. (5)) following the actor update formula in Eq. (8).

**Proposition 3.2.** For the Markov chain defined in Eq. (11), we have

$$d_{TV}(\mathbb{P}(\hat{s}_t \in \cdot | \hat{s}_0 = s), \nu_{\theta}(\cdot)) \leq \gamma^t, \, \forall t \geq 0, \forall s \in \mathcal{S}.$$

Proposition 3.2 states that the distribution of  $\hat{s}_t$  converges to  $\nu_{\theta}$  geometrically with rate  $\gamma$ , a crucial property for managing the Markovian noise arising from the actor's Markovian sampling in Eq. (9).

We summarize the above-described actor-critic algorithm in Algorithm 1. The "continuous" refers to the general setting of continuous state-action spaces. "single-timescale" refers to the fact that the stepsizes  $\alpha$  and  $\beta$  are kept in constant proportion. In addition, the terminology "singlesample" follows Olshevsky & Gharesifard (2023), which refers to the fact that at each iteration, the critic and the actor are each updated using a single sample. Note that Olshevsky & Gharesifard (2023), who consider the discounted reward setting, assume access to samples from the discounted state visitation distribution and the stationary distribution for updating the actor and critic, respectively. This assumption requires a simulator capable of resetting to arbitrary states. In the simpler *time-average reward* setting (Wu et al., 2020; Chen & Zhao, 2024), the policy gradient depends solely on the stationary distribution, allowing the actor to utilize the same samples as the critic. In contrast, discounted reward setting requires the policy gradient to be computed with respect to the visitation distribution, as shown in Eq. (5), which is more challenging. To this end, the Markovian sampling strategy introduced in Eq. (9) becomes necessary to track this distribution. Consequently, Algorithm 1 inherently supports online learning and applies to continuous control tasks.

As shown in Line 4 and Line 5 of Algorithm 1, we adopt Markovian sampling for both the critic and the actor. Specifically, the transition tuple for the critic is generated by the following Markov chain

$$s_0 \xrightarrow{(\pi_{\theta_0}, P)} s_1 \xrightarrow{(\pi_{\theta_1}, P)} \cdots \xrightarrow{(\pi_{\theta_{t-1}}, P)} s_t \xrightarrow{(\pi_{\theta_t}, P)} s_{t+1}.$$
(12)

while the actor's transition tuple is generated by the Markov chain

$$\hat{s}_0 \xrightarrow{(\pi_{\theta_0}, \hat{P})} \hat{s}_1 \xrightarrow{(\pi_{\theta_1}, \hat{P})} \cdots \xrightarrow{(\pi_{\theta_{t-1}}, \hat{P})} \hat{s}_t \xrightarrow{(\pi_{\theta_t}, \hat{P})} \hat{s}_{t+1}.$$
(13)

Note that the above Markov chains (time-inhomogeneous) differ from the one defined in Eq. (11) (time-homogeneous), as they involve a varying policy  $\pi_{\theta_t}$ . This poses a major challenge for analyzing Algorithm 1, since a single sample is insufficient to accurately approximate the stationary distribution of the state under a fixed policy. Previous studies simplified their analysis by assuming i.i.d. samples drawn from the stationary distribution for the critic and from the visitation distribution for the actor. However, such sampling is infeasible in practice because both of them are unknown. In contrast, our approach is more practical since samples can be drawn directly from the Markov chain.

In Algorithm 1 Line 7, a projection  $(proj(\cdot))$  is introduced to keep the critic norm-bounded by  $\bar{\omega}$ , which is widely adopted in the literature (Wu et al., 2020; Chen et al., 2021; Olshevsky & Gharesifard, 2023; Chen & Zhao, 2024) for analysis. Note that the projection can be handled easily, relaxed using its non-expansive property in our analysis.

#### 4. Assumptions

Before presenting the main results, we will discuss several standard assumptions that are common in the literature of analyzing actor-critic with linear function approximation (Wu et al., 2020; Xu et al., 2020b; Shen et al., 2023; Chen et al., 2021; Olshevsky & Gharesifard, 2023; Chen & Zhao, 2024).

By taking the expectation of  $\omega_{t+1}$  in Eq. (6) with respect to the stationary distribution, and conditioning on  $\omega_t$ , we have

$$\mathbb{E}_{\boldsymbol{\theta}}[\boldsymbol{\omega}_{t+1} \,|\, \boldsymbol{\omega}_t] = \boldsymbol{\omega}_t + \beta(\boldsymbol{b}_{\boldsymbol{\theta}} - \boldsymbol{A}_{\boldsymbol{\theta}}\boldsymbol{\omega}_t),$$

where

$$\begin{aligned} \boldsymbol{A}_{\boldsymbol{\theta}} &:= \mathbb{E}_{(s,a,s')} \big[ \boldsymbol{\phi}(s) \big( \boldsymbol{\phi}(s) - \gamma \boldsymbol{\phi}(s') \big)^{\top} \big], \quad (14) \\ \boldsymbol{b}_{\boldsymbol{\theta}} &:= \mathbb{E}_{(s,a)} \big[ \boldsymbol{r}(s,a) \boldsymbol{\phi}(s) \big], \end{aligned}$$

and  $s \sim \mu_{\theta}(\cdot), a \sim \pi_{\theta}(\cdot | s), s' \sim P(\cdot | s, a)$  is the subsequent state of the (s, a). It can be easily shown that (Sutton & Barto, 2018) the TD limiting point  $\omega^*(\theta)$  satisfies:

$$\boldsymbol{A}_{\boldsymbol{\theta}}\boldsymbol{\omega}^*(\boldsymbol{\theta}) = \boldsymbol{b}_{\boldsymbol{\theta}}.$$
 (16)

**Assumption 4.1.** For any  $\theta$ , the matrix  $A_{\theta}$  defined in Eq. (14) is positive definite and its minimal eigenvalue can be lower bounded by  $\lambda$ .

Assumption 4.1 is commonly adopted in analyzing actorcritic (TD learning) with linear function approximation (Bhandari et al., 2018; Wu et al., 2020; Qiu et al., 2021; Chen et al., 2021; Olshevsky & Gharesifard, 2023; Chen & Zhao, 2024). It is explained as exploration since  $A_{\theta}$  can be rank deficient without sufficient exploration (Olshevsky & Gharesifard, 2023; Chen & Zhao, 2024). Assumption 4.1 further guarantees the problem's solvability since with this assumption, we have  $\omega^*(\theta) = A_{\theta}^{-1}b_{\theta}$ . In addition, we can choose  $\bar{\omega} = \bar{r}\lambda^{-1}$  so that all  $\omega^*$  lie within the projection radius  $\bar{\omega}$  because  $||b_{\theta}|| \leq \bar{r}$  and  $||A_{\theta}^{-1}|| \leq \lambda^{-1}$ , which justifies the projection operator introduced in Line 7 of Algorithm 1.

Assumption 4.2 (Uniform ergodicity). For any  $\theta$ , denote  $\mu_{\theta}(\cdot)$  as the stationary distribution induced by the policy  $\pi_{\theta}(\cdot | s)$  and the transition kernel  $P(\cdot | s, a)$ . For the following Markov chain (augmented with action) generated by the policy  $\pi_{\theta}$  and transition kernel P, i.e.,

$$s_0 \xrightarrow{(\pi_{\theta}, P)} s_1 \xrightarrow{(\pi_{\theta}, P)} \cdots \xrightarrow{(\pi_{\theta}, P)} s_t \xrightarrow{(\pi_{\theta}, P)} s_{t+1}, \quad (17)$$

there exist m > 0 and  $\rho \in (0, 1)$  such that

$$d_{TV} \big( \mathbb{P}(s_{\tau} \in \cdot \mid s_0 = s), \mu_{\boldsymbol{\theta}}(\cdot) \big) \le m \rho^{\tau}, \forall \tau \ge 0, \forall s \in \mathcal{S}$$

Assumption 4.2 assumes the Markov chain is geometrically mixing. It is commonly employed to characterize the noise induced by Markovian sampling in RL algorithms (Bhandari et al., 2018; Wu et al., 2020; Chen et al., 2021; Chen & Zhao, 2024). This is the counterpart of Proposition 3.2 (which is proved for analyzing the induced Markovian noise associated with the actor update). It is assumed since P is a general transition kernel that lacks the  $\gamma$ -contraction property of the transition kernel  $\hat{P}$  established in Proposition 3.2.

To justify this assumption in the continuous space, we note that all the distributions specified by the Ornstein–Uhlenbeck (OU) process satisfy this property. The OU process converges to a Gaussian distribution with the exponential mixing time. Moreover, it can also be shown that this property holds for more general diffusion processes (Del Moral & Villemonais, 2018).

Assumption 4.3 (Lipschitz continuity of policy). Let  $\pi_{\theta}(a \mid s)$  be a policy parameterized by  $\theta \in \mathcal{X}_{\Theta}$  with bounded support. There exist positive constants  $B, L_l$  and

L such that for any  $\theta$ ,  $\theta_1$ ,  $\theta_2 \in \mathcal{X}_{\Theta}$ ,  $s \in S$ , and  $a \in A$ , it holds that:

- (a)  $\|\nabla \log \pi_{\theta}(a \mid s)\| \leq B$ ,
- (b)  $\|\nabla \log \pi_{\theta_1}(a \,|\, s) \nabla \log \pi_{\theta_2}(a \,|\, s)\| \le L_l \|\theta_1 \theta_2\|,$
- (c)  $|\pi_{\theta_1}(a | s) \pi_{\theta_2}(a | s)| \le L \|\theta_1 \theta_2\|.$

Assumption 4.3 states the regularity of the policy which is standard in the literature of actor-critic methods (Xu et al., 2020a; Wu et al., 2020; Chen et al., 2021; Olshevsky & Gharesifard, 2023; Chen & Zhao, 2024). These conditions are sufficiently general to be satisfied by a wide range of distributions, including the uniform distribution, the truncated Gaussian distribution, and the Beta distribution with  $\alpha, \beta > 1$ .

With Assumption 4.3, we show in the following proposition that the policy  $\pi_{\theta}$  is Lipschitz continuous with respect to its parameter  $\theta$  in terms of the total variation distance.

**Proposition 4.4.** There exists a positive constant  $L_{\pi}$  such that for any  $\theta_1, \theta_2 \in \mathcal{X}_{\Theta}$ , it holds that

$$d_{TV}(\pi_{\boldsymbol{\theta}_1}(\cdot \mid s), \pi_{\boldsymbol{\theta}_2}(\cdot \mid s)) \le L_{\pi} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|.$$
(18)

We observed that Proposition 4.4 plays a key ingredient in the overall proof. With this proposition, we establish a bound on the distance between stationary distributions, as detailed in Lemma B.1 within Preliminary Lemmas in Appendix B, extending previous results from the tabular case to the continuous setting. This further facilitates the derivation of corresponding results for the discounted state visitation distribution, as presented in Lemma B.3.

#### 5. Main Results

We define the following uniform upper bound for the linear function approximation error of the critic:

$$\epsilon_{\mathrm{app}} := \sup_{\boldsymbol{\theta}} \sqrt{\mathbb{E}_{s \sim \nu_{\boldsymbol{\theta}}}(\boldsymbol{\phi}(s)^{\top} \boldsymbol{\omega}^{*}(\boldsymbol{\theta}) - V_{\boldsymbol{\theta}}(s))^{2}}.$$
 (19)

The error  $\epsilon_{app}$  is zero if  $V_{\theta}$  is indeed a linear function for any  $\theta$ . Naturally, it can be expected that the learning errors of Algorithm 1 depend on  $\epsilon_{app}$ .

We define the following integer  $\tau_{mix}$  that will be useful in the statement of the theorems:

$$\tau_{\min} := \min\left\{ i \ge 0 \mid m\rho^{i-1} \le \frac{1}{\sqrt{T}} \land \gamma^{i-1} \le \frac{1}{\sqrt{T}} \right\},\,$$

where  $m, \rho$  are constants defined in Assumption 4.2 and  $\gamma$  is the discounted factor. Therefore, we choose

$$\tau_{\rm mix} = \mathcal{O}(\log T)$$

such that  $m\rho^{\tau_{\min}-1} \leq 1/\sqrt{T}$  and  $\gamma^{\tau_{\min}-1} \leq 1/\sqrt{T}$ . The integer  $\tau_{\min}$  represents the mixing time of the ergodic Markov chain defined in Eq. (11) and Eq. (17), which will be used to control the Markovian noise in the analysis.

We define  $\Delta_t = \omega_t - \omega_t^*$  with  $\omega_t^* = \omega^*(\theta_t)$  to measure the critic error while  $\nabla J(\theta_t)$  serves as a measure of the actor error since for a general non-convex problem, our objective is to demonstrate that  $\nabla J(\theta_t)$  converges to zero.

**Theorem 5.1.** Consider Algorithm 1 with  $\alpha = c/\sqrt{T}$ ,  $\beta = 1/\sqrt{T}$ , where c is a constant depending on problem parameters. Suppose Assumptions 4.1-4.3 hold, we have for  $T \ge 2\tau_{\text{mix}}$ ,

$$\frac{1}{T - \tau_{\min}} \sum_{t=\tau_{\min}}^{T-1} \mathbb{E} \|\Delta_t\|^2 = \mathcal{O}\left(\frac{\log^2 T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\mathrm{app}}),$$
$$\frac{1}{T - \tau_{\min}} \sum_{t=\tau_{\min}}^{T-1} \mathbb{E} \|\nabla J(\boldsymbol{\theta}_t)\|^2 = \mathcal{O}\left(\frac{\log^2 T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\mathrm{app}}).$$

Theorem 5.1 establishes the finite-time convergence of Algorithm 1. If the critic approximation error  $\epsilon_{app}$  is zero, we see that both the critic error and the actor error diminish at a sub-linear rate of  $\tilde{\mathcal{O}}(T^{-1/2})$ . The additional logarithmic term  $(\log^2 T)$  is incurred by the mixing time of the Markov chain, which can be eliminated under i.i.d. sampling as will be shown in Proof Sketch. In terms of sample complexity, to obtain an  $\epsilon$ -approximate stationary point, it takes a number of  $\tilde{\mathcal{O}}(\epsilon^{-2})$  samples, which is typically the sample complexity of single-timescale actor-critic (Chen et al., 2021; Olshevsky & Gharesifard, 2023; Chen & Zhao, 2024).

The vanilla version of Algorithm 1 is introduced in the classic textbook (Sutton & Barto, 2018) as a canonical actorcritic algorithm with linear function approximation. As a canonical algorithm, its convergence has been a focal point of research, extensively studied across diverse settings, e.g., two-timescale (Wu et al., 2020), single-timescale (Chen et al., 2021; Olshevsky & Gharesifard, 2023; Chen & Zhao, 2024), time-average reward setting (Wu et al., 2020; Chen & Zhao, 2024), discounted setting (Chen et al., 2021; Olshevsky & Gharesifard, 2023). Notably, among all the aforementioned studies, this work is the first to address the important yet challenging setting of continuous state and action spaces. Among the widely used discounted singletimescale approaches considered, our method is the first to employ Markovian sampling for both the critic and the actor in contrast to the artificial i.i.d. sampling (Chen et al., 2021; Olshevsky & Gharesifard, 2023). Therefore, our work compares favorably by closing two significant gaps left by prior studies.

#### 5.1. Proof Sketch

To better illustrate our technical contribution, we provide a proof sketch to elucidate the significance of each error term and offer insights into the methods used to address them.

The key difference between single-timescale and twotimescale (double-loop) actor-critic lies in the strong coupling of the critic and actor errors. Unlike the two-timescale approach, which sequentially analyzes the convergence of critic error and actor error, the single-timescale setting requires simultaneous treatment of both errors. To address this, we propose a novel Lyapunov analysis framework and outline the proof of Theorem 5.1 in three steps. Step 1 derives an implicit upper bound for the critic error, treating it as an intermediate result. Step 2 performs a similar implicit analysis for the actor error. Step 3 combines these results into a novel Lyapunov function, whose convergence implies the simultaneous convergence of the critic and actor.

**Step 1: An implicit bound for critic error**. Using the critic update rule, we decompose the squared critic error by (see Eq. (27))

$$\mathbb{E} \|\Delta_{t+1}\|^{2} \leq \mathbb{E} \|\Delta_{t}\|^{2} + \underbrace{2\beta^{2}\mathbb{E} \|\boldsymbol{f}(O_{t},\boldsymbol{\omega}_{t})\|^{2}}_{I_{1}} + \underbrace{2\mathbb{E} \|\boldsymbol{\omega}_{t}^{*} - \boldsymbol{\omega}_{t+1}^{*}\|^{2}}_{I_{2}} + \underbrace{2\beta\mathbb{E} \langle\Delta_{t}, \boldsymbol{\bar{f}}(\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t})\rangle}_{I_{3}} + \underbrace{2\beta\mathbb{E} \langle\Delta_{t}, \boldsymbol{f}(O_{t},\boldsymbol{\omega}_{t}) - \boldsymbol{\bar{f}}(\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t})\rangle}_{I_{4}} + \underbrace{2\mathbb{E} \langle\Delta_{t}, \boldsymbol{\omega}_{t}^{*} - \boldsymbol{\omega}_{t+1}^{*}\rangle}_{I_{5}},$$

where  $f(O, \omega) = \delta(O, \omega)\phi(s)$  is the update term of the critic and  $\bar{f}$  is its mean value defined in Eq. (20).

 $I_1$  reflects the variance of the critic update which can be bounded by  $O(1/\sqrt{T})$  due to its bounded update.

 $I_2$  represents the difference between the moving critic target  $\omega_t^*$ , which can be controlled due to its Lipschitz continuity shown in Lemma B.6.

 $I_3$  is the inner product between the critic error  $\Delta_t$  and its mean-path update  $\bar{f}$ . It can be bounded by  $-2\lambda\beta\mathbb{E}||\Delta_t||^2$ under Assumption 4.1 since  $\omega^*$  is the solution of Eq. (16). Note that this bound combined with first term  $\mathbb{E}||\Delta_t||^2$  is  $(1 - 2\lambda\beta)\mathbb{E}||\Delta_t||^2$  which implies a contraction of the critic error because the coefficient is less than 1.

 $I_4$  represents the Markovian noise term, capturing the deviation between the critic's actual update f and its mean-path  $\bar{f}$ . To analyze this deviation, we aim to show that the sample  $O_t$  is close to its stationary distribution, as the error term  $I_3$  vanishes when  $O_t$  is drawn from the stationary distribution. First, we demonstrate that the sample  $O_t$  from the original Markov chain defined in Eq. (12) is close to the sample from the auxiliary Markov chain in Eq. (21), as their differences are limited to the last  $\tau$  steps. The total variation distance between these samples is controlled by the actor's change, i.e.,  $\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-\tau}\|$ , as established in Lemma B.2. By choosing  $\tau = \tau_{mix} = \mathcal{O}(\log T)$  and noting that the actor's update speed is  $\mathcal{O}(1/\sqrt{T})$ , the distance between  $\boldsymbol{\theta}_t$ and  $\theta_{t-\tau}$  is bounded by  $\mathcal{O}(\tau_{\min}/\sqrt{T})$ . Consequently, the accumulated deviation over the last  $au_{mix}$  steps is bounded by  $\mathcal{O}(\tau_{\text{mix}}^2/\sqrt{T}) = \mathcal{O}(\log^2 T/\sqrt{T})$ , explaining the logarithmic term in the convergence rate. Next, we show that the Markov noise of the sample from the auxiliary Markov chain approaches its stationary distribution after  $au_{
m mix}$  steps under a fixed policy, leveraging the uniform ergodicity assumption in Assumption 4.2. This highlights the role of Assumption 4.2 in analyzing single-sample single/two-timescale algorithms with Markovian sampling. The complete analysis of this Markovian noise is shown in Lemma C.1.

