

# Global Optimality of Single-Timescale Actor-Critic under Continuous State-Action Space: A Study on Linear Quadratic Regulator

Xuyang Chen<sup>1</sup>, Jingliang Duan<sup>2</sup>, Lin Zhao<sup>1</sup>

<sup>1</sup>National University of Singapore

<sup>2</sup>University of Science and Technology Beijing

chenxuyang@u.nus.edu, duanj@ustb.edu.cn, elezhli@nus.edu.sg

## Abstract

Actor-critic methods have achieved state-of-the-art performance in various challenging tasks. However, theoretical understandings of their performance remain elusive and challenging. Existing studies mostly focus on practically uncommon variants such as double-loop or two-timescale stepsize actor-critic algorithms for simplicity. These results certify local convergence on finite state- or action-space only. We push the boundary to investigate the classic single-sample single-timescale actor-critic on continuous (infinite) state-action space, where we employ the canonical linear quadratic regulator (LQR) problem as a