 $I_5$  tracks both the critic error  $\Delta_t$  and the difference between the drifting critic targets  $\omega_t^*$ . It can be bounded by the critic error  $\Delta_t$  and the actor error  $\nabla J(\boldsymbol{\theta}_t)$  after error decomposition. In contrast, the two-timescale setting can prove that  $I_4$  converges to zero. To see why this is the case, note that from the Lipschitz continuity of the critic target  $\omega_t^*$  shown in Lemma B.6, error term  $\omega_t^* - \omega_{t+1}^*$  can be bounded by the update of the actor  $\theta$ , i.e.,  $\theta_t - \theta_{t+1}$ . Since the update step size for the actor is  $\alpha$  while the contraction of the critic error is at a rate  $1 - 2\lambda\beta$ , a ratio term  $\alpha/\beta$ appears by moving the term  $-2\lambda\beta\mathbb{E}\|\Delta_t\|^2$  to the left side of the above inequality and dividing its coefficient. Therefore, one can leave other terms in  $I_4$  as constant and bound it by  $\mathcal{O}(\alpha/\beta)$ . Since  $\lim_{t\to\infty} \alpha_t/\beta_t = 0$  in two-timescale approach, thereby directly establish the convergence of the critic. However,  $\lim_{t\to\infty} \alpha_t/\beta_t = c$  in single-timescale approach which is why we can only bound  $I_4$  by  $\Delta_t$  and  $\nabla J(\boldsymbol{\theta})$  and make an implicit upper bound for the critic error. The final bound is summarized in Theorem D.1.

**Step 2:** An implicit bound for actor error. Using the actor update rule and the smoothness property of  $J(\theta)$  (Lemma B.8), we decompose the squared actor error by (see Eq. (31))

$$(1-\gamma)\mathbb{E}\|\nabla J(\boldsymbol{\theta}_{t})\|^{2} \leq \frac{1}{\alpha} \left(\mathbb{E}\left[J(\boldsymbol{\theta}_{t+1}) - J(\boldsymbol{\theta}_{t})\right]\right) \\ + \underbrace{\frac{\alpha L_{g}}{2}\mathbb{E}\|\boldsymbol{h}(\widehat{O}_{t},\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t})\|^{2}}_{I_{1}} - \underbrace{\mathbb{E}\langle\nabla J(\boldsymbol{\theta}_{t}), \bar{\boldsymbol{g}}(\boldsymbol{\omega}_{t}^{*},\boldsymbol{\theta}_{t})\rangle}_{I_{2}} \\ + \underbrace{\mathbb{E}\langle\nabla J(\boldsymbol{\theta}_{t}), \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t}) - \boldsymbol{h}(\widehat{O}_{t},\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t})\rangle}_{I_{3}} \\ + \underbrace{\mathbb{E}\langle\nabla J(\boldsymbol{\theta}_{t}), \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}^{*},\boldsymbol{\theta}_{t}) - \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t})\rangle}_{I_{4}},$$

where  $h(O, \omega, \theta) = \delta(O, \omega) \nabla \log \pi_{\theta}(a | s)$  is the update term of the actor,  $\bar{h}$  is its mean value, and  $\bar{g}(\omega_t^*, \theta_t)$  defined in Eq. (20) represents the approximation error of the optimal critic  $\omega_t^*$ . The first term on the right-hand side of the above inequality compares the actor's performances between consecutive updates, which can be eliminated by telescoping.

 $I_1$  reflects the variance of the actor update which can be controlled by  $\mathcal{O}(1/\sqrt{T})$  due to its bounded update.

 $I_2$  is the inner product between actor error and the approximation error of the optimal critic  $\omega_t^*$ . This term is control by the approximation error  $\mathcal{O}(\epsilon_{app})$  defined in Eq. (19).

 $I_3$  represents the Markovian noise term, capturing the deviation between the actor's actual update h and its mean-path  $\bar{h}$ . Similar to the critic analysis, this noise is controlled by showing that the original Markov chain defined in Eq. (13) is close to the auxiliary Markov chain in Eq. (22). Additionally, samples from the auxiliary Markov chain approach their stationary distribution after  $\tau_{mix}$  steps, leveraging the uniform ergodicity property established in Proposition 3.2. This error term is bounded in Lemma C.3.

 $I_4$  tracks the inner product between the actor error  $\nabla J(\theta)$ and the critic error ( $\Delta_t = \omega - \omega_t^*$ ). In two-timescale actor-critic, this term goes to zero due to the convergence of the critic. However, in single-timescale approach, we can only bound this term by  $\nabla J(\theta_t)$  and  $\Delta_t$  which will be treated together later. This give an implicit upper bound for the actor error. The final result of the above inequality is summarized in Theorem D.2.

**Step 3: A novel Lyapunov analysis**. From Step 1 and Step 2, we get two inequalities about the coupled critic error and actor error. Here we bring them together by defining the following Lyapunov function

$$\mathbb{L}_t = \frac{2B}{1-\gamma} \mathbb{E} \|\Delta_t\|^2 + \frac{1-\gamma}{2B} \mathbb{E} \|\nabla J(\boldsymbol{\theta}_t)\|^2$$

where  $2B/(1-\gamma)$  is the scaling coefficient which balances the contribution of the critic error and the actor error. Combine the results in Step 1 and Step 2 (Eq. (26) and Eq. (30)) gives an unified inequality of  $\mathbb{L}_t$ . We then define the total error as  $\mathcal{L} := 1/(T - \tau_{\min}) \sum_{t=\tau_{\min}}^{T-1} \mathbb{L}_t$ . Telescoping from  $t = T - \tau_{\min}$  to T - 1, it can be shown that (see Eq. (33))

$$\mathcal{L} \leq \left(\frac{2L_c Bc}{\lambda} + \frac{1}{2}\right) \mathcal{L} + \mathcal{O}\left(\frac{\log^2 T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\mathrm{app}}),$$

where  $c = \alpha/\beta$  is the stepsize ratio between the actor and the critic,  $\gamma$  is the discounted factor,  $\lambda$  is the maximum eigenvalue of  $A_{\theta}$  defined in Assumption 4.1, and  $L_c$  is the Lipschitz constant characterized in Lemma B.6. Therefore, choosing  $c < \lambda/4BL_c$  (see Eq. (34)), we have

$$\mathcal{L} = \mathcal{O}\left(\frac{\log^2 T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\rm app}),$$

which implies the convergence of the critic error and the actor error simultaneously. Therefore, we finish the proof of Theorem 5.1.

## 6. Conclusion

In this paper, we provide a finite-time convergence analysis for the single-sample, single-timescale actor-critic algorithm in continuous state-action spaces. We propose a novel Lyapunov analysis framework, which allows a less conservative analysis under the same set of assumptions adopted in existing studies. Our analysis offers new insights into the coupled dynamics of actor-critic updates. Unlike prior works that assume artificial decoupling between the actor and critic, our results capture the interdependencies that arise naturally in practical implementations. Moreover, our framework and analytical techniques can serve as a foundation for studying other single-timescale reinforcement learning algorithms in continuous domains.

#### Acknowledgements

This work was supported by the Singapore Ministry of Education Tier 1 Academic Research Fund (A-8001174-00-00) and Tier 2 Academic Research Funds (T2EP20123-0037, T2EP20224-0035).

#### **Impact Statement**

This paper presents work whose goal is to advance the field of Machine Learning. There are no potential societal consequences of our work.

#### References

- Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98):1–76, 2021.
- Bhandari, J., Russo, D., and Singal, R. A finite time analysis of temporal difference learning with linear function approximation. In *Conference on learning theory*, pp. 1691–1692. PMLR, 2018.
- Bhatnagar, S., Sutton, R. S., Ghavamzadeh, M., and Lee, M. Natural actor–critic algorithms. *Automatica*, 45(11): 2471–2482, 2009.
- Chen, T., Sun, Y., and Yin, W. Closing the gap: Tighter analysis of alternating stochastic gradient methods for bilevel

problems. Advances in Neural Information Processing Systems, 34:25294–25307, 2021.

- Chen, X. and Zhao, L. Finite-time analysis of singletimescale actor-critic. Advances in Neural Information Processing Systems, 36, 2024.
- Chen, X., Duan, J., Liang, Y., and Zhao, L. Global convergence of two-timescale actor-critic for solving linear quadratic regulator. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 7087–7095, 2023.
- Chen, X., Duan, J., and Zhao, L. Global optimality of singletimescale actor-critic under continuous state-action space: a study on linear quadratic regulator. In *Proceedings* of the Thirty-Third International Joint Conference on Artificial Intelligence, pp. 3816–3824, 2024.
- Del Moral, P. and Villemonais, D. Exponential mixing properties for time inhomogeneous diffusion processes with killing. 2018.
- Fu, Z., Yang, Z., and Wang, Z. Single-timescale actor-critic provably finds globally optimal policy. arXiv preprint arXiv:2008.00483, 2020.
- Gaur, M., Bedi, A., Wang, D., and Aggarwal, V. Closing the gap: Achieving global convergence (last iterate) of actorcritic under markovian sampling with neural network parametrization. In *International Conference on Machine Learning*, pp. 15153–15179. PMLR, 2024.
- Heidergott, B. and Hordijk, A. Taylor series expansions for stationary markov chains. *Advances in Applied Probability*, 35(4):1046–1070, 2003.
- Hoeller, D., Rudin, N., Sako, D., and Hutter, M. Anymal parkour: Learning agile navigation for quadrupedal robots. *Science Robotics*, 9(88):eadi7566, 2024. doi: 10.1126/scirobotics.adi7566.
- Hong, M., Wai, H.-T., Wang, Z., and Yang, Z. A twotimescale stochastic algorithm framework for bilevel optimization: Complexity analysis and application to actorcritic. *SIAM Journal on Optimization*, 33(1):147–180, 2023.
- Kakade, S. M. A natural policy gradient. Advances in neural information processing systems, 14, 2001.
- Kaufmann, E., Bauersfeld, L., Loquercio, A., Müller, M., Koltun, V., and Scaramuzza, D. Champion-level drone racing using deep reinforcement learning. *Nature*, 620 (7976):982–987, 2023. ISSN 1476-4687. doi: 10.1038/ s41586-023-06419-4.

- Konda, V. and Tsitsiklis, J. Actor-critic algorithms. *Advances in neural information processing systems*, 12, 1999.
- Konda, V. R. and Tsitsiklis, J. N. Onactor-critic algorithms. SIAM journal on Control and Optimization, 42(4):1143– 1166, 2003.
- Kumar, H., Koppel, A., and Ribeiro, A. On the sample complexity of actor-critic method for reinforcement learning with function approximation. *Machine Learning*, 112(7): 2433–2467, 2023.
- Lazaridis, A., Fachantidis, A., and Vlahavas, I. Deep reinforcement learning: A state-of-the-art walkthrough. *Journal of Artificial Intelligence Research*, 69:1421–1471, 2020.
- Mitrophanov, A. Y. Sensitivity and convergence of uniformly ergodic markov chains. *Journal of Applied Probability*, 42(4):1003–1014, 2005.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pp. 1928– 1937. PMLR, 2016.
- Olshevsky, A. and Gharesifard, B. A small gain analysis of single timescale actor critic. *SIAM Journal on Control and Optimization*, 61(2):980–1007, 2023.
- Qiu, S., Yang, Z., Ye, J., and Wang, Z. On finite-time convergence of actor-critic algorithm. *IEEE Journal on Selected Areas in Information Theory*, 2(2):652–664, 2021.
- Radosavovic, I., Xiao, T., Zhang, B., Darrell, T., Malik, J., and Sreenath, K. Real-world humanoid locomotion with reinforcement learning. *Science Robotics*, 9(89): eadi9579, 2024. doi: 10.1126/scirobotics.adi9579.
- Shen, H., Zhang, K., Hong, M., and Chen, T. Towards understanding asynchronous advantage actor-critic: Convergence and linear speedup. *IEEE Transactions on Signal Processing*, 71:2579–2594, 2023.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Sutton, R. S., McAllester, D., Singh, S., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information* processing systems, 12, 1999.

- Vinyals, O., Babuschkin, I., Czarnecki, W. M., Mathieu, M., Dudzik, A., Chung, J., Choi, D. H., Powell, R., Ewalds, T., Georgiev, P., et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *nature*, 575 (7782):350–354, 2019.
- Wang, L., Cai, Q., Yang, Z., and Wang, Z. Neural policy gradient methods: Global optimality and rates of convergence. arXiv preprint arXiv:1909.01150, 2019.
- Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
- Wu, Y. F., Zhang, W., Xu, P., and Gu, Q. A finite-time analysis of two time-scale actor-critic methods. *Advances* in Neural Information Processing Systems, 33:17617– 17628, 2020.
- Xu, P., Gao, F., and Gu, Q. An improved convergence analysis of stochastic variance-reduced policy gradient. In *Uncertainty in Artificial Intelligence*, pp. 541–551. PMLR, 2020a.
- Xu, T., Wang, Z., and Liang, Y. Improving sample complexity bounds for (natural) actor-critic algorithms. *Ad*vances in Neural Information Processing Systems, 33: 4358–4369, 2020b.
- Xu, T., Wang, Z., and Liang, Y. Non-asymptotic convergence analysis of two time-scale (natural) actor-critic algorithms. *arXiv preprint arXiv:2005.03557*, 2020c.
- Yang, Z., Chen, Y., Hong, M., and Wang, Z. Provably global convergence of actor-critic: A case for linear quadratic regulator with ergodic cost. *Advances in neural information processing systems*, 32, 2019.
- Zhang, K., Koppel, A., Zhu, H., and Basar, T. Global convergence of policy gradient methods to (almost) locally optimal policies. *SIAM Journal on Control and Optimiza*tion, 58(6):3586–3612, 2020a.
- Zhang, S., Liu, B., Yao, H., and Whiteson, S. Provably convergent two-timescale off-policy actor-critic with function approximation. In *International Conference on Machine Learning*, pp. 11204–11213. PMLR, 2020b.
- Zhou, M. and Lu, J. Single timescale actor-critic method to solve the linear quadratic regulator with convergence guarantees. *Journal of Machine Learning Research*, 24 (222):1–34, 2023.
- Zou, S., Xu, T., and Liang, Y. Finite-sample analysis for sarsa with linear function approximation. *Advances in neural information processing systems*, 32, 2019.

# **Supplementary Material**

# **Table of Contents**

<ul> <li>A Notation</li> <li>B Preliminary Lemmas</li> <li>C Markovian Noise</li> <li>D Proof of Main Theorem <ul> <li>D.1 An implicit bound for critic error</li> <li>D.2 An implicit bound for actor error</li> <li>D.3 A novel Lyapunov analysis</li> </ul> </li> <li>E Proof of Propositions</li> <li>F Proof of Preliminary Lemmas</li> <li>G Proof of Markovian Noise</li> </ul>			
<ul> <li>B Preliminary Lemmas</li> <li>C Markovian Noise</li> <li>D Proof of Main Theorem <ul> <li>D.1 An implicit bound for critic error</li> <li>D.2 An implicit bound for actor error</li> <li>D.3 A novel Lyapunov analysis</li> </ul> </li> <li>E Proof of Propositions</li> <li>F Proof of Preliminary Lemmas</li> <li>G Proof of Markovian Noise</li> </ul>	A	Notation	11
<ul> <li>C Markovian Noise</li> <li>Proof of Main Theorem <ul> <li>D.1 An implicit bound for critic error</li> <li>D.2 An implicit bound for actor error</li> <li>D.3 A novel Lyapunov analysis</li> <li>C Proof of Propositions</li> </ul> </li> <li>F Proof of Preliminary Lemmas</li> <li>G Proof of Markovian Noise</li> </ul>	B	Preliminary Lemmas	12
<ul> <li>D Proof of Main Theorem</li> <li>D.1 An implicit bound for critic error</li></ul>	С	Markovian Noise	13
<ul> <li>D.1 An implicit bound for critic error</li></ul>	D	Proof of Main Theorem	14
<ul> <li>D.2 An implicit bound for actor error</li></ul>		D.1 An implicit bound for critic error	14
<ul> <li>D.3 A novel Lyapunov analysis</li> <li>F Proof of Prepininary Lemmas</li> <li>G Proof of Markovian Noise</li> </ul>		D.2 An implicit bound for actor error	16
<ul> <li>E Proof of Propositions</li> <li>F Proof of Preliminary Lemmas</li> <li>G Proof of Markovian Noise</li> </ul>		D.3 A novel Lyapunov analysis	17
<ul><li>F Proof of Preliminary Lemmas</li><li>G Proof of Markovian Noise</li></ul>	E	Proof of Propositions	19
G Proof of Markovian Noise	F	Proof of Preliminary Lemmas	21
	G	Proof of Markovian Noise	25

# A. Notation

In the following, we will analyze the convergence of the above algorithm. We define the following notations:

$$f(O, \omega) := (r(s, a) + \gamma \phi(s')^{\top} \omega - \phi(s)^{\top} \omega) \phi(s)$$
  

$$\bar{f}(\omega, \theta) := \mathbb{E}_{O \sim (\mu_{\theta}, \pi_{\theta}, P)} [f(O, \omega)],$$
  

$$h(O, \omega, \theta) := (r(s, a) + \gamma \phi(s')^{\top} \omega - \phi(s)^{\top} \omega) \nabla \log \pi_{\theta}(a \mid s)$$
  

$$\bar{h}(\omega, \theta) := \mathbb{E}_{O \sim (\nu_{\theta}, \pi_{\theta}, P)} [h(O, \omega, \theta)],$$
  

$$g(O, \omega, \theta) := ((\gamma \phi(s') - \phi(s))^{\top} \omega - (\gamma V_{\theta}(s') - V_{\theta}(s))) \nabla \log \pi_{\theta}(a \mid s),$$
  

$$\bar{g}(\omega, \theta) := \mathbb{E}_{O \sim (\nu_{\theta}, \pi_{\theta}, P)} [g(O, \omega, \theta)].$$
  
(20)

We make use of the following auxiliary Markov chain to deal with the Markovian noise.

#### Auxiliary Markov Chain for the Critic:

$$s_{t-\tau} \xrightarrow{\boldsymbol{\theta}_{t-\tau}} a_{t-\tau} \xrightarrow{P} s_{t-\tau+1} \xrightarrow{\boldsymbol{\theta}_{t-\tau}} \tilde{a}_{t-\tau+1} \xrightarrow{P} \tilde{s}_{t-\tau+2} \xrightarrow{\boldsymbol{\theta}_{t-\tau}} \tilde{a}_{t-\tau+2} \cdots \xrightarrow{P} \tilde{s}_t \xrightarrow{\boldsymbol{\theta}_{t-\tau}} \tilde{a}_t \xrightarrow{P} \tilde{s}_{t+1}.$$
(21)

## Auxiliary Markov Chain for the Actor:

$$\hat{s}_{t-\tau} \xrightarrow{\boldsymbol{\theta}_{t-\tau}} \hat{a}_{t-\tau} \xrightarrow{\hat{P}} \hat{s}_{t-\tau+1} \xrightarrow{\boldsymbol{\theta}_{t-\tau}} \bar{a}_{t-\tau+1} \xrightarrow{\hat{P}} \bar{s}_{t-\tau+2} \xrightarrow{\boldsymbol{\theta}_{t-\tau}} \bar{a}_{t-\tau+2} \cdots \xrightarrow{\hat{P}} \bar{s}_t \xrightarrow{\boldsymbol{\theta}_{t-\tau}} \bar{a}_t \xrightarrow{\hat{P}} \bar{s}_{t+1}.$$
(22)

In the sequel, we denote by  $\tilde{O}_t := (\tilde{s}_t, \tilde{a}_t, \tilde{s}_{t+1})$  the tuple generated from the auxiliary Markov chain in Eq. (21) and  $\bar{O}_t := (\bar{s}_t, \bar{a}_t, \bar{s}_{t+1})$  the tuple generated from the auxiliary Markov chain in Eq. (22). In comparison,  $O_t := (s_t, a_t, s_{t+1})$  and  $\hat{O}_t := (\hat{s}_t, \hat{a}_t, \hat{s}_{t+1})$  denotes the tuple generated by Algorithm 1. We use O' as a shorthand for an independent sample

from stationary distribution  $s \sim \mu_{\theta}, a \sim \pi_{\theta}, s' \sim P$  and use O'' as a shorthand for an independent sample from discounted state visitation distribution  $s \sim \nu_{\theta}, a \sim \pi_{\theta}, s' \sim P$ .

Throughout the proof, we define  $\delta_t := \delta(O_t, \omega_t)$ , and  $\overline{\delta} = \overline{r} + 2\overline{\omega}$  is the uniform upper bound for  $\delta$ . We also define a filtration  $\mathcal{F}_t = \sigma(s_0, \widehat{s}_0, a_0, \widehat{a}_0, s_1, \widehat{s}_1, a_1, \widehat{a}_1, \cdots, s_t, \widehat{s}_t)$ , where  $\sigma(\cdot)$  denotes the  $\sigma$ -algebra generated by the random variables.

## **B.** Preliminary Lemmas

In this section, we present several preliminary lemmas, encompassing three aspects: extending previous work to continuous settings (Lemma B.1, Lemma B.2, Lemma B.5, Lemma B.6), establishing the corresponding statistical properties for actor samples (Lemma B.3, Lemma B.4), and stating previously established results (Lemma B.7, Lemma B.8, Lemma B.9).

**Lemma B.1.** For any  $\theta_1$  and  $\theta_2$ , it holds that

$$d_{TV}(\mu_{\theta_1}, \mu_{\theta_2}) \leq 2L_{\pi} \left( \lceil \log_{\rho} m^{-1} \rceil + \frac{1}{1-\rho} \right) \|\theta_1 - \theta_2\|,$$
  
$$d_{TV}(\mu_{\theta_1} \otimes \pi_{\theta_1}, \mu_{\theta_2} \otimes \pi_{\theta_2}) \leq 2L_{\pi} \left( 1 + \lceil \log_{\rho} m^{-1} \rceil + \frac{1}{1-\rho} \right) \|\theta_1 - \theta_2\|,$$
  
$$d_{TV}(\mu_{\theta_1} \otimes \pi_{\theta_1} \otimes P, \mu_{\theta_2} \otimes \pi_{\theta_2} \otimes P) \leq 2L_{\pi} \left( 1 + \lceil \log_{\rho} m^{-1} \rceil + \frac{1}{1-\rho} \right) \|\theta_1 - \theta_2\|.$$

**Lemma B.2.** Given time indexes t and  $\tau$  such that  $t \ge \tau > 0$ , consider the auxiliary Markov chain in Eq. (21). Conditioning on  $\mathcal{F}_{t-\tau}$ , we have

$$\begin{aligned} d_{TV} \big( \mathbb{P}(s_{t+1} \in \cdot), \mathbb{P}(\widetilde{s}_{t+1} \in \cdot) \big) &\leq d_{TV} \big( \mathbb{P}(O_t \in \cdot), \mathbb{P}(\widetilde{O}_t \in \cdot) \big), \\ d_{TV} \big( \mathbb{P}(O_t \in \cdot), \mathbb{P}(\widetilde{O}_t \in \cdot) \big) &= d_{TV} \big( \mathbb{P}((s_t, a_t) \in \cdot), \mathbb{P}((\widetilde{s}_t, \widetilde{a}_t) \in \cdot) \big), \\ d_{TV} \big( \mathbb{P}((s_t, a_t) \in \cdot), \mathbb{P}((\widetilde{s}_t, \widetilde{a}_t) \in \cdot) \big) &\leq d_{TV} \big( \mathbb{P}(s_t \in \cdot), \mathbb{P}(\widetilde{s}_t \in \cdot) \big) + L_{\pi} \mathbb{E} \big[ \| \boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-\tau} \| \big] \end{aligned}$$

**Lemma B.3.** For any  $\theta_1$  and  $\theta_2$ , it holds that

$$d_{TV}(\nu_{\theta_1}, \nu_{\theta_2}) \leq \frac{2L_{\pi}}{1-\gamma} \| \theta_1 - \theta_2 \|,$$
  
$$d_{TV}(\nu_{\theta_1} \otimes \pi_{\theta_1}, \nu_{\theta_2} \otimes \pi_{\theta_2}) \leq 2L_{\pi} \left( 1 + \frac{1}{1-\gamma} \right) \| \theta_1 - \theta_2 \|,$$
  
$$d_{TV}(\nu_{\theta_1} \otimes \pi_{\theta_1} \otimes P, \nu_{\theta_2} \otimes \pi_{\theta_2} \otimes P) \leq 2L_{\pi} \left( 1 + \frac{1}{1-\gamma} \right) \| \theta_1 - \theta_2 \|.$$

**Lemma B.4.** Given time indexes t and  $\tau$  such that  $t \ge \tau > 0$ , consider the auxiliary Markov chain in Eq. (22). Conditioning on  $\mathcal{F}_{t-\tau}$ , we have

$$\begin{aligned} d_{TV} \big( \mathbb{P}(\hat{s}_{t+1} \in \cdot), \mathbb{P}(\bar{s}_{t+1} \in \cdot) \big) &\leq d_{TV} \big( \mathbb{P}(\widehat{O}_t \in \cdot), \mathbb{P}(\bar{O}_t \in \cdot) \big), \\ d_{TV} \big( \mathbb{P}(\widehat{O}_t \in \cdot), \mathbb{P}(\bar{O}_t \in \cdot) \big) &= d_{TV} \big( \mathbb{P}((\hat{s}_t, \hat{a}_t) \in \cdot), \mathbb{P}((\bar{s}_t, \bar{a}_t) \in \cdot) \big), \\ d_{TV} \big( \mathbb{P}((\hat{s}_t, \hat{a}_t) \in \cdot), \mathbb{P}((\bar{s}_t, \bar{a}_t) \in \cdot) \big) &\leq d_{TV} \big( \mathbb{P}(\hat{s}_t \in \cdot), \mathbb{P}(\bar{s}_t \in \cdot)) \big) + L_{\pi} \mathbb{E} \big[ \| \boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-\tau} \| \big] \end{aligned}$$

**Lemma B.5.** For any  $\theta_1, \theta_2$ , we have

$$|J(\boldsymbol{\theta}_1) - J(\boldsymbol{\theta}_2)| \le L_J \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|,$$

where  $L_J = 4\bar{r}L_{\pi}(1 + (1 - \gamma)^{-1}).$ 

**Lemma B.6.** There exists a constant  $L_c > 0$  such that

 $\|\boldsymbol{\omega}^*(\boldsymbol{\theta}_1) - \boldsymbol{\omega}^*(\boldsymbol{\theta}_2)\| \leq L_c \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|, \forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^d,$ 

where  $L_c = (8\lambda^{-2}\bar{r} + 4\lambda^{-1}\bar{r})L_{\pi}(1 + \lceil \log_{\rho} m^{-1} \rceil + 1/(1-\rho)).$ 

**Lemma B.7.** For any  $\theta, \theta' \in \mathbb{R}^d$ , there exists constant  $L_{\mu}$  such that  $\|\nabla \mu_{\theta} - \nabla \mu_{\theta'}\| \leq L_{\mu} \|\theta - \theta'\|$ , where  $\mu_{\theta}(s)$  is the stationary distribution under the policy  $\pi_{\theta}$ .

**Lemma B.8** ((Zhang et al., 2020a), Lemma 3.2). For the performance function  $J(\theta)$ , there exists a constant  $L_g > 0$  such that for all  $\theta_1, \theta_2 \in \mathbb{R}^d$ , it holds that

$$\|\nabla J(\boldsymbol{\theta}_1) - \nabla J(\boldsymbol{\theta}_2)\| \le L_q \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|,\tag{23}$$

which further implies

$$J(\boldsymbol{\theta}_2) \ge J(\boldsymbol{\theta}_1) + \langle \nabla J(\boldsymbol{\theta}_1), \boldsymbol{\theta}_2 - \boldsymbol{\theta}_1 \rangle - \frac{L_g}{2} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|^2.$$
(24)

**Lemma B.9** ((Chen et al., 2021), Proposition 8). For any  $\theta_1, \theta_2 \in \mathbb{R}^d$ , we have

$$\|
abla oldsymbol{\omega}^*(oldsymbol{ heta}_1) - 
abla oldsymbol{\omega}^*(oldsymbol{ heta}_2)\| \leq L_s \|oldsymbol{ heta}_1 - oldsymbol{ heta}_2\|_2$$

where  $L_s$  is a positive constant.

# C. Markovian Noise

We then the following Markovian noise term

$$\Lambda(O, \boldsymbol{\omega}, \boldsymbol{\theta}) = \langle \boldsymbol{\omega} - \boldsymbol{\omega}^*, \boldsymbol{f}(O, \boldsymbol{\omega}) - \bar{\boldsymbol{f}}(\boldsymbol{\omega}, \boldsymbol{\theta}) \rangle,$$
  

$$\Gamma(O, \boldsymbol{\omega}, \boldsymbol{\theta}) = \langle \boldsymbol{\omega} - \boldsymbol{\omega}^*, (\nabla \boldsymbol{\omega}^*)^\top (\bar{\boldsymbol{h}}(\boldsymbol{\omega}^*, \boldsymbol{\theta}) - \boldsymbol{h}(O, \boldsymbol{\omega}^*, \boldsymbol{\theta})) \rangle,$$
  

$$\Xi(O, \boldsymbol{\omega}, \boldsymbol{\theta}) = \langle \nabla J(\boldsymbol{\theta}), \bar{\boldsymbol{h}}(\boldsymbol{\omega}, \boldsymbol{\theta}) - \boldsymbol{h}(O, \boldsymbol{\omega}, \boldsymbol{\theta}) \rangle.$$
(25)

**Lemma C.1.** For any  $t \ge \tau_{\text{mix}}$ , the Markovian noise in the critic update, denoted by  $\Lambda(O_t, \boldsymbol{\omega}_t, \boldsymbol{\theta}_t)$ , satisfies

$$\mathbb{E}\big[\Lambda(O_t,\boldsymbol{\omega}_t,\boldsymbol{\theta}_t)\big] \leq M_1 \frac{1}{\sqrt{T}},$$

where  $M_1 = (8\bar{\omega}\bar{\delta}L_{\pi}(1+\lceil \log_{\rho}m^{-1}\rceil+(1-\rho)^{-1})+2\bar{\delta}L_c)\bar{\delta}B\tau_{\min}c + (8\bar{\omega}+2\bar{\delta})\bar{\delta}\tau_{\min}+4\bar{\omega}L_{\pi}B\bar{\delta}^2\tau_{\min}^2c + 4\bar{\omega}\bar{\delta}.$ Lemma C.2. For any  $t \ge \tau_{\min} > 0$ , it holds that

$$\mathbb{E}\big[\Gamma(\widehat{O}_t, \boldsymbol{\omega}_t, \boldsymbol{\theta}_t)\big] \leq M_2 \frac{1}{\sqrt{T}},$$

where

$$M_{2} = (2\bar{\delta}BL_{c}^{2} + 4\bar{\delta}\bar{\omega}BL_{s} + 4\bar{\omega}L_{c}L_{\bar{h}})\tau_{\mathrm{mix}}\bar{\delta}Bc + 2\bar{\delta}BL_{c}\tau_{\mathrm{mix}}\bar{\delta} + 4\bar{\omega}\bar{\delta}BL_{c}L_{\pi}\tau_{\mathrm{mix}}^{2}\bar{\delta}Bc + 4\bar{\omega}\bar{\delta}BL_{c}$$
$$L_{\bar{h}} = \bar{\delta}L_{l} + 2BL_{c} + 4\bar{\delta}BL_{\pi}\left(1 + \frac{1}{1 - \gamma}\right).$$

**Lemma C.3.** For any  $t \ge \tau_{mix} > 0$ , it can be shown that

$$\mathbb{E}\big[\Xi(\widehat{O}_t, \boldsymbol{\omega}_t, \boldsymbol{\theta}_t)\big] \le M_3 \frac{1}{\sqrt{T}},$$

where

$$M_{3} = (2\delta BL_{g} + 2L_{J}L_{\bar{h}})\delta Bc\tau_{\rm mix} + 4BL_{J}\delta\tau_{\rm mix} + 2\delta^{2}B^{2}L_{J}L_{\pi}c\tau_{\rm mix}^{2} + 2\delta BL_{J},$$
  
$$L_{\bar{h}} = \bar{\delta}L_{l} + 2BL_{c} + 4\bar{\delta}BL_{\pi}\left(1 + \frac{1}{1 - \gamma}\right).$$

# **D.** Proof of Main Theorem

## D.1. An implicit bound for critic error

**Theorem D.1.** Choose  $\alpha_t = c/\sqrt{T}$ ,  $\beta_t = 1/\sqrt{T}$ , for any  $\tau_{mix} \leq t < T$ , we have

$$\mathbb{E}\|\Delta_t\|^2 \le \frac{1}{\lambda\beta} (\mathbb{E}\|\Delta_t\|^2 - \mathbb{E}\|\Delta_{t+1}\|^2) + \frac{2c(1-\gamma)}{\lambda} L_c \mathbb{E}\|\Delta_t\| \|\nabla J(\boldsymbol{\theta}_t)\| + \mathcal{O}(\frac{\log^2 T}{\sqrt{T}}) + \mathcal{O}(\epsilon_{\mathrm{app}}).$$
(26)

Proof. From the update rule of the critic in Line 6 of Algorithm 1, we have

$$\begin{split} \|\boldsymbol{\omega}_{t+1} - \boldsymbol{\omega}_{t+1}^*\| &= \|proj_{\bar{\omega}}(\boldsymbol{\omega}_t + \beta\delta_t\boldsymbol{\phi}(s_t)) - \boldsymbol{\omega}_{t+1}^*\| \\ &= \|proj_{\bar{\omega}}(\boldsymbol{\omega}_t + \beta\delta_t\boldsymbol{\phi}(s_t)) - proj_{\bar{\omega}}(\boldsymbol{\omega}_{t+1}^*)\| \\ &\stackrel{(1)}{\leq} \|\boldsymbol{\omega}_t + \beta\delta_t\boldsymbol{\phi}(s_t) - \boldsymbol{\omega}_{t+1}^*\| \\ &= \|\boldsymbol{\omega}_t - \boldsymbol{\omega}_t^* + \beta\delta_t\boldsymbol{\phi}(s_t) + \boldsymbol{\omega}_t^* - \boldsymbol{\omega}_{t+1}^*\|, \end{split}$$

where (1) holds because the projection function  $proj_{\bar{\omega}}(\cdot)$  is 1-Lipschitz continuous. It follows that

$$\begin{split} \|\Delta_{t+1}\|^2 &= \|\Delta_t + \beta \delta_t \phi(s_t) + \omega_t^* - \omega_{t+1}^* \|^2 \\ &= \|\Delta_t\|^2 + \|\beta \delta_t \phi(s_t) + \omega_t^* - \omega_{t+1}^* \|^2 \\ &+ 2 \langle \Delta_t, \beta \delta_t \phi(s_t) \rangle + 2 \langle \Delta_t, \omega_t^* - \omega_{t+1}^* \rangle \\ &= \|\Delta_t\|^2 + \|\beta f(O_t, \omega_t) + \omega_t^* - \omega_{t+1}^* \|^2 \\ &+ 2\beta \langle \Delta_t, f(O_t, \omega_t) - \bar{f}(\omega_t, \theta_t) \rangle \\ &+ 2\beta \langle \Delta_t, \bar{f}(\omega_t, \theta_t) \rangle + 2 \langle \Delta_t, \omega_t^* - \omega_{t+1}^* \rangle \\ &\leq \|\Delta_t\|^2 + 2\beta^2 \|f(O_t, \omega_t)\|^2 + 2\|\omega_t^* - \omega_{t+1}^* \|^2 \\ &+ 2\beta \langle \Delta_t, f(O_t, \omega_t) - \bar{f}(\omega_t, \theta_t) \rangle \\ &+ 2\beta \langle \Delta_t, \bar{f}(\omega_t, \theta_t) \rangle + 2 \langle \Delta_t, \omega_t^* - \omega_{t+1}^* \rangle, \end{split}$$

where f and  $\overline{f}$  are defined in Eq. (20).

Taking expectation up to  $s_{t+1}$ , we have

$$\mathbb{E}\|\Delta_{t+1}\|^{2} \leq \mathbb{E}\|\Delta_{t}\|^{2} + \underbrace{2\beta^{2}\mathbb{E}\|\boldsymbol{f}(O_{t},\boldsymbol{\omega}_{t})\|^{2}}_{I_{1}} + \underbrace{2\mathbb{E}\|\boldsymbol{\omega}_{t}^{*}-\boldsymbol{\omega}_{t+1}^{*}\|^{2}}_{I_{2}} + \underbrace{2\beta\mathbb{E}\langle\Delta_{t},\boldsymbol{f}(\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t})\rangle}_{I_{3}} + \underbrace{2\beta\mathbb{E}\langle\Delta_{t},\boldsymbol{f}(O_{t},\boldsymbol{\omega}_{t})-\boldsymbol{\bar{f}}(\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t})\rangle}_{I_{4}} + \underbrace{2\mathbb{E}\langle\Delta_{t},\boldsymbol{\omega}_{t}^{*}-\boldsymbol{\omega}_{t+1}^{*}\rangle}_{I_{5}}.$$

$$(27)$$

In the sequel, we will tackle  $I_1, I_2, I_3, I_4, I_5$  respectively.

For term  $I_1$ , since  $\|\boldsymbol{f}(O_t, \boldsymbol{\omega}_t)\| \leq \bar{\delta}$ , we have

$$I_1 = 2\beta^2 \mathbb{E} \| \boldsymbol{f}(O_t, \boldsymbol{\omega}_t) \|^2 \le 2\beta^2 \bar{\delta}^2$$

For term  $I_2$ , from Lemma B.6, it can be shown that

$$I_2 = 2\mathbb{E} \|\boldsymbol{\omega}_t^* - \boldsymbol{\omega}_{t+1}^*\|^2 \le 2L_c^2 \mathbb{E} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|^2 \le 2\alpha^2 \bar{\delta}^2 B^2 L_c^2.$$

For term  $I_3$ , we have

$$\begin{split} \langle \Delta_t, \boldsymbol{f}(\boldsymbol{\omega}_t, \boldsymbol{\theta}_t) \rangle = & \langle \Delta_t, \boldsymbol{f}(\boldsymbol{\omega}_t, \boldsymbol{\theta}_t) - \boldsymbol{f}(\boldsymbol{\omega}_t^*, \boldsymbol{\theta}_t) \rangle \\ = & \langle \Delta_t, \mathbb{E}[(\gamma \boldsymbol{\phi}(s') - \boldsymbol{\phi}(s))^\top (\boldsymbol{\omega}_t - \boldsymbol{\omega}_t^*) \boldsymbol{\phi}(s)] \rangle \\ = & \Delta_t^\top \mathbb{E}[\boldsymbol{\phi}(s)(\gamma \boldsymbol{\phi}(s') - \boldsymbol{\phi}(s)] \Delta_t \\ = & - \Delta_t^\top \boldsymbol{A}_{\boldsymbol{\theta}} \Delta_t \\ \leq & - \lambda \|\Delta_t\|^2. \end{split}$$

It follows that

$$I_3 \le -2\lambda\beta \mathbb{E} \|\Delta_t\|^2.$$

For term  $I_4$ , according to Lemma C.1, it holds that

$$I_4 = 2\beta \mathbb{E} \left[ \Lambda(O_t, \boldsymbol{\omega}_t, \boldsymbol{\theta}_t) \right] \le 2\beta M_1 \frac{1}{\sqrt{T}}.$$

For term  $I_5$ , we will instead give an implicit upper bound. It can be shown that

$$\begin{split} \mathbb{E}\langle\Delta_{t}, \boldsymbol{\omega}_{t}^{*} - \boldsymbol{\omega}_{t+1}^{*}\rangle &= \mathbb{E}\langle\Delta_{t}, \boldsymbol{\omega}_{t}^{*} - \boldsymbol{\omega}_{t+1}^{*} + (\nabla\boldsymbol{\omega}_{t}^{*})^{\top}(\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t})\rangle + \mathbb{E}\langle\Delta_{t}, -(\nabla\boldsymbol{\omega}_{t}^{*})^{\top}(\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t})\rangle \\ &\leq \frac{1}{2} \mathbb{E}\|\Delta_{t}\|\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t}\|^{2} + \alpha \mathbb{E}\langle\Delta_{t}, -(\nabla\boldsymbol{\omega}_{t}^{*})^{\top}\boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t})\rangle \\ &\leq \alpha^{2}\overline{\delta}^{2}B^{2}L_{s}\overline{\boldsymbol{\omega}} + \alpha \mathbb{E}\langle\Delta_{t}, -(\nabla\boldsymbol{\omega}_{t}^{*})^{\top}\boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t})\rangle \\ &= \alpha^{2}\overline{\delta}^{2}B^{2}L_{s}\overline{\boldsymbol{\omega}} + \alpha \mathbb{E}\langle\Delta_{t}, -(\nabla\boldsymbol{\omega}_{t}^{*})^{\top}(\boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t}) - \boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t}))\rangle \\ &+ \alpha \mathbb{E}\langle\Delta_{t}, -(\nabla\boldsymbol{\omega}_{t}^{*})^{\top}\boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t})\rangle \\ &\leq \alpha^{2}\overline{\delta}^{2}B^{2}L_{s}\overline{\boldsymbol{\omega}} + 2\alpha BL_{c}\mathbb{E}\|\Delta_{t}\|^{2} + \alpha \mathbb{E}\langle\Delta_{t}, -(\nabla\boldsymbol{\omega}_{t}^{*})^{\top}\boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t})\rangle \\ &= \alpha^{2}\overline{\delta}^{2}B^{2}L_{s}\overline{\boldsymbol{\omega}} + 2\alpha BL_{c}\mathbb{E}\|\Delta_{t}\|^{2} + \alpha \mathbb{E}\langle\Delta_{t}, (\nabla\boldsymbol{\omega}_{t}^{*})^{\top}(\bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t}) - \boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t}))\rangle \\ &+ \alpha \mathbb{E}\langle\Delta_{t}, (\nabla\boldsymbol{\omega}_{t}^{*})^{\top}(\bar{\boldsymbol{g}}(\boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t}) - \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t}))\rangle + \alpha \mathbb{E}\langle\Delta_{t}, -(\nabla\boldsymbol{\omega}_{t}^{*})^{\top}\bar{\boldsymbol{g}}(\boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t})\rangle, \\ & J_{1} \end{pmatrix}$$

where (1) follows from the smoothness of the optimal critic shown in Lemma B.9. We will analyze  $J_1$ ,  $J_2$ , and  $J_3$  individually.

For term  $J_1$ , from the Markovian noise analysis in Lemma C.2, we have

$$J_1 = \mathbb{E}\big[\Gamma(\widehat{O}_t, \boldsymbol{\omega}_t, \boldsymbol{\theta}_t)\big] \le M_2 \frac{1}{\sqrt{T}}.$$

For term  $J_2$ , from the policy gradient theorem in Eq. (5), we obtain

$$\bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}^{*},\boldsymbol{\theta}_{t}) - \bar{\boldsymbol{g}}(\boldsymbol{\omega}_{t}^{*},\boldsymbol{\theta}_{t}) = \mathbb{E}_{(s,a,s')\sim(\nu_{\boldsymbol{\theta}_{t}},\pi_{\boldsymbol{\theta}_{t}},P)}[(r(s,a) + \gamma V_{\boldsymbol{\theta}_{t}}(s') - V_{\boldsymbol{\theta}_{t}}(s))\nabla\log\pi_{\boldsymbol{\theta}_{t}}(a\,|\,s)] = (1-\gamma)\nabla J(\boldsymbol{\theta}_{t}).$$
(28)

It follows that

$$J_2 = \mathbb{E} \langle \Delta_t, -(\nabla \boldsymbol{\omega}_t^*)^\top (1-\gamma) \nabla J(\boldsymbol{\theta}_t) \rangle \leq (1-\gamma) L_c \mathbb{E} \| \Delta_t \| \| \nabla J(\boldsymbol{\theta}_t) \|.$$

For term  $J_3$ , we first show that

$$\bar{\boldsymbol{g}}(\boldsymbol{\omega}_{t}^{*},\boldsymbol{\theta}_{t}) \leq \sqrt{\mathbb{E}_{(s,a,s')\sim(\boldsymbol{\nu}_{\boldsymbol{\theta}_{t}},\pi_{\boldsymbol{\theta}_{t}},P)} \|\boldsymbol{g}(O_{t},\boldsymbol{\omega}_{t}^{*},\boldsymbol{\theta}_{t})\|^{2}} \\
\leq \sqrt{\mathbb{E}\left[B^{2}((\gamma\boldsymbol{\phi}(s')^{\top}\boldsymbol{\omega}_{t}^{*}-\gamma V_{\boldsymbol{\theta}_{t}}(s'))-(\boldsymbol{\phi}(s)^{\top}\boldsymbol{\omega}_{t}^{*}-V_{\boldsymbol{\theta}}(s)))^{2}\right]} \\
\leq \sqrt{\mathbb{E}\left[2B^{2}\left(\gamma^{2}(\boldsymbol{\phi}(s')^{\top}\boldsymbol{\omega}_{t}^{*}-V_{\boldsymbol{\theta}_{t}}(s'))^{2}+(\boldsymbol{\phi}(s)^{\top}\boldsymbol{\omega}_{t}^{*}-V_{\boldsymbol{\theta}_{t}}(s))^{2}\right)\right]} \\
\leq 2B\sqrt{\mathbb{E}\left[(\boldsymbol{\phi}(s)^{\top}\boldsymbol{\omega}_{t}^{*}-V_{\boldsymbol{\theta}_{t}}(s))^{2}\right]} \\
= 2B\epsilon_{\mathrm{app}}.$$
(29)

Then we have

 $J_3 \leq 4\bar{\omega}BL_c\epsilon_{\mathrm{app}}.$ 

Combining  $J_1, J_2$ , and  $J_3$ , we get

$$I_5 \leq 2\alpha^2 \bar{\delta}^2 B^2 L_s \bar{\omega} + 4\alpha B L_c \mathbb{E} \|\Delta_t\|^2 + 8\alpha \bar{\omega} B L_c \epsilon_{\mathrm{app}} + 2\alpha (1-\gamma) L_c \mathbb{E} \|\Delta_t\| \|\nabla J(\boldsymbol{\theta}_t)\| + 2\alpha M_2 \frac{1}{\sqrt{T}}.$$

Plugging  $I_1, I_2, I_3, I_4$  and  $I_5$  into Eq. (27), we obtain

$$\begin{split} \mathbb{E}\|\Delta_{t+1}\|^2 &\leq \mathbb{E}\|\Delta_t\|^2 + 2\beta^2 \bar{\delta}^2 + 2\alpha^2 \bar{\delta}^2 B^2 L_c^2 + 2\beta M_1 \frac{1}{\sqrt{T}} - 2\lambda\beta \mathbb{E}\|\Delta_t\|^2 + 2\alpha^2 \bar{\delta}^2 B^2 L_s \bar{\omega} \\ &\quad + 4\alpha B L_c \mathbb{E}\|\Delta_t\|^2 + 8\alpha \bar{\omega} B L_c \epsilon_{\mathrm{app}} + 2\alpha (1-\gamma) L_c \mathbb{E}\|\Delta_t\| \|\nabla J(\boldsymbol{\theta}_t)\| + 2\alpha M_2 \frac{1}{\sqrt{T}} \\ &\stackrel{(1)}{\leq} (1-\lambda\beta) \mathbb{E}\|\Delta_t\|^2 + 2\alpha (1-\gamma) L_c \mathbb{E}\|\Delta_t\| \|\nabla J(\boldsymbol{\theta}_t)\| \\ &\quad + (2\bar{\delta}^2 + 2c^2 \bar{\delta}^2 B^2 L_c^2 + 2M_1 + 2c^2 \bar{\delta}^2 B^2 L_s \bar{\omega} + 2cM_2) \frac{1}{T} + 8\alpha \bar{\omega} B L_c \epsilon_{\mathrm{app}}, \end{split}$$

where (1) holds as the step size ratio c is chosen to satisfy  $4\alpha BL_c \leq \lambda\beta$ .

Rearranging the above inequality, we obtain

$$\begin{aligned} \mathbb{E}\|\Delta_t\|^2 &\leq \frac{1}{\lambda\beta} (\mathbb{E}\|\Delta_t\|^2 - \mathbb{E}\|\Delta_{t+1}\|^2) + \frac{2c(1-\gamma)}{\lambda} L_c \mathbb{E}\|\Delta_t\| \|\nabla J(\boldsymbol{\theta}_t)\| \\ &+ \lambda^{-1} (2\bar{\delta}^2 + 2c^2\bar{\delta}^2 B^2 L_c^2 + 2M_1 + 2c^2\bar{\delta}^2 B^2 L_s \bar{\omega} + 2cM_2) \frac{1}{\sqrt{T}} + 8c\bar{\omega}BL_c \epsilon_{\mathrm{app}}. \end{aligned}$$

By leveraging the  $\mathcal{O}(\cdot)$  notation, we can further summarise our implicit analysis for the critic as

$$\mathbb{E}\|\Delta_t\|^2 \leq \frac{1}{\lambda\beta} (\mathbb{E}\|\Delta_t\|^2 - \mathbb{E}\|\Delta_{t+1}\|^2) + \frac{2c(1-\gamma)}{\lambda} L_c \mathbb{E}\|\Delta_t\| \|\nabla J(\boldsymbol{\theta}_t)\| + \mathcal{O}(\frac{\log^2 T}{\sqrt{T}}) + \mathcal{O}(\epsilon_{\mathrm{app}}),$$

where the term  $\log^2 T$  arises from the presence of  $\tau_{\min}^2$  in  $M_1$  and  $M_2$ . Therefore, we finish the proof of Theorem D.1.

## D.2. An implicit bound for actor error

**Theorem D.2.** Choose  $\alpha_t = c/\sqrt{T}$ ,  $\beta_t = 1/\sqrt{T}$ , for any  $\tau_{mix} \leq t < T$ , we have

$$(1-\gamma)\mathbb{E}\|\nabla J(\boldsymbol{\theta}_t)\|^2 \le \frac{1}{\alpha} (\mathbb{E}[J(\boldsymbol{\theta}_{t+1}) - J(\boldsymbol{\theta}_t)]) + 2B\mathbb{E}\|\nabla J(\boldsymbol{\theta}_t)\|\|\Delta_t\| + \mathcal{O}(\frac{\log^2 T}{\sqrt{T}}) + \mathcal{O}(\epsilon_{\mathrm{app}}).$$
(30)

Proof. From the update rule of actor in Line 8 of Algorithm 1 and Lemma B.8, we have

$$\begin{split} J(\boldsymbol{\theta}_{t+1}) &\geq J(\boldsymbol{\theta}_{t}) + \langle \nabla J(\boldsymbol{\theta}_{t}), \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t} \rangle - \frac{L_{g}}{2} \|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}_{t+1}\|^{2} \\ &= J(\boldsymbol{\theta}_{t}) + \alpha \langle \nabla J(\boldsymbol{\theta}_{t}), \boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t}) \rangle - \frac{L_{g}}{2} \alpha^{2} \|\boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t})\|^{2} \\ &= J(\boldsymbol{\theta}_{t}) + \alpha \langle \nabla J(\boldsymbol{\theta}_{t}), \boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t}) - \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t}) \rangle + \alpha \langle \nabla J(\boldsymbol{\theta}_{t}), \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t}) \rangle - \frac{L_{g}}{2} \alpha^{2} \|\boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t})\|^{2} \\ &= J(\boldsymbol{\theta}_{t}) + \alpha \langle \nabla J(\boldsymbol{\theta}_{t}), \boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t}) - \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t}) \rangle + \alpha \langle \nabla J(\boldsymbol{\theta}_{t}), \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t}) - \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t}) \rangle \\ &+ \alpha \langle \nabla J(\boldsymbol{\theta}_{t}), \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t}) - \bar{\boldsymbol{g}}(\boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t}) \rangle + \alpha \langle \nabla J(\boldsymbol{\theta}_{t}), \bar{\boldsymbol{g}}(\boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t}) \rangle - \frac{L_{g}}{2} \alpha^{2} \|\boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t})\|^{2} \\ & \stackrel{(1)}{=} J(\boldsymbol{\theta}_{t}) + \alpha \langle \nabla J(\boldsymbol{\theta}_{t}), \boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t}) - \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t}) \rangle + \alpha \langle \nabla J(\boldsymbol{\theta}_{t}), \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t}) - \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t}) \rangle \\ &+ \alpha (1 - \gamma) \| \nabla J(\boldsymbol{\theta}_{t}) \|^{2} + \alpha \langle \nabla J(\boldsymbol{\theta}_{t}), \bar{\boldsymbol{g}}(\boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t}) \rangle - \frac{L_{g}}{2} \alpha^{2} \|\boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t})\|^{2}, \end{split}$$

where (1) follows from Eq. (28).

Rearranging the above inequality and taking expectation, we have

$$(1-\gamma)\mathbb{E}\|\nabla J(\boldsymbol{\theta}_{t})\|^{2} \leq \frac{1}{\alpha} (\mathbb{E}[J(\boldsymbol{\theta}_{t+1}) - J(\boldsymbol{\theta}_{t})]) + \underbrace{\frac{\alpha L_{g}}{2} \mathbb{E}\|\boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t})\|^{2}}_{I_{1}} - \underbrace{\mathbb{E}\langle\nabla J(\boldsymbol{\theta}_{t}), \bar{\boldsymbol{g}}(\boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t})\rangle}_{I_{2}}_{I_{2}} + \underbrace{\mathbb{E}\langle\nabla J(\boldsymbol{\theta}_{t}), \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t}) - \boldsymbol{h}(\widehat{O}_{t}, \boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t})\rangle}_{I_{3}} + \underbrace{\mathbb{E}\langle\nabla J(\boldsymbol{\theta}_{t}), \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}^{*}, \boldsymbol{\theta}_{t}) - \bar{\boldsymbol{h}}(\boldsymbol{\omega}_{t}, \boldsymbol{\theta}_{t})\rangle}_{I_{4}}.$$

$$(31)$$

In the sequel, we will analyze  $I_1, I_2, I_3, I_4$  one by one. For term  $I_1$ , since  $h(\widehat{O}_t, \omega_t, \theta_t) \leq \overline{\delta}B$ , we have

$$I_1 = \frac{\alpha L_g}{2} \mathbb{E} \| \boldsymbol{h}(\widehat{O}_t, \boldsymbol{\omega}_t, \boldsymbol{\theta}_t) \|^2 \le \frac{\alpha \delta^2 B^2 L_g}{2}.$$

For term  $I_2$ , from Eq. (29), we have

$$I_2 = \mathbb{E} \langle \nabla J(\boldsymbol{\theta}_t), \bar{\boldsymbol{g}}(\boldsymbol{\omega}_t^*, \boldsymbol{\theta}_t) \rangle \leq 2BL_J \epsilon_{\mathrm{app}}.$$

For term  $I_3$ , from Lemma C.3, we obtain

$$I_3 = \mathbb{E}\big[\Xi(\widehat{O}_t, \boldsymbol{\omega}_t, \boldsymbol{\theta}_t)\big] \le M_3 \frac{1}{\sqrt{T}}.$$

For term  $I_4$ , it holds that

$$I_4 = \mathbb{E}\langle \nabla J(\boldsymbol{\theta}_t), \bar{\boldsymbol{h}}(\boldsymbol{\omega}_t^*, \boldsymbol{\theta}_t) - \bar{\boldsymbol{h}}(\boldsymbol{\omega}_t, \boldsymbol{\theta}_t) \rangle \leq 2B\mathbb{E} \|\nabla J(\boldsymbol{\theta}_t)\| \|\Delta_t\|$$

Plugging  $I_1, I_2, I_3$  and  $I_4$  into Eq. (31), we have

$$(1-\gamma)\mathbb{E}\|\nabla J(\boldsymbol{\theta}_t)\|^2 \leq \frac{1}{\alpha} (\mathbb{E}[J(\boldsymbol{\theta}_{t+1}) - J(\boldsymbol{\theta}_t)]) + \frac{\alpha \bar{\delta}^2 B^2 L_g}{2} + 2BL_J \epsilon_{\mathrm{app}} + M_3 \frac{1}{\sqrt{T}} + 2B\mathbb{E}\|\nabla J(\boldsymbol{\theta}_t)\| \|\Delta_t\|.$$

By leveraging the  $\mathcal{O}(\cdot)$  notation, we can further summarise our implicit analysis for the actor as

$$(1-\gamma)\mathbb{E}\|\nabla J(\boldsymbol{\theta}_t)\|^2 \leq \frac{1}{\alpha} (\mathbb{E}[J(\boldsymbol{\theta}_{t+1}) - J(\boldsymbol{\theta}_t)]) + 2B\mathbb{E}\|\nabla J(\boldsymbol{\theta}_t)\|\|\Delta_t\| + \mathcal{O}(\frac{\log^2 T}{\sqrt{T}}) + \mathcal{O}(\epsilon_{\mathrm{app}}),$$

where the term  $\log^2 T$  arises from the presence of  $\tau_{\text{mix}}^2$  in  $M_3$ . Therefore, we complete the proof of Theorem D.2.

## D.3. A novel Lyapunov analysis

**Theorem D.3.** Choose  $\alpha_t = c/\sqrt{T}$ ,  $\beta_t = 1/\sqrt{T}$ , for any  $T \ge 2\tau_{\text{mix}}$ , we have

$$\frac{1}{T - \tau_{\min}} \sum_{t=\tau_{\min}}^{T-1} \mathbb{E} \left\| \Delta_t \right\|^2 = \mathcal{O}\left(\frac{\log^2 T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\text{app}}),$$

$$\frac{1}{T - \tau_{\min}} \sum_{t=\tau_{\min}}^{T-1} \mathbb{E} \left\| \nabla J(\boldsymbol{\theta}_t) \right\|^2 = \mathcal{O}\left(\frac{\log^2 T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\text{app}}).$$
(32)

Proof. Define

$$\mathbb{L}_t = \frac{2B}{1-\gamma} \mathbb{E} \|\Delta_t\|^2 + \frac{1-\gamma}{2B} \mathbb{E} \|\nabla J(\boldsymbol{\theta}_t)\|^2,$$

the sum of Eq. (26) and Eq. (30) yields

$$\begin{split} \mathbb{L}_{t} &\leq \frac{2B}{\lambda\beta(1-\gamma)} (\mathbb{E}\|\Delta_{t}\|^{2} - \mathbb{E}\|\Delta_{t+1}\|^{2}) + \frac{4L_{c}Bc}{\lambda} \mathbb{E}\|\Delta_{t}\|\|\nabla J(\boldsymbol{\theta}_{t})\| \\ &+ \frac{1}{2B\alpha} (\mathbb{E}[J(\boldsymbol{\theta}_{t+1}) - J(\boldsymbol{\theta}_{t})]) + \mathbb{E}\|\Delta_{t}\|\|\nabla J(\boldsymbol{\theta}_{t})\| + \mathcal{O}(\frac{\log^{2}T}{\sqrt{T}}) + \mathcal{O}(\epsilon_{\mathrm{app}}) \\ &\leq \left(\frac{2L_{c}Bc}{\lambda} + \frac{1}{2}\right) \mathbb{L}_{t} + \frac{2B}{\lambda\beta(1-\gamma)} (\mathbb{E}\|\Delta_{t}\|^{2} - \mathbb{E}\|\Delta_{t+1}\|^{2}) \\ &+ \frac{1}{2B\alpha} (\mathbb{E}[J(\boldsymbol{\theta}_{t+1}) - J(\boldsymbol{\theta}_{t})]) + \mathcal{O}\left(\frac{\log^{2}T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\mathrm{app}}), \end{split}$$

where the last inequality follows from  $\mathbb{E} \|\Delta_t\| \|\nabla J(\boldsymbol{\theta}_t)\| \leq 1/2\mathbb{L}_t$ . Since  $\mathcal{L} := 1/(T - \tau_{\min}) \sum_{t=\tau_{\min}}^{T-1} \mathbb{L}_t$ , it can be shown that

$$\begin{aligned} \mathcal{L} &\leq \left(\frac{2L_cBc}{\lambda} + \frac{1}{2}\right) \mathcal{L} + \frac{2B}{\lambda\beta(1-\gamma)(T-\tau_{\min})} \sum_{t=\tau_{\min}}^{T-1} (\mathbb{E}\|\Delta_t\|^2 - \mathbb{E}\|\Delta_{t+1}\|^2) \\ &+ \frac{1}{2B\alpha(T-\tau_{\min})} \sum_{t=\tau_{\min}}^{T-1} (\mathbb{E}[J(\boldsymbol{\theta}_{t+1}) - J(\boldsymbol{\theta}_t)]) + \mathcal{O}\left(\frac{\log^2 T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\mathrm{app}}) \\ &\leq \left(\frac{2L_cBc}{\lambda} + \frac{1}{2}\right) \mathcal{L} + \frac{8B\bar{\omega}^2}{\lambda\beta(1-\gamma)(T-\tau_{\min})} + \frac{\bar{r}}{B\alpha(T-\tau_{\min})} + \mathcal{O}\left(\frac{\log^2 T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\mathrm{app}}). \end{aligned}$$

Choose  $T \geq 2\tau_{\rm mix},$  we have  $T-\tau_{\rm mix} \geq 1/2T,$  which implies

$$\begin{aligned} \mathcal{L} &\leq \left(\frac{2L_cBc}{\lambda} + \frac{1}{2}\right) \mathcal{L} + \left(\frac{16B\bar{\omega}^2}{\lambda(1-\gamma)} + \frac{2\bar{r}}{Bc}\right) \frac{1}{\sqrt{T}} + \mathcal{O}\left(\frac{\log^2 T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\mathrm{app}}) \\ &= \left(\frac{2L_cBc}{\lambda} + \frac{1}{2}\right) \mathcal{L} + \mathcal{O}\left(\frac{\log^2 T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\mathrm{app}}). \end{aligned}$$

Overall, we have

$$\mathcal{L} \le \left(\frac{2L_c Bc}{\lambda} + \frac{1}{2}\right) \mathcal{L} + \mathcal{O}\left(\frac{\log^2 T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\rm app}).$$
(33)

To make  $\mathcal{L}$  convergence, we need  $2L_cBc/\lambda + 1/2 < 1$ , which can be achieved by choosing

$$c < \frac{\lambda}{4BL_c}.$$
(34)

It follows that

$$\mathcal{L} = \mathcal{O}\left(\frac{\log^2 T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\mathrm{app}}),$$

which implies

$$\frac{1}{T - \tau_{\min}} \sum_{t=\tau_{\min}}^{T-1} \mathbb{E} \|\Delta_t\|^2 = \mathcal{O}\left(\frac{\log^2 T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\text{app}}),$$
$$\frac{1}{T - \tau_{\min}} \sum_{t=\tau_{\min}}^{T-1} \mathbb{E} \|\nabla J(\boldsymbol{\theta}_t)\|^2 = \mathcal{O}\left(\frac{\log^2 T}{\sqrt{T}}\right) + \mathcal{O}(\epsilon_{\text{app}}).$$

Therefore, we complete our proof.

## **E. Proof of Propositions**

## **Proof of Proposition 3.1**.

*Proof.* We show that  $\nu_{\theta}$  is the stationary distribution of the Markov chain induced by  $\widehat{\mathcal{P}}$  by showing that  $\nu_{\theta}$  is a fixed point of the operator  $\widehat{\mathcal{P}}$ , i.e.,

$$\widehat{\mathcal{P}}\nu_{\theta} = \nu_{\theta}.$$

Define operator  $\mathcal{P}^t$  by iterative application of the operator  $\mathcal{P}$ :

$$(\mathcal{P}^t f)(s') = \int_{\mathcal{S}} \int_{\mathcal{A}} \pi_{\theta}(a \mid s) P(s' \mid s, a) \mathcal{P}^{t-1} f(s) \, dads.$$

From the definition of the operator  $\mathcal{P}^t$ , we can rewrite  $\nu_{\theta}$  in Eq. (4) as

$$\nu_{\boldsymbol{\theta}}(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathcal{P}^t \eta(s)$$

Then we have

$$\widehat{\mathcal{P}}\nu_{\theta}(s) = (1 - \gamma)\eta(s) + \gamma \mathcal{P}\nu_{\theta}(s).$$
(35)

For term  $\mathcal{P}\nu_{\theta}(s)$ , it holds that

$$\mathcal{P}\nu_{\theta}(s) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^{t} \mathcal{P}^{t+1} \eta(s)$$
$$= (1-\gamma) \sum_{k=1}^{\infty} \gamma^{k-1} \mathcal{P}^{k} \eta(s)$$
$$= \frac{1-\gamma}{\gamma} \sum_{k=1}^{\infty} \gamma^{k} \mathcal{P}^{k} \eta(s)$$
$$= \frac{1-\gamma}{\gamma} (\frac{\nu_{\theta}(s)}{1-\gamma} - \eta(s))$$
$$= \frac{\nu_{\theta}(s)}{\gamma} - \frac{1-\gamma}{\gamma} \eta(s).$$

Plugging the above result to Eq. (35), we obtain

$$\begin{aligned} \widehat{\mathcal{P}}\nu_{\boldsymbol{\theta}}(s) &= (1-\gamma)\eta(s) + \gamma(\frac{\nu_{\boldsymbol{\theta}}(s)}{\gamma} - \frac{1-\gamma}{\gamma}\eta(s)) \\ &= (1-\gamma)\eta(s) + \nu_{\boldsymbol{\theta}}(s) - (1-\gamma)\eta(s) \\ &= \nu_{\boldsymbol{\theta}}(s). \end{aligned}$$

Suppose  $\widehat{\mathcal{P}}$  has two fixed points f and g, then we have

$$\widehat{\mathcal{P}}f = f, \quad \widehat{\mathcal{P}}g = g.$$

Recall that for two probability distributions  $\mu$  and  $\nu$  on S, the total variation distance is defined as

$$d_{\rm TV}(\mu,\nu) = \frac{1}{2} \int_{\mathcal{S}} |\mu(s) - \nu(s)| \, ds.$$

Since we have

 $\widehat{\mathcal{P}}f - \widehat{\mathcal{P}}g = \gamma(\mathcal{P}f - \mathcal{P}g),$ 

it follows that

$$d_{\mathrm{TV}}(\mathcal{P}f, \mathcal{P}g) = \gamma d_{TV}(\mathcal{P}f, \mathcal{P}g).$$

From Eq. (10), we know that  $\mathcal{P}$  is also a Markov kernel which does not increase the total variation distance. Therefore, it can be shown that

$$d_{\mathrm{TV}}(f,g) = d_{\mathrm{TV}}(\mathcal{P}f,\mathcal{P}g) = \gamma d_{TV}(\mathcal{P}f,\mathcal{P}g) \le \gamma d_{\mathrm{TV}}(f,g).$$

Therefore, we get

 $d_{\rm TV}(f,g) = 0,$ 

which means f and g are same distributions. Hence we finish our proof.

#### **Proof of Proposition 3.2**.

*Proof.* Recall that for two probability distributions  $\mu$  and  $\nu$  on S, the total variation distance is defined as

$$d_{\rm TV}(\mu,\nu) = \frac{1}{2} \int_{\mathcal{S}} |\mu(s) - \nu(s)| \, ds.$$

For two distribution f and g, we have

$$\widehat{\mathcal{P}}f - \widehat{\mathcal{P}}g = \gamma(\mathcal{P}f - \mathcal{P}g),$$

where the operator  $\widehat{\mathcal{P}}$  and  $\mathcal{P}$  are defined in the proof of Proposition 3.1.

It follows that

$$d_{\mathrm{TV}}(\widehat{\mathcal{P}}f,\widehat{\mathcal{P}}g) = \gamma d_{TV}(\mathcal{P}f,\mathcal{P}g).$$

From Eq. (10), we know that  $\mathcal{P}$  is also a Markov kernel which does not increase the total variation distance. Therefore, it can be shown that

$$d_{\mathrm{TV}}(\widehat{\mathcal{P}}f,\widehat{\mathcal{P}}g) = \gamma d_{TV}(\mathcal{P}f,\mathcal{P}g) \le \gamma d_{\mathrm{TV}}(f,g).$$

As shown in Proposition 3.1,  $\nu_{\theta}$  is the stationary distribution of the Markov chain induced by  $\hat{\mathcal{P}}$ . For any initial distribution f, we have

$$d_{\mathrm{TV}}(\widehat{\mathcal{P}}^t f, \nu_{\theta}) = d_{\mathrm{TV}}(\widehat{\mathcal{P}}^t f, \widehat{\mathcal{P}}^t \nu_{\theta})$$
  
$$\leq \gamma^t d_{\mathrm{TV}}(f, \nu_{\theta})$$
  
$$\leq \gamma^t.$$

Thus, it completes the proof.

#### **Proof of Proposition 4.4**.

Proof. From the definition of the total variation distance, we have

$$d_{\mathrm{TV}}(\pi_{\boldsymbol{\theta}_{1}}(\cdot \mid s) - \pi_{\boldsymbol{\theta}_{2}}(\cdot \mid s)) = \frac{1}{2} \int_{\mathcal{A}} |\pi_{\boldsymbol{\theta}_{1}}(a \mid s) - \pi_{\boldsymbol{\theta}_{2}}(a \mid s)| \, da$$
$$= \frac{1}{2} \int_{\bar{\mathcal{A}}} |\pi_{\boldsymbol{\theta}_{1}}(a \mid s) - \pi_{\boldsymbol{\theta}_{2}}(a \mid s)| \, da$$
$$\leq \frac{1}{2} \int_{\bar{\mathcal{A}}} L \|\boldsymbol{\theta}_{1} - \boldsymbol{\theta}_{2}\| \, da$$
$$\leq \frac{1}{2} \bar{\mathcal{A}} L \|\boldsymbol{\theta}_{1} - \boldsymbol{\theta}_{2}\|,$$

where  $\bar{\mathcal{A}}$  is the bounded support of  $\pi_{\theta}(a \mid s)$  which satisfies  $\int_{\bar{\mathcal{A}}} da = \bar{A}$ . Define  $L_{\pi} := 1/2\bar{A}L$ , which completes the proof.

# F. Proof of Preliminary Lemmas

## Proof of Lemma B.1.

*Proof.* This is a minor adjustment to the proof of Lemma 3 in Zou et al. (2019), extending it to continuous settings. For any  $\theta_1$  and  $\theta_2$ , define the transition densities respectively as follows:

$$P_{\boldsymbol{\theta}_i}(s \,|\, ds') = \int_{\mathcal{A}} P(ds' \,|\, s, a) \pi_{\boldsymbol{\theta}_i}(a \,|\, s), \quad i = 1, 2$$

Following from Theorem 3.1 in (Mitrophanov, 2005), we obtain

$$d_{\mathrm{TV}}(\mu_{\boldsymbol{\theta}_1}, \mu_{\boldsymbol{\theta}_2}) \leq (\lceil \log_{\rho} m^{-1} \rceil + \frac{1}{1-\rho}) \| P_{\boldsymbol{\theta}_1} - P_{\boldsymbol{\theta}_2} \|_{\mathrm{op}},$$

where  $\|\cdot\|_{\text{op}}$  is the operator norm defined in (Mitrophanov, 2005):  $\|A\|_{\text{op}} := \sup_{\|q\|_{\text{TV}}=1} \|qA\|_{\text{TV}}$ , and  $\|\cdot\|_{\text{TV}}$  denotes the total-variation norm. Then we have

$$\begin{split} \|P_{\theta_{1}} - P_{\theta_{2}}\|_{\mathrm{op}} &= \sup_{\|q\|_{\mathrm{TV}}=1} \|\int_{\mathcal{S}} q(ds)(P_{\theta_{1}} - P_{\theta_{2}})(s \mid \cdot)\|_{\mathrm{TV}} \\ &= \sup_{\|q\|_{\mathrm{TV}}=1} \int_{\mathcal{S}} \int_{\mathcal{S}} |\int_{\mathcal{S}} q(ds)(P_{\theta_{1}} - P_{\theta_{2}})(s \mid ds')| \\ &\leq \sup_{\|q\|_{\mathrm{TV}}=1} \int_{\mathcal{S}} \int_{\mathcal{S}} |q(ds)| \left| (P_{\theta_{1}} - P_{\theta_{2}})(s \mid ds') \right| \\ &= \sup_{\|q\|_{\mathrm{TV}}=1} \int_{\mathcal{S}} \int_{\mathcal{S}} |q(ds)| \left| \int_{\mathcal{A}} P(ds' \mid s, a)(\pi_{\theta_{1}}(da \mid s) - \pi_{\theta_{2}}(da \mid s)) \right| \\ &= \sup_{\|q\|_{\mathrm{TV}}=1} \int_{\mathcal{S}} \int_{\mathcal{S}} |q(ds)| \int_{\mathcal{A}} P(ds' \mid s, a)|(\pi_{\theta_{1}}(da \mid s) - \pi_{\theta_{2}}(da \mid s))| \\ &= \sup_{\|q\|_{\mathrm{TV}}=1} \int_{\mathcal{S}} \int_{\mathcal{S}} |q(ds)| \int_{\mathcal{A}} |(\pi_{\theta_{1}}(da \mid s) - \pi_{\theta_{2}}(da \mid s))| \\ &\leq 2L_{\pi} \|\theta_{1} - \theta_{2}\|. \end{split}$$

Therefore, we have

$$d_{TV}(\mu_{\boldsymbol{\theta}_1}, \mu_{\boldsymbol{\theta}_2}) \leq 2L_{\pi}(\lceil \log_{\rho} m^{-1} \rceil + \frac{1}{1-\rho}) \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|.$$

For the second inequality, we have

$$\begin{split} d_{TV}(\mu_{\theta_{1}} \otimes \pi_{\theta_{1}}, \mu_{\theta_{2}} \otimes \pi_{\theta_{2}}) &= \frac{1}{2} \int_{\mathcal{S}} \int_{\mathcal{A}} |\mu_{\theta_{1}}(ds)\pi_{\theta_{1}}(da \,|\, s) - \mu_{\theta_{2}}(ds)\pi_{\theta_{2}}(da \,|\, s)| \\ &\leq \frac{1}{2} \int_{\mathcal{S}} \int_{\mathcal{A}} |\mu_{\theta_{1}}(ds)(\pi_{\theta_{1}}(da \,|\, s) - \pi_{\theta_{2}}(da \,|\, s))| \\ &\quad + \frac{1}{2} \int_{\mathcal{S}} \int_{\mathcal{A}} |(\mu_{\theta_{1}}(ds) - \mu_{\theta_{2}}(ds))\pi_{\theta_{2}}(da \,|\, s))| \\ &= d_{TV}(\pi_{\theta_{1}}, \pi_{\theta_{2}}) + d_{TV}(\mu_{\theta_{1}}, \mu_{\theta_{2}}) \\ &\leq L_{\pi} \|\theta_{1} - \theta_{2}\| + 2L_{\pi}(\lceil \log_{\rho} m^{-1} \rceil + \frac{1}{1 - \rho})\|\theta_{1} - \theta_{2}\| \\ &= 2L_{\pi}(1 + \lceil \log_{\rho} m^{-1} \rceil + \frac{1}{1 - \rho})\|\theta_{1} - \theta_{2}\|. \end{split}$$

For the third inequality, we have

$$\begin{aligned} d_{\mathrm{TV}}(\mu_{\theta_{1}} \otimes \pi_{\theta_{1}} \otimes \mathcal{P}, \mu_{\theta_{2}} \otimes \pi_{\theta_{2}} \otimes \mathcal{P}) \\ &= \frac{1}{2} \int_{S} \int_{\mathcal{A}} \int_{S} |\mu_{\theta_{1}}(ds) \pi_{\theta_{1}}(da \mid s) P(ds' \mid s, a) - \mu_{\theta_{2}}(ds) \pi_{\theta_{2}}(da \mid s) P(ds' \mid s, a)| \\ &= \frac{1}{2} \int_{S} \int_{\mathcal{A}} |\mu_{\theta_{1}}(ds) \pi_{\theta_{1}}(da \mid s) - \mu_{\theta_{2}}(ds) \pi_{\theta_{2}}(da \mid s)| \\ &= d_{TV}(\mu_{\theta_{1}} \otimes \pi_{\theta_{1}}, \mu_{\theta_{2}} \otimes \pi_{\theta_{2}}), \end{aligned}$$

which concludes the proof.

#### Proof of Lemma B.2.

*Proof.* This is a slight modification of the proof of Lemma B.2 in Wu et al. (2020), which extends it to continuous settings. From the fact that

$$\mathbb{P}(s_{t+1} \in \cdot) = \int_{\mathcal{S}} \int_{\mathcal{A}} \mathbb{P}(s_t = ds, a_t = da, s_{t+1} \in \cdot),$$

we have

$$\begin{split} d_{\mathrm{TV}}(\mathbb{P}(s_{t+1} \in \cdot), \mathbb{P}(\tilde{s}_{t+1} \in \cdot)) \\ &= \frac{1}{2} \int_{\mathcal{S}} \int_{\mathcal{A}} \mathbb{P}(s_t = ds, a_t = da, s_{t+1} = ds') - \int_{\mathcal{S}} \int_{\mathcal{A}} \mathbb{P}(\tilde{s}_t = ds, \tilde{a}_t = da, \tilde{s}_{t+1} = ds')| \\ &\leq \frac{1}{2} \int_{\mathcal{S}} \int_{\mathcal{S}} \int_{\mathcal{A}} |\mathbb{P}(s_t = ds, a_t = da, s_{t+1} = ds') - \mathbb{P}(\tilde{s}_t = ds, \tilde{a}_t = da, \tilde{s}_{t+1} = ds')| \\ &= \frac{1}{2} \int_{\mathcal{S}} \int_{\mathcal{S}} \int_{\mathcal{A}} |\mathbb{P}(O_t = (ds, da, ds')) - \mathbb{P}(\tilde{O}_t = (ds, da, ds'))| \\ &= d_{TV}(\mathbb{P}(O_t \in \cdot), \mathbb{P}(\tilde{O} \in \cdot)), \end{split}$$

where the last equality requires the exchange of integral which is guaranteed by Fubini's theorem since  $\mathbb{P}$  is an absolute integrable function.

For the second equality, we have

$$\begin{split} &d_{TV}(\mathbb{P}(O_t \in \cdot), \mathbb{P}(\tilde{O}_t \in \cdot)) \\ &= \int_{\mathcal{S}} \int_{\mathcal{A}} \int_{\mathcal{S}} |\mathbb{P}(O_t = (ds, da, ds')) - \mathbb{P}(\tilde{O}_t = (ds, da, ds'))| \\ &= \frac{1}{2} \int_{\mathcal{S}} \int_{\mathcal{A}} \int_{\mathcal{S}} |P(ds'|s, a)\mathbb{P}((s_t, a_t) = (ds, da)) - P(ds'|s, a)\mathbb{P}((\tilde{s}_t, \tilde{a}_t) = (ds, da))| \\ &= \frac{1}{2} \int_{\mathcal{S}} \int_{\mathcal{A}} \int_{\mathcal{S}} P(ds'|s, a)|\mathbb{P}((s_t, a_t) = (ds, da)) - \mathbb{P}((\tilde{s}_t, \tilde{a}_t) = (ds, da))| \\ &= \frac{1}{2} \int_{\mathcal{S}} \int_{\mathcal{A}} |\mathbb{P}((s_t, a_t) = (ds, da)) - \mathbb{P}((\tilde{s}_t, \tilde{a}_t) = (ds, da))| \\ &= d_{TV}(\mathbb{P}((s_t, a_t) \in \cdot), \mathbb{P}((\tilde{s}_t, \tilde{a}_t) \in \cdot)). \end{split}$$

For the third inequality, since  $\theta_t$  is dependent on  $s_t$ , it holds that

$$\begin{split} d_{\mathrm{TV}}(\mathbb{P}((s_{t},a_{t})\in\cdot),\mathbb{P}(\tilde{s}_{t},\tilde{a}_{t})\in\cdot)) \\ &= \frac{1}{2}\int_{\mathcal{S}}\int_{\mathcal{A}}|\mathbb{P}(s_{t}=ds,a_{t}=da)-\mathbb{P}(\tilde{s}_{t}=ds,\tilde{a}_{t}=da)| \\ &= \frac{1}{2}\int_{\mathcal{S}}\int_{\mathcal{A}}|\int_{\theta}\mathbb{P}(s_{t}=ds)\mathbb{P}(\theta_{t}=d\theta\mid s_{t}=s)\mathbb{P}(a_{t}=da\mid s_{t}=s,\theta_{t}=\theta)-\mathbb{P}(\tilde{s}_{t}=ds,\tilde{a}_{t}=da)| \\ &= \frac{1}{2}\int_{\mathcal{S}}\int_{\mathcal{A}}|\mathbb{P}(s_{t}=ds)\int_{\theta}\mathbb{P}(\theta_{t}=d\theta\mid s_{t}=s)\pi_{\theta_{t}}(da\mid s)-\mathbb{P}(\tilde{s}_{t}=ds)\pi_{\theta_{t-\tau}}(da\mid s)| \\ &= \frac{1}{2}\int_{\mathcal{S}}\int_{\mathcal{A}}|\mathbb{P}(s_{t}=ds)\mathbb{E}[\pi_{\theta_{t}}(da\mid s)\mid s_{t}=s]-\mathbb{P}(\tilde{s}_{t}=ds)\pi_{\theta_{t-\tau}}(da\mid s)| \\ &= \frac{1}{2}\int_{\mathcal{S}}\int_{\mathcal{A}}|\mathbb{P}(s_{t}=ds)\mathbb{E}[\pi_{\theta_{t}}(da\mid s)\mid s_{t}=s]-\mathbb{P}(s_{t}=ds)\pi_{\theta_{t-\tau}}(da\mid s)| \\ &+ \frac{1}{2}\int_{\mathcal{S}}\int_{\mathcal{A}}|\mathbb{P}(s_{t}=ds)\pi_{\theta_{t-\tau}}(da\mid s)-\mathbb{P}(\tilde{s}_{t}=ds)\pi_{\theta_{t-\tau}}(da\mid s)| \\ &= \frac{1}{2}\int_{\mathcal{S}}\mathbb{P}(s_{t}=ds)\int_{\mathcal{A}}|\mathbb{E}[\pi_{\theta_{t}}(da\mid s)|s_{t}=s]-\pi_{\theta_{t-\tau}}(da\mid s)| \\ &= \frac{1}{2}\int_{\mathcal{S}}\mathbb{P}(s_{t}=ds)\int_{\mathcal{A}}|\mathbb{E}[\pi_{\theta_{t}}(da\mid s)|s_{t}=s]-\pi_{\theta_{t-\tau}}(da\mid s)| \\ &= \frac{1}{2}\int_{\mathcal{S}}\mathbb{P}(s_{t}=ds)\int_{\mathcal{A}}|\mathbb{E}[\pi_{\theta_{t}}(da\mid s)|s_{t}=s]-\pi_{\theta_{t-\tau}}(da\mid s)| \\ &= \frac{1}{2}\int_{\mathcal{S}}\mathbb{P}(s_{t}\in\cdot),\mathbb{P}(\tilde{s}_{t}\in\cdot)) \\ &\leq L_{\pi}\mathbb{E}[|\theta_{t}-\theta_{t-\tau}||+d_{TV}(\mathbb{P}(s_{t}\in\cdot),\mathbb{P}(\tilde{s}_{t}\in\cdot))). \end{split}$$

Therefore, we finish our proof.

### Proof of Lemma B.3.

*Proof.* Following the same proof as shown in Lemma B.1. The final results are derived by substituting the results of Lemma B.1 with m = 1 and  $\rho = \gamma$ , as outlined in Proposition 3.2.

#### Proof of Lemma B.4.

*Proof.* By the same proof as shown in Lemma B.2.

#### Proof of Lemma B.5.

Proof. By definition, we have

$$J(\theta_1) - J(\theta_2) = \mathbb{E}[r(s^1, a^1) - r(s^2, a^2)],$$

where  $s^i \sim \nu_{{m heta}_i}, a^i \sim \pi_{{m heta}_i}.$  Therefore, it holds that

$$J(\boldsymbol{\theta}_1) - J(\boldsymbol{\theta}_2) = \mathbb{E}[r(s^1, a^1) - r(s^1, a^1)]$$
  

$$\leq 2\bar{r}d_{TV}(\nu_{\boldsymbol{\theta}_1} \otimes \pi_{\boldsymbol{\theta}_1}, \nu_{\boldsymbol{\theta}_2} \otimes \pi_{\boldsymbol{\theta}_2})$$
  

$$\leq 4\bar{r}L_{\pi}(1 + \frac{1}{1 - \gamma}) \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|$$
  

$$= L_J \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|.$$

### Proof of Lemma B.6.

*Proof.* From Eq. (16), we have

$$A_{\theta}\omega^*(\theta) = b_{\theta}$$

where  $A_{\theta} := \mathbb{E}_{(s,a,s')}[\phi(s)(\phi(s) - \gamma \phi(s'))^{\top})]$  and  $b_{\theta} := \mathbb{E}_{(s,a)}[r(s,a)\phi(s)]$ . The expectation is taken over the stationary distribution  $s \sim \mu_{\theta}$ , the action  $a \sim \pi_{\theta}(\cdot | s)$ , and the transition probability kernel  $s' \sim P(\cdot | s, a)$ .

Denote  $\omega_1^*, \omega_2^*, \hat{\omega}_1$  as the unique solutions of the following equations respectively:

$$oldsymbol{A}_{oldsymbol{ heta}_1}oldsymbol{\omega}_1^*=oldsymbol{b}_{oldsymbol{ heta}_1},\quad oldsymbol{A}_{oldsymbol{ heta}_2}\hat{oldsymbol{\omega}}_1=oldsymbol{b}_1,\quad oldsymbol{A}_{oldsymbol{ heta}_2}oldsymbol{\omega}_2^*=oldsymbol{b}_2$$

First we bound  $\|\boldsymbol{\omega}_1^* - \hat{\boldsymbol{\omega}}_1\|$ . By definition, we have

$$\|\boldsymbol{\omega}_1^* - \hat{\boldsymbol{\omega}}_1\| \leq \|\boldsymbol{A}_{\boldsymbol{ heta}_1}^{-1} - \boldsymbol{A}_{\boldsymbol{ heta}_2}^{-1}\|\|\boldsymbol{b}_{\boldsymbol{ heta}_1}\|.$$

It can be shown that

$$A_{\theta_1}^{-1} - A_{\theta_2}^{-1} = A_{\theta_1}^{-1} (A_{\theta_2} - A_{\theta_1}) A_{\theta_2}^{-1},$$

which implies

$$\|m{\omega}_1^* - \hat{m{\omega}}_1\| \le \|m{A}_{m{ heta}_1}^{-1}\|\|m{A}_{m{ heta}_1} - m{A}_{m{ heta}_2}\|\|m{A}_{m{ heta}_2}^{-1}\|\|m{b}_{m{ heta}_1}\|$$

Then we bound  $\|\hat{\boldsymbol{\omega}}_1 - \boldsymbol{\omega}_2^*\|$ :

$$\|\hat{oldsymbol{\omega}}_1-oldsymbol{\omega}_2^*\|\leq \|oldsymbol{A}_{oldsymbol{ heta}_2}^{-1}\|\|oldsymbol{b}_{oldsymbol{ heta}_1}-oldsymbol{b}_{oldsymbol{ heta}_2}\|.$$

By Assumption 4.1, the eigenvalues of  $A_{\theta_i}$  are bounded from below by  $\lambda > 0$ , therefore  $||A_{\theta_i}^{-1}|| \le \lambda^{-1}$ . Also  $||b_{\theta_i}|| \le \bar{r}$ , due to the assumption that  $|r(s, a)| \le \bar{r}$ , and  $||\phi(s)|| \le 1$ . To bound  $||A_{\theta_1} - A_{\theta_2}||$  and  $||b_{\theta_1} - b_{\theta_2}||$ , we first note that

$$\begin{aligned} \|\boldsymbol{A}_{\boldsymbol{\theta}_{1}} - \boldsymbol{A}_{\boldsymbol{\theta}_{2}}\| &\leq 2 \sup_{s,s' \in \mathcal{S}} \|\boldsymbol{\phi}(s)(\boldsymbol{\phi}(s) - \gamma \boldsymbol{\phi}(s')^{\top})\| \cdot 2d_{TV}(\mathbb{P}(O^{1} \in \cdot), \mathbb{P}(O^{2} \in \cdot)) \\ &\leq 4d_{TV}(\mathbb{P}(O^{1} \in \cdot), \mathbb{P}(O^{2} \in \cdot)), \end{aligned}$$

and

$$\begin{aligned} \|b_{\boldsymbol{\theta}_1} - b_{\boldsymbol{\theta}_2}\| &\leq \|\mathbb{E}[r(s^1, a^1)\boldsymbol{\phi}(s^1)] - \mathbb{E}[r(s^2, a^2)\boldsymbol{\phi}(s^2)]\| \\ &\leq 2\bar{r}d_{\mathrm{TV}}(\mathbb{P}(O^1 \in \cdot), \mathbb{P}(O^2 \in \cdot)), \end{aligned}$$

where  $O^i$  is the tuple obtained by  $s^i \sim \mu_{\theta_i}, a^i \sim \pi_{\theta_i}(\cdot | s^i)$ , and  $s' \sim P(\cdot | s^i, a^i)$ . And the total variation norm can be bounded by Lemma B.1 as:

$$d_{\mathrm{TV}}(\mathbb{P}(O^1 \in \cdot), \mathbb{P}(O^2 \in \cdot)) \le 2L_{\pi}(1 + \lceil \log_{\rho} m^{-1} \rceil + \frac{1}{1-\rho}) \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|$$

Collecting the above results, we have

$$\begin{aligned} \|\boldsymbol{\omega}_{2}^{*}-\boldsymbol{\omega}_{1}^{*}\| &\leq \|\boldsymbol{\omega}_{1}^{*}-\hat{\boldsymbol{\omega}}_{1}\|+\|\hat{\boldsymbol{\omega}}_{1}-\boldsymbol{\omega}_{2}^{*}\|\\ &\leq (8\lambda^{-2}\bar{r}+4\lambda^{-1}\bar{r})L_{\pi}\left(1+\lceil\log_{\rho}m^{-1}\rceil+\frac{1}{1-\rho}\right)\|\boldsymbol{\theta}_{1}-\boldsymbol{\theta}_{2}\|\end{aligned}$$

and we set  $L_c := (8\lambda^{-2}\bar{r} + 4\lambda^{-1}\bar{r})L_{\pi}(1 + \lceil \log_{\rho} m^{-1} \rceil + 1/(1-\rho))$  to obtain the final result.

## Proof of Lemma B.7.

*Proof.* Lemma B.7 is adopted as an assumption in Chen et al. (2021) and Chen & Zhao (2024), but it directly follows from Heidergott & Hordijk (2003), as pointed out by Olshevsky & Gharesifard (2023).

### Proof of Lemma B.8.

*Proof.* See the proof in Lemma 3.2 of Zhang et al. (2020a).

#### Proof of Lemma B.9.

*Proof.* See the proof in Proposition 8 of Chen et al. (2021).

# G. Proof of Markovian Noise

## Proof of Lemma C.1.

*Proof.* We will divide the proof of this lemma into five steps.

**Step 1.** show that for any  $\theta_1, \theta_1, \omega$ , and tuple O(s, a, s'), we have

$$\Lambda(O,\boldsymbol{\omega},\boldsymbol{\theta}_1) - \Lambda(O,\boldsymbol{\omega},\boldsymbol{\theta}_2) \le (8\bar{\boldsymbol{\omega}}\bar{\delta}L_{\pi}(1 + \lceil \log_{\rho}m^{-1}\rceil + \frac{1}{1-\rho}) + 2\bar{\delta}L_c)\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|.$$
(36)

By the definition of  $\Lambda(O, \boldsymbol{\omega}, \boldsymbol{\theta})$  in Eq. (25), we have

$$\begin{split} \Lambda(O,\boldsymbol{\omega},\boldsymbol{\theta}_1) - \Lambda(O,\boldsymbol{\omega},\boldsymbol{\theta}_2) = & \langle \boldsymbol{\omega} - \boldsymbol{\omega}_1^*, \boldsymbol{f}(O,\boldsymbol{\omega}) - \bar{\boldsymbol{f}}(\boldsymbol{\omega},\boldsymbol{\theta}_1) \rangle - \langle \boldsymbol{\omega} - \boldsymbol{\omega}_2^*, \boldsymbol{f}(O,\boldsymbol{\omega}) - \bar{\boldsymbol{f}}(\boldsymbol{\omega},\boldsymbol{\theta}_2) \rangle \\ \leq & \underbrace{\left| \langle \boldsymbol{\omega} - \boldsymbol{\omega}_1^*, \boldsymbol{f}(O,\boldsymbol{\omega}) - \bar{\boldsymbol{f}}(\boldsymbol{\omega},\boldsymbol{\theta}_1) \rangle - \langle \boldsymbol{\omega} - \boldsymbol{\omega}_1^*, \boldsymbol{f}(O,\boldsymbol{\omega}) - \bar{\boldsymbol{f}}(\boldsymbol{\omega},\boldsymbol{\theta}_2) \rangle \right|}_{I_1} \\ & + \underbrace{\left| \langle \boldsymbol{\omega} - \boldsymbol{\omega}_1^*, \boldsymbol{f}(O,\boldsymbol{\omega}) - \bar{\boldsymbol{f}}(\boldsymbol{\omega},\boldsymbol{\theta}_2) \rangle - \langle \boldsymbol{\omega} - \boldsymbol{\omega}_2^*, \boldsymbol{f}(O,\boldsymbol{\omega}) - \bar{\boldsymbol{f}}(\boldsymbol{\omega},\boldsymbol{\theta}_2) \rangle \right|}_{I_2}. \end{split}$$

For term  $I_1$ , we have

$$\begin{split} I_{1} &= \left| \langle \boldsymbol{\omega} - \boldsymbol{\omega}_{1}^{*}, \bar{\boldsymbol{f}}(\boldsymbol{\omega}, \boldsymbol{\theta}_{2}) - \bar{\boldsymbol{f}}(\boldsymbol{\omega}, \boldsymbol{\theta}_{1}) \rangle \right| \\ &\leq 2\bar{\boldsymbol{\omega}} \left\| \bar{\boldsymbol{f}}(\boldsymbol{\omega}, \boldsymbol{\theta}_{2}) - \bar{\boldsymbol{f}}(\boldsymbol{\omega}, \boldsymbol{\theta}_{1}) \right\| \\ &\leq 4\bar{\boldsymbol{\omega}}\bar{\delta}d_{TV}(\mu_{\boldsymbol{\theta}_{1}} \otimes \pi_{\boldsymbol{\theta}_{1}} \otimes \mathcal{P}, \mu_{\boldsymbol{\theta}_{2}} \otimes \pi_{\boldsymbol{\theta}_{2}} \otimes \mathcal{P}) \\ &\stackrel{(1)}{\leq} 8\bar{\boldsymbol{\omega}}\bar{\delta}L_{\pi}(1 + \lceil \log_{\rho}m^{-1} \rceil + \frac{1}{1-\rho}) \|\boldsymbol{\theta}_{1} - \boldsymbol{\theta}_{2}\|, \end{split}$$

where (1) comes from Lemma B.1.

For term  $I_2$ , we have

$$egin{aligned} &I_2 = \left| \langle oldsymbol{\omega}_2^* - oldsymbol{\omega}_1^*, oldsymbol{f}(O,oldsymbol{\omega}) - oldsymbol{ar{f}}(oldsymbol{\omega},oldsymbol{ heta}_2) 
angle 
ight| \ &\leq 2ar{\delta} ig\| oldsymbol{\omega}_1^* - oldsymbol{\omega}_2^* ig\| \ &\leq 2ar{\delta} L_c ig\| oldsymbol{ heta}_1 - oldsymbol{ heta}_2 ig\|, \end{aligned}$$

where (1) follows from Lemma B.6. Combining  $I_1$  and  $I_2$ , we have

$$\Lambda(O,\boldsymbol{\omega},\boldsymbol{\theta}_1) - \Lambda(O,\boldsymbol{\omega},\boldsymbol{\theta}_2) \le (8\bar{\boldsymbol{\omega}}\bar{\delta}L_{\pi}(1 + \lceil \log_{\rho}m^{-1}\rceil + \frac{1}{1-\rho}) + 2\bar{\delta}L_c)\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|.$$

**Step 2.** show that for any  $\theta$ ,  $\omega_1$ ,  $\omega_2$ , and tuple O(s, a, s'), we have

$$\Lambda(O,\boldsymbol{\omega}_1,\boldsymbol{\theta}) - \Lambda(O,\boldsymbol{\omega}_2,\boldsymbol{\theta}) \le (8\bar{\boldsymbol{\omega}} + 2\bar{\boldsymbol{\delta}}) \|\boldsymbol{\omega}_1 - \boldsymbol{\omega}_2\|.$$
(37)

According to the definition, we have

$$\begin{split} \Lambda(O,\omega_1,\boldsymbol{\theta}) - \Lambda(O,\omega_2,\boldsymbol{\theta}) &= \langle \omega_1 - \omega^*, \boldsymbol{f}(O,\omega_1) - \bar{\boldsymbol{f}}(\omega_1,\boldsymbol{\theta}) \rangle - \mathbb{E} \langle \omega_2 - \omega^*, \boldsymbol{f}(O,\omega_2) - \bar{\boldsymbol{f}}(\omega_2,\boldsymbol{\theta}) \rangle \\ &\leq \underbrace{\left| \langle \omega_1 - \omega^*, \boldsymbol{f}(O,\omega_1) - \bar{\boldsymbol{f}}(\omega_1,\boldsymbol{\theta}) \rangle - \langle \omega_1 - \omega^*, \boldsymbol{f}(O,\omega_2) - \bar{\boldsymbol{f}}(\omega_2,\boldsymbol{\theta}) \rangle \right|}_{I_1} \\ &+ \underbrace{\left| \langle \omega_1 - \omega^*, \boldsymbol{f}(O,\omega_2) - \bar{\boldsymbol{f}}(\omega_2,\boldsymbol{\theta}) \rangle - \langle \omega_2 - \omega^*, \boldsymbol{f}(O,\omega_2) - \bar{\boldsymbol{f}}(\omega_2,\boldsymbol{\theta}) \rangle \right|}_{I_2}. \end{split}$$

For term  $I_1$ , we have

$$\begin{split} I_1 &= \left| \langle \boldsymbol{\omega}_1 - \boldsymbol{\omega}^*, \boldsymbol{f}(O, \boldsymbol{\omega}_1) - \bar{\boldsymbol{f}}(\boldsymbol{\omega}_1, \boldsymbol{\theta}) - (\boldsymbol{f}(O, \boldsymbol{\omega}_2) - \bar{\boldsymbol{f}}(\boldsymbol{\omega}_2, \boldsymbol{\theta})) \rangle \right| \\ &\leq 2\bar{\omega} (\|\boldsymbol{f}(O, \boldsymbol{\omega}_1) - \boldsymbol{f}(O, \boldsymbol{\omega}_2)\| + \|\bar{\boldsymbol{f}}(\boldsymbol{\omega}_2, \boldsymbol{\theta}) - \bar{\boldsymbol{f}}(\boldsymbol{\omega}_1, \boldsymbol{\theta})\|) \\ &\leq 8\bar{\omega} \|\boldsymbol{\omega}_1 - \boldsymbol{\omega}_2\|. \end{split}$$

For term  $I_2$ , we have

$$I_2 = \left| \langle \boldsymbol{\omega}_2 - \boldsymbol{\omega}_1, \boldsymbol{f}(O, \boldsymbol{\omega}_2) - \bar{\boldsymbol{f}}(\boldsymbol{\omega}_2, \boldsymbol{\theta}) \rangle \right| \le 2\bar{\delta} \| \boldsymbol{\omega}_1 - \boldsymbol{\omega}_2 \|.$$

Combining  $I_1$  and  $I_2$ , we have

$$\Lambda(O,\boldsymbol{\omega}_1,\boldsymbol{\theta}) - \Lambda(O,\boldsymbol{\omega}_2,\boldsymbol{\theta}) \le (8\bar{\boldsymbol{\omega}} + 2\bar{\boldsymbol{\delta}}) \|\boldsymbol{\omega}_1 - \boldsymbol{\omega}_2\|.$$

**Step 3:** show that for tuples  $O_t = (s_t, a_t, s_{t+1})$  and  $\widetilde{O}_t = (\widetilde{s}_t, \widetilde{a}_t, \widetilde{s}_{t+1})$ . Conditioning on  $\mathcal{F}_{t-\tau}$ , we have

$$\mathbb{E}\left[\Lambda(O_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau}) - \Lambda(\widetilde{O}_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau})\right] \le 4\bar{\omega}\bar{\delta}L_{\pi}\sum_{k=t-\tau}^t \mathbb{E}\|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{t-\tau}\|.$$
(38)

By the definition of total variation distance, we have

$$\mathbb{E}\left[\Lambda(O_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau}) - \Lambda(\widetilde{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})\right] \leq \mathbb{E}\langle\boldsymbol{\omega}_{t-\tau} - \boldsymbol{\omega}_{t-\tau}^{*}, \boldsymbol{f}(O_{t},\boldsymbol{\omega}_{t-\tau}) - \boldsymbol{f}(\widetilde{O}_{t},\boldsymbol{\omega}_{t-\tau})\rangle \\ \leq 4\bar{\omega}\bar{\delta}d_{TV}(\mathbb{P}(O_{t}\in\cdot|\mathcal{F}_{t-\tau}), \mathbb{P}(\widetilde{O}_{t}\in\cdot|\mathcal{F}_{t-\tau})).$$
(39)

By Lemma B.2, we get

$$d_{TV}(\mathbb{P}(O_t \in \cdot | \mathcal{F}_{t-\tau}), \mathbb{P}(\widetilde{O}_t \in \cdot | \mathcal{F}_{t-\tau})) = d_{TV}(\mathbb{P}((s_t, a_t) \in \cdot | \mathcal{F}_{t-\tau}), \mathbb{P}((\widetilde{s}_t, \widetilde{a}_t) \in \cdot | \mathcal{F}_{t-\tau})) \leq d_{TV}(\mathbb{P}(s_t \in \cdot | \mathcal{F}_{t-\tau}), \mathbb{P}(\widetilde{s}_t \in \cdot | \mathcal{F}_{t-\tau})) + L_{\pi}\mathbb{E} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-\tau}\| \leq d_{TV}(\mathbb{P}(O_{t-1} \in \cdot | \mathcal{F}_{t-\tau}), \mathbb{P}(\widetilde{O}_{t-1} \in \cdot | \mathcal{F}_{t-\tau})) + L_{\pi}\mathbb{E} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-\tau}\|.$$

Repeat the above argument from t to  $t - \tau$ , we have

$$d_{TV}(\mathbb{P}(O_t \in \cdot | \mathcal{F}_{t-\tau}), \mathbb{P}(\widetilde{O}_t \in \cdot | \mathcal{F}_{t-\tau})) \le L_{\pi} \sum_{k=t-\tau}^{t} \mathbb{E} \| \boldsymbol{\theta}_k - \boldsymbol{\theta}_{t-\tau} \|.$$
(40)

Plugging Eq. (40) into Eq. (39), we get

$$\mathbb{E}\big[\Lambda(O_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau}) - \Lambda(\widetilde{O}_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau})\big] \le 4\bar{\omega}\bar{\delta}L_{\pi}\sum_{k=t-\tau}^t \mathbb{E}\|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{t-\tau}\|.$$

**Step 4:** show that conditioning on  $\mathcal{F}_{t-\tau}$ , we have

$$\mathbb{E}\left[\Lambda(\widetilde{O}_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau})\right] \le 4\bar{\omega}\bar{\delta}m\rho^{\tau-1}.$$
(41)

It can be shown that

$$\mathbb{E}\big[\Lambda(O'_{t-\tau},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})|\mathcal{F}_{t-\tau}\big]=0.$$

Then we have

$$\mathbb{E}[\Lambda(\widetilde{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})] = \mathbb{E}[\Lambda(\widetilde{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau}) - \Lambda(O'_{t-\tau},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})] \\ = \mathbb{E}[\langle \boldsymbol{\omega}_{t-\tau} - \boldsymbol{\omega}^{*}_{t-\tau}, \boldsymbol{f}(\widetilde{O}_{t},\boldsymbol{\omega}_{t-\tau}) - \boldsymbol{f}(O'_{t-\tau},\boldsymbol{\omega}_{t-\tau})] \\ \leq 4\bar{\omega}\bar{\delta}d_{TV}(\mathbb{P}(\widetilde{O}_{t} = \cdot|\mathcal{F}_{t-\tau}), \mu_{\boldsymbol{\theta}_{t-\tau}} \otimes \pi_{\boldsymbol{\theta}_{t-\tau}} \otimes \mathcal{P}) \\ \stackrel{(1)}{\leq} 4\bar{\omega}\bar{\delta}m\rho^{\tau-1},$$

where (1) follows from Assumption 4.2.

**Step 5:** show that for  $t \ge \tau_{mix}$ , we have

$$\mathbb{E}\big[\Lambda(O_t,\boldsymbol{\omega}_t,\boldsymbol{\theta}_t)\big] \leq M_1 \frac{1}{\sqrt{T}},$$

where  $M_1 = 2\bar{\delta}L_{\pi}(1 + \lceil \log_{\rho} m^{-1} \rceil + (1 - \rho)^{-1})\bar{\delta}B\tau_{\min}c + 4\bar{\delta}\tau_{\min} + \bar{\delta}L_{\pi}\bar{\delta}B\tau_{\min}^2c + 2\bar{\delta}.$ Combining Eq. (36), Eq. (37), Eq. (38), and Eq. (41), we have

$$\begin{split} \mathbb{E}\left[\Lambda(O_{t},\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t})\right] &= \mathbb{E}\left[\Lambda(O_{t},\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t}) - \Lambda(O_{t},\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t-\tau})\right] + \mathbb{E}\left[\Lambda(O_{t},\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t-\tau}) - \Lambda(O_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})\right] \\ &+ \mathbb{E}\left[\Lambda(O_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau}) - \Lambda(\widetilde{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})\right] + \mathbb{E}\left[\Lambda(\widetilde{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})\right] \\ &\leq (8\bar{\omega}\bar{\delta}L_{\pi}(1+\lceil\log_{\rho}m^{-1}\rceil+\frac{1}{1-\rho}) + 2\bar{\delta}L_{c})\|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}_{t-\tau}\| + (8\bar{\omega}+2\bar{\delta})\|\boldsymbol{\omega}_{t} - \boldsymbol{\omega}_{t-\tau}\| \\ &+ 4\bar{\omega}\bar{\delta}L_{\pi}\sum_{k=t-\tau}^{t}\mathbb{E}\|\boldsymbol{\theta}_{k} - \boldsymbol{\theta}_{t-\tau}\| + 4\bar{\omega}\bar{\delta}m\rho^{\tau-1} \\ &\leq (8\bar{\omega}\bar{\delta}L_{\pi}(1+\lceil\log_{\rho}m^{-1}\rceil+\frac{1}{1-\rho}) + 2\bar{\delta}L_{c})\sum_{k=t-\tau}^{t-1}\alpha\bar{\delta}B + (8\bar{\omega}+2\bar{\delta})\sum_{k=t-\tau}^{t-1}\beta\bar{\delta} \\ &+ 4\bar{\omega}\bar{\delta}L_{\pi}\sum_{k=t-\tau}^{t}\sum_{i=t-\tau}^{k-1}\alpha\bar{\delta}B + 4\bar{\omega}\bar{\delta}m\rho^{\tau-1} \\ &\leq (8\bar{\omega}\bar{\delta}L_{\pi}(1+\lceil\log_{\rho}m^{-1}\rceil+\frac{1}{1-\rho}) + 2\bar{\delta}L_{c})\tau\bar{\delta}B\frac{c}{\sqrt{T}} + (8\bar{\omega}+2\bar{\delta})\tau\bar{\delta}\frac{1}{\sqrt{T}} \\ &+ 4\bar{\omega}\bar{\delta}L_{\pi}\tau^{2}\bar{\delta}B\frac{c}{\sqrt{T}} + 4\bar{\omega}\bar{\delta}m\rho^{\tau-1} \\ &\leq ((8\bar{\omega}\bar{\delta}L_{\pi}(1+\lceil\log_{\rho}m^{-1}\rceil+\frac{1}{1-\rho}) + 2\bar{\delta}L_{c})\bar{\delta}B\tau_{\mathrm{mix}}c + (8\bar{\omega}+2\bar{\delta})\bar{\delta}\tau_{\mathrm{mix}} + 4\bar{\omega}L_{\pi}B\bar{\delta}^{2}\tau_{\mathrm{mix}}^{2}c + 4\bar{\omega}\bar{\delta})\frac{1}{\sqrt{T}} \end{split}$$

where (1) comes from the update rule of the critic and the actor, (2) is followed by choosing  $\tau = \tau_{mix}$ . Therefore, we conclude our proof.

#### Proof of Lemma C.2

*Proof.* We will divide the proof of this lemma into five steps.

**Step 1:** show that for any  $O, \boldsymbol{\omega}, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2$ , we have

$$\|\Gamma(O,\boldsymbol{\omega},\boldsymbol{\theta}_1) - \Gamma(O,\boldsymbol{\omega},\boldsymbol{\theta}_2)\| \le (2\bar{\delta}BL_c^2 + 4\bar{\delta}\bar{\omega}BL_s + 4\bar{\omega}L_cL_{\bar{h}})\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|,\tag{42}$$

where  $L_{\bar{h}} = \bar{\delta}L_l + 2BL_c + 4\bar{\delta}BL_{\pi}(1 + (1 - \gamma)^{-1}).$ 

Since  $\Gamma(O, \boldsymbol{\omega}, \boldsymbol{\theta}) = \langle \boldsymbol{\omega} - \boldsymbol{\omega}^*, (\nabla \boldsymbol{\omega}^*)^\top (\bar{\boldsymbol{h}}(\boldsymbol{\omega}^*, \boldsymbol{\theta}) - h(O, \boldsymbol{\omega}^*, \boldsymbol{\theta})) \rangle$ , we represent  $\bar{\boldsymbol{h}}(\boldsymbol{\omega}^*, \boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}}[h(O, \boldsymbol{\omega}^*, \boldsymbol{\theta})]$ , where  $\mathbb{E}_{\boldsymbol{\theta}}$  is the shorthand of  $\mathbb{E}_{O \sim (\nu_{\boldsymbol{\theta}}, \pi_{\boldsymbol{\theta}}, \mathcal{P})}$ . In the following, we will show that each term in  $\Gamma(O, \boldsymbol{\omega}, \boldsymbol{\theta})$  is Lipschitz with respect to  $\boldsymbol{\theta}$ .

Term  $\boldsymbol{\omega}$  is not related to  $\boldsymbol{\theta}$ , term  $\boldsymbol{\omega}^*(\boldsymbol{\theta})$  is  $L_c$ -Lipschitz according to Lemma B.6, and term  $\nabla \boldsymbol{\omega}^*(\boldsymbol{\theta})$  is  $L_s$ -Lipschitz according to Lemma B.9.

For term  $h(O, \omega^*, \theta)$ , we have

$$\begin{aligned} \|\boldsymbol{h}(O,\boldsymbol{\omega}_{1}^{*},\boldsymbol{\theta}_{1}) - \boldsymbol{h}(O,\boldsymbol{\omega}_{2}^{*},\boldsymbol{\theta}_{2})\| &\leq \|\boldsymbol{h}(O,\boldsymbol{\omega}_{1}^{*},\boldsymbol{\theta}_{1}) - \boldsymbol{h}(O,\boldsymbol{\omega}_{1}^{*},\boldsymbol{\theta}_{2})\| + \|\boldsymbol{h}(O,\boldsymbol{\omega}_{1}^{*},\boldsymbol{\theta}_{2}) - \boldsymbol{h}(O,\boldsymbol{\omega}_{2}^{*},\boldsymbol{\theta}_{2})\| \\ &\leq \bar{\delta}\|\nabla\log\pi_{\boldsymbol{\theta}_{1}}(a\,|\,s) - \nabla\log\pi_{\boldsymbol{\theta}_{2}}(a\,|\,s)\| + 2B\|\boldsymbol{\omega}_{1}^{*} - \boldsymbol{\omega}_{2}^{*}\| \\ &\leq (\bar{\delta}L_{l} + 2BL_{c})\|\boldsymbol{\theta}_{1} - \boldsymbol{\theta}_{2}\|. \end{aligned}$$

Hence we have  $h(O, \omega^*, \theta)$  is  $L_h$ -Lipschitz, where  $L_h = \bar{\delta}L_l + 2BL_c$ .

For term  $\mathbb{E}_{\boldsymbol{\theta}}[\boldsymbol{h}(O, \boldsymbol{\omega}^*, \boldsymbol{\theta})]$ , we have

$$\begin{split} \|\mathbb{E}_{\theta_{1}}[h(O,\omega_{1}^{*},\theta_{1})] - \mathbb{E}_{\theta_{2}}[h(O,\omega_{2}^{*},\theta_{2})]\| \\ \leq \|\mathbb{E}_{\theta_{1}}[h(O,\omega_{1}^{*},\theta_{1})] - \mathbb{E}_{\theta_{1}}[h(O,\omega_{2}^{*},\theta_{2})]\| + \|\mathbb{E}_{\theta_{1}}[h(O,\omega_{2}^{*},\theta_{2})] - \mathbb{E}_{\theta_{2}}[h(O,\omega_{2}^{*},\theta_{2})]\| \\ \leq \mathbb{E}_{\theta_{1}}\|h(O,\omega_{1}^{*},\theta_{1}) - h(O,\omega_{2}^{*},\theta_{2})\| + \|\mathbb{E}_{\theta_{1}}[h(O,\omega_{2}^{*},\theta_{2})] - \mathbb{E}_{\theta_{2}}[h(O,\omega_{2}^{*},\theta_{2})]\| \\ \leq L_{h}\|\theta_{1} - \theta_{2}\| + \|\mathbb{E}_{\theta_{1}}[h(O,\omega_{2}^{*},\theta_{2})] - \mathbb{E}_{\theta_{2}}[h(O,\omega_{2}^{*},\theta_{2})]\| \\ \leq L_{h}\|\theta_{1} - \theta_{2}\| + 2\bar{\delta}Bd_{TV}(\nu_{\theta_{1}} \otimes \pi_{\theta_{1}} \otimes P, \nu_{\theta_{2}} \otimes \pi_{\theta_{2}} \otimes P) \\ \stackrel{(1)}{\leq} (L_{h} + 4\bar{\delta}BL_{\pi}(1 + \frac{1}{1 - \gamma}))\|\theta_{1} - \theta_{2}\| \\ = L_{\bar{h}}\|\theta_{1} - \theta_{2}\|, \end{split}$$

where (1) follows from Lemma B.3 and  $L_{\bar{h}} = \bar{\delta}L_l + 2BL_c + 4\bar{\delta}BL_{\pi}(1 + (1 - \gamma)^{-1}).$ 

Then we have  $\boldsymbol{\omega} - \boldsymbol{\omega}_{\boldsymbol{\theta}}^*$  is  $2\bar{\omega}$ -bounded and  $L_c$ -Lipschitz;  $\nabla \boldsymbol{\omega}_{\boldsymbol{\theta}}^*$  is  $L_c$ -bounded and  $L_s$ -Lipschitz;  $\mathbb{E}_{\boldsymbol{\theta}}[\boldsymbol{h}(O, \boldsymbol{\omega}^*, \boldsymbol{\theta})] - \boldsymbol{h}(O, \boldsymbol{\omega}^*, \boldsymbol{\theta})$  is  $2\bar{\delta}B$ -bounded and  $2L_{\bar{h}}$ -Lipschitz. By the triangle inequality, we have

$$\|\Gamma(O,\boldsymbol{\omega},\boldsymbol{\theta}_1) - \Gamma(O,\boldsymbol{\omega},\boldsymbol{\theta}_2)\| \le (2\bar{\delta}BL_c^2 + 4\bar{\delta}\bar{\omega}BL_s + 4\bar{\omega}L_cL_{\bar{h}})\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|.$$

**Step 2:** show that for any  $O, \omega_1, \omega_2, \theta$ , we have

$$\|\Gamma(O,\boldsymbol{\omega}_1,\boldsymbol{\theta}) - \Gamma(O,\boldsymbol{\omega}_2,\boldsymbol{\theta})\| \le 2\bar{\delta}BL_c\|\boldsymbol{\omega}_1 - \boldsymbol{\omega}_2\|.$$
(43)

It can be shown that

$$\|\Gamma(O,\boldsymbol{\omega}_1,\boldsymbol{\theta}) - \Gamma(O,\boldsymbol{\omega}_2,\boldsymbol{\theta})\| = \langle \boldsymbol{\omega}_1 - \boldsymbol{\omega}_2, (\nabla \boldsymbol{\omega}^*)^\top (\bar{\boldsymbol{h}}(\boldsymbol{\omega}^*,\boldsymbol{\theta}) - \boldsymbol{h}(O,\boldsymbol{\omega}^*,\boldsymbol{\theta})) \rangle \le 2\bar{\delta}BL_c \|\boldsymbol{\omega}_1 - \boldsymbol{\omega}_2\|.$$

Step 3: show that for tuples  $\hat{O}_t = (\hat{s}_t, \hat{a}_t, \hat{s}_{t+1})$  and  $\bar{O}_t = (\bar{s}_t, \bar{a}_t, \bar{s}_{t+1})$ . Conditioning on  $\mathcal{F}_{t-\tau}$ , we have

$$\|\mathbb{E}[\Gamma(\widehat{O}_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau}) - \Gamma(\overline{O}_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau})]\| \le 4\bar{\omega}\bar{\delta}BL_cL_{\pi}\sum_{k=t-\tau}^t \mathbb{E}\|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{t-\tau}\|.$$
(44)

By definition of  $\Gamma(O, \boldsymbol{\omega}, \boldsymbol{\theta})$  in Eq. (25), we have

$$\begin{aligned} &\|\mathbb{E}[\Gamma(\widehat{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})-\Gamma(\bar{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})]\|\\ &=\|\mathbb{E}[\langle\boldsymbol{\omega}_{t-\tau}-\boldsymbol{\omega}_{t-\tau}^{*},(\nabla\boldsymbol{\omega}_{t-\tau}^{*})^{\top}(\boldsymbol{h}(\widehat{O}_{t},\boldsymbol{\omega}_{t-\tau}^{*},\boldsymbol{\theta}_{t-\tau})-\boldsymbol{h}(\bar{O}_{t},\boldsymbol{\omega}_{t-\tau}^{*},\boldsymbol{\theta}_{t-\tau}))\rangle]\|\\ &\leq 4\bar{\omega}\bar{\delta}BL_{c}d_{TV}(\mathbb{P}(\widehat{O}_{t}\in\cdot|\mathcal{F}_{t-\tau}),\mathbb{P}(\bar{O}_{t}\in\cdot|\mathcal{F}_{t-\tau})),\end{aligned}$$

$$(45)$$

where the inequality comes from the definition of total variation distance.

By Lemma B.4, we get

$$\begin{aligned} &d_{TV}(\mathbb{P}(\widehat{O}_t \in \cdot | \mathcal{F}_{t-\tau}), \mathbb{P}(\bar{O}_t \in \cdot | \mathcal{F}_{t-\tau})) \\ &= d_{TV}(\mathbb{P}((\hat{s}_t, \hat{a}_t) \in \cdot | \mathcal{F}_{t-\tau}), \mathbb{P}((\bar{s}_t, \bar{a}_t) \in \cdot | \mathcal{F}_{t-\tau})) \\ &\leq d_{TV}(\mathbb{P}(\hat{s}_t \in \cdot | \mathcal{F}_{t-\tau}), \mathbb{P}(\bar{s}_t \in \cdot | \mathcal{F}_{t-\tau})) + L_{\pi} \mathbb{E} \| \boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-\tau} \| \\ &\leq d_{TV}(\mathbb{P}(\widehat{O}_{t-1} \in \cdot | \mathcal{F}_{t-\tau}), \mathbb{P}(\bar{O}_{t-1} \in \cdot | \mathcal{F}_{t-\tau})) + L_{\pi} \mathbb{E} \| \boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-\tau} \|. \end{aligned}$$

Repeat the above argument from t to  $t - \tau$ , we have

$$d_{TV}(\mathbb{P}(\widehat{O}_t \in \cdot | \mathcal{F}_{t-\tau}), \mathbb{P}(\overline{O}_t \in \cdot | \mathcal{F}_{t-\tau})) \le L_{\pi} \sum_{k=t-\tau}^{t} \mathbb{E} \| \boldsymbol{\theta}_k - \boldsymbol{\theta}_{t-\tau} \|.$$
(46)

Plugging Eq. (46) into Eq. (45), we have

$$\|\mathbb{E}[\Gamma(\widehat{O}_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau}) - \Gamma(\overline{O}_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau})]\| \le 4\bar{\omega}\bar{\delta}BL_cL_{\pi}\sum_{k=t-\tau}^t \mathbb{E}\|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{t-\tau}\|$$

**Step 4:** show that conditioning on  $\mathcal{F}_{t-\tau}$ , we have

$$\|\mathbb{E}[\Gamma(\bar{O}_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau})]\| \le 4\bar{\omega}\bar{\delta}BL_c\gamma^{\tau-1}.$$
(47)

It can be shown that

$$\begin{aligned} \|\mathbb{E}[\Gamma(\bar{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})]\| \stackrel{(1)}{=} \|\mathbb{E}[\Gamma(\bar{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau}) - \Gamma(O_{t-\tau}'',\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})]\| \\ \stackrel{(2)}{\leq} 4\bar{\omega}\bar{\delta}BL_{c}d_{TV}(\mathbb{P}(\bar{O}_{t}\in\cdot|\mathcal{F}_{t-\tau}),\nu_{\boldsymbol{\theta}_{t-\tau}}\otimes\pi_{\boldsymbol{\theta}_{t-\tau}}\otimes P), \end{aligned}$$

where (1) is due to the fact that  $O_{t-\tau}''$  is from the discounted state visitation distribution which satisfies  $\mathbb{E}[\Gamma(O_{t-\tau}', \omega_{t-\tau}, \theta_{t-\tau})|\mathcal{F}_{t-\tau}] = 0$  and (2) follows from the definition of total variation distance. From Proposition 3.2, we know that

$$d_{TV}(\mathbb{P}(\bar{s}_t \in \cdot), \nu_{\boldsymbol{\theta}_{t-\tau}}) \leq \gamma^{\tau-1}.$$

Therefore, we have

$$\begin{aligned} \|\mathbb{E}[\Gamma(\bar{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})]\| &\leq 4\bar{\omega}\delta BL_{c}d_{TV}(\mathbb{P}(\bar{O}_{t}\in\cdot|\mathcal{F}_{t-\tau}),\nu_{\boldsymbol{\theta}_{t-\tau}}\otimes\pi_{\boldsymbol{\theta}_{t-\tau}}\otimes P) \\ &= 4\bar{\omega}\bar{\delta}BL_{c}d_{TV}(\mathbb{P}((\bar{s}_{t},\bar{a}_{t})\in\cdot|\mathcal{F}_{t-\tau},\mu_{\boldsymbol{\theta}_{t-\tau}}\otimes\pi_{\boldsymbol{\theta}_{t-\tau}}) \\ &= 4\bar{\omega}\bar{\delta}BL_{c}d_{TV}(\mathbb{P}(\bar{s}_{t}\in\cdot|\mathcal{F}_{t-\tau}),\mu_{\boldsymbol{\theta}_{t-\tau}}) \\ &\leq 4\bar{\omega}\bar{\delta}BL_{c}\gamma^{\tau-1}. \end{aligned}$$

**Step 5:** show that for  $t \ge \tau_{mix}$ , we have

$$\mathbb{E}\big[\Gamma(\widehat{O}_t, \boldsymbol{\omega}_t, \boldsymbol{\theta}_t)\big] \leq M_2 \frac{1}{\sqrt{T}},$$

where  $M_2 = (2\bar{\delta}BL_c^2 + 4\bar{\delta}\bar{\omega}BL_s + 4\bar{\omega}L_cL_{\bar{h}})\bar{\delta}Bc\tau_{\text{mix}} + 2\bar{\delta}BL_c\bar{\delta}\tau_{\text{mix}} + 4\bar{\omega}\bar{\delta}BL_cL_{\pi}\bar{\delta}Bc\tau_{\text{mix}}^2 + 4\bar{\omega}\bar{\delta}BL_c$ . Combining Eq. (42), Eq. (43), Eq. (44), and Eq. (47), we have

$$\begin{split} \mathbb{E}\big[\Gamma(\widehat{O}_{t},\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t})\big] =& \mathbb{E}\big[\Gamma(\widehat{O}_{t},\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t}) - \Gamma(\widehat{O}_{t},\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t-\tau})\big] + \mathbb{E}\big[\Gamma(\widehat{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})\big] \\ &+ \mathbb{E}\big[\Gamma(\widehat{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau}) - \Gamma(\overline{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})\big] + \mathbb{E}\big[\Gamma(\overline{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})\big] \\ &\leq (2\bar{\delta}BL_{c}^{2} + 4\bar{\delta}\bar{\omega}BL_{s} + 4\bar{\omega}L_{c}L_{\bar{h}})\big\|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}_{t-\tau}\big\| + 2\bar{\delta}BL_{c}\big\|\boldsymbol{\omega}_{t} - \boldsymbol{\omega}_{t-\tau}\big\| \\ &+ 4\bar{\omega}\bar{\delta}BL_{c}L_{\pi}\sum_{k=t-\tau}^{t}\mathbb{E}\big\|\boldsymbol{\theta}_{k} - \boldsymbol{\theta}_{t-\tau}\big\| + 4\bar{\omega}\bar{\delta}BL_{c}\gamma^{\tau-1} \\ &\leq (2\bar{\delta}BL_{c}^{2} + 4\bar{\delta}\bar{\omega}BL_{s} + 4\bar{\omega}L_{c}L_{\bar{h}})\sum_{k=t-\tau}^{t-1}\alpha\bar{\delta}B + 2\bar{\delta}BL_{c}\sum_{k=t-\tau}^{t-1}\beta\bar{\delta} \\ &+ 4\bar{\omega}\bar{\delta}BL_{c}L_{\pi}\sum_{k=t-\tau}^{t}\sum_{i=t-\tau}^{k-1}\alpha\bar{\delta}B + 4\bar{\omega}\bar{\delta}BL_{c}\gamma^{\tau-1} \\ &\leq (2\bar{\delta}BL_{c}^{2} + 4\bar{\delta}\bar{\omega}BL_{s} + 4\bar{\omega}L_{c}L_{\bar{h}})\tau\bar{\delta}B\frac{c}{\sqrt{T}} + 2\bar{\delta}BL_{c}\tau\bar{\delta}\frac{1}{\sqrt{T}} \\ &+ 4\bar{\omega}\bar{\delta}BL_{c}L_{\pi}\tau^{2}\bar{\delta}B\frac{c}{\sqrt{T}} + 4\bar{\omega}\bar{\delta}BL_{c}\gamma^{\tau-1} \\ &\leq ((2\bar{\delta}BL_{c}^{2} + 4\bar{\delta}\bar{\omega}BL_{s} + 4\bar{\omega}L_{c}L_{\bar{h}})\bar{\delta}Bc\tau_{\mathrm{mix}} + 2\bar{\delta}BL_{c}\bar{\delta}\tau_{\mathrm{mix}} + 4\bar{\omega}\bar{\delta}BL_{c}L_{\pi}\bar{\delta}Bc\tau_{\mathrm{mix}}^{2} + 4\bar{\omega}\bar{\delta}BL_{c}\big)\frac{1}{\sqrt{T}}, \end{split}$$

where (1) comes from the update rule of the critic and the actor, (2) is followed by choosing  $\tau = \tau_{mix}$ . Thus we conclude our proof.

Proof of Lemma C.3

Proof. We will divide the proof of this lemma into four steps.

**Step 1:** show that for any  $O, \theta_1, \theta_2$ , we have

$$\|\Xi(O,\boldsymbol{\omega},\boldsymbol{\theta}_1) - \Xi(O,\boldsymbol{\omega},\boldsymbol{\theta}_2)\| \le (2\bar{\delta}BL_g + 2L_JL_{\bar{h}})\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|.$$
(48)

Since  $\Xi(O, \omega, \theta) = \langle \nabla J(\theta), \mathbb{E}_{\theta}[h(O, \omega, \theta)] - h(O, \omega, \theta) \rangle$ , we will show that each term in  $\Xi(O, \omega, \theta)$  is Lipschitz.

For the term  $\nabla J(\theta)$ , we know it's  $L_J$ -bounded and  $L_g$ -Lipschitz. For term  $\mathbb{E}_{\theta}[h(O, \omega, \theta)] - h(O, \omega, \theta)$ , by the same argument shown in the proof of Lemma C.2, it's  $2\overline{\delta}B$ -bounded and  $2L_{\overline{h}}$ -Lipschitz. By the triangle inequality, we have

$$\|\Xi(O,\boldsymbol{\omega},\boldsymbol{\theta}_1) - \Xi(O,\boldsymbol{\omega},\boldsymbol{\theta}_2)\| \le (2\delta BL_g + 2L_J L_{\bar{h}})\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|.$$

**Step 2:** show that for any  $O, \omega_1, \omega_2, \theta$ , we have

$$\|\Xi(O,\boldsymbol{\omega}_1,\boldsymbol{\theta}) - \Xi(O,\boldsymbol{\omega}_2,\boldsymbol{\theta})\| \le 4BL_J \|\boldsymbol{\omega}_1 - \boldsymbol{\omega}_2\|.$$
(49)

It follows that

$$\begin{aligned} \|\Xi(O,\omega_1,\boldsymbol{\theta}) - \Xi(O,\omega_2,\boldsymbol{\theta})\| &= |\langle \nabla J(\boldsymbol{\theta}), \boldsymbol{h}(O,\omega_1,\boldsymbol{\theta}) - \boldsymbol{h}(O,\omega_2,\boldsymbol{\theta})\rangle| \\ &+ |\langle \nabla J(\boldsymbol{\theta}), \mathbb{E}_{\boldsymbol{\theta}}[\boldsymbol{h}(O,\omega_1,\boldsymbol{\theta})] - \mathbb{E}_{\boldsymbol{\theta}}[\boldsymbol{h}(O,\omega_2,\boldsymbol{\theta})]\rangle| \\ &\leq 2BL_J \|\omega_1 - \omega_2\| + 2BL_J \|\omega_1 - \omega_2\| \\ &= 4BL_J \|\omega_1 - \omega_2\|. \end{aligned}$$

**Step 3:** show that for tuples  $\hat{O}_t = (\hat{s}_t, \hat{a}_t, \hat{s}_{t+1})$  and  $\bar{O}_t = (\bar{s}_t, \bar{a}_t, \bar{s}_{t+1})$ . Conditioning on  $\mathcal{F}_{t-\tau}$ , we have

$$\|\Xi(\widehat{O}_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau}) - \Xi(\overline{O}_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau})\| \le \bar{\delta}BL_J L_{\pi} \sum_{k=t-\tau}^t \mathbb{E}\|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{t-\tau}\|.$$
(50)

By definition of  $\Xi(O, \boldsymbol{\omega}, \boldsymbol{\theta})$ , we have

$$\begin{aligned} \|\mathbb{E}[\Xi(\widehat{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau}) - \Xi(\overline{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})]\| &= \|\mathbb{E}[\langle \nabla J(\boldsymbol{\theta}_{t-\tau}),\boldsymbol{h}(\widehat{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau}) - \boldsymbol{h}(\overline{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})]\| \\ &\leq 2\bar{\delta}BL_{J}d_{TV}(\mathbb{P}(\widehat{O}_{t} \in \cdot|\mathcal{F}_{t-\tau}),\mathbb{P}(\overline{O}_{t} \in \cdot|\mathcal{F}_{t-\tau})), \end{aligned}$$

where the inequality comes from the definition of total variation distance. The total variation distance between  $\hat{O}_t$  and  $\bar{O}_t$  has been computed in Eq. (46). Plugging Eq. (46) into the above inequality, we get

$$\|\mathbb{E}[\Xi(\widehat{O}_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau}) - \Xi(\overline{O}_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau})]\| \le 2\bar{\delta}BL_J L_{\pi} \sum_{k=t-\tau}^t \mathbb{E}\|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{t-\tau}\|.$$

**Step 4:** show that conditioning on  $\mathcal{F}_{t-\tau}$ , we have

$$\|\mathbb{E}[\Xi(\bar{O}_t, \boldsymbol{\omega}_{t-\tau}, \boldsymbol{\theta}_{t-\tau})]\| \le 2\bar{\delta}BL_J \gamma^{\tau-1}.$$
(51)

It holds that

$$\begin{split} \|\mathbb{E}[\Xi(\bar{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})]\| \stackrel{(1)}{=} \|\mathbb{E}[\Xi(\bar{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau}) - \Xi(O_{t-\tau}'',\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})]\| \\ \stackrel{(2)}{\leq} 2\bar{\delta}BL_{J}d_{TV}(\mathbb{P}(\bar{O}_{t}\in\cdot|\mathcal{F}_{t-\tau}),\nu_{\boldsymbol{\theta}_{t-\tau}}\otimes\pi_{\boldsymbol{\theta}_{t-\tau}}\otimes P) \\ = 2\bar{\delta}BL_{J}d_{TV}(\mathbb{P}((\bar{s}_{t},\bar{a}_{t})\in\cdot|\mathcal{F}_{t-\tau}),\nu_{\boldsymbol{\theta}_{t-\tau}}\otimes\pi_{\boldsymbol{\theta}_{t-\tau}}) \\ = 2\bar{\delta}BL_{J}d_{TV}(\mathbb{P}(\bar{s}_{t}\in\cdot|\mathcal{F}_{t-\tau}),\nu_{\boldsymbol{\theta}_{t-\tau}}) \\ \stackrel{(3)}{\leq} 2\bar{\delta}BL_{J}\gamma^{\tau-1}, \end{split}$$

where (1) is due to the fact that  $O_{t-\tau}''$  is sampled from the discounted state visitation distribution which satisfies  $\mathbb{E}[\Xi(O_{t-\tau}', \omega_{t-\tau}, \theta_{t-\tau}) | \mathcal{F}_{t-\tau}] = 0$ , (2) follows from the definition of total variation distance, and (3) comes from Proposition 3.2.

**Step 5:** show that for  $t \ge \tau_{\min}$ , we have

$$\mathbb{E}\left[\Xi(\widehat{O}_t, \boldsymbol{\omega}_t, \boldsymbol{\theta}_t)\right] \leq M_3 \frac{1}{\sqrt{T}},$$

where  $M_3 = (2\bar{\delta}BL_g + 2L_JL_{\bar{h}})\bar{\delta}Bc\tau_{\min} + 4BL_J\bar{\delta}\tau_{\min} + 2\bar{\delta}^2B^2L_JL_{\pi}c\tau_{\min}^2 + 2\bar{\delta}BL_J.$ 

Combining Eq. (48), Eq. (49), Eq. (50), and Eq. (51), we can decompose the Markovian bias as

$$\begin{split} \mathbb{E}\big[\Xi(\widehat{O}_{t},\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t})\big] &= \mathbb{E}\big[\Xi(\widehat{O}_{t},\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t}) - \Xi(\widehat{O}_{t},\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t-\tau})\big] + \mathbb{E}\big[\Xi(\widehat{O}_{t},\boldsymbol{\omega}_{t},\boldsymbol{\theta}_{t-\tau}) - \Xi(\widehat{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})\big] \\ &\quad + \mathbb{E}\big[\Xi(\widehat{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau}) - \Xi(\overline{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})\big] + \mathbb{E}\big[\Xi(\overline{O}_{t},\boldsymbol{\omega}_{t-\tau},\boldsymbol{\theta}_{t-\tau})\big] \\ &\leq (2\bar{\delta}BL_{g} + 2L_{J}L_{\bar{h}}) \|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}_{t-\tau}\| + 4BL_{J}\|\boldsymbol{\omega}_{1} - \boldsymbol{\omega}_{2}\| \\ &\quad + 2\bar{\delta}BL_{J}L_{\pi} \sum_{k=t-\tau}^{t} \mathbb{E}\|\boldsymbol{\theta}_{k} - \boldsymbol{\theta}_{t-\tau}\| + 2\bar{\delta}BL_{J}\gamma^{\tau-1} \\ &\stackrel{(1)}{\leq} (2\bar{\delta}BL_{g} + 2L_{J}L_{\bar{h}}) \sum_{k=t-\tau}^{t-1} \alpha\bar{\delta}B + 4BL_{J} \sum_{k=t-\tau}^{t-1} \beta\bar{\delta} + 2\bar{\delta}BL_{J}L_{\pi} \sum_{k=t-\tau}^{t} \sum_{i=t-\tau}^{k-1} \alpha\bar{\delta}B + 2\bar{\delta}BL_{J}\gamma^{\tau-1} \\ &\stackrel{(2)}{\leq} ((2\bar{\delta}BL_{g} + 2L_{J}L_{\bar{h}})\bar{\delta}Bc\tau_{\mathrm{mix}} + 4BL_{J}\bar{\delta}\tau_{\mathrm{mix}} + 2\bar{\delta}^{2}B^{2}L_{J}L_{\pi}c\tau_{\mathrm{mix}}^{2} + 2\bar{\delta}BL_{J})\frac{1}{\sqrt{T}}, \end{split}$$

where (1) owes to the update rule of the actor and (2) is followed by choosing  $\tau = \tau_{mix}$ . Hence we conclude our proof.  $\Box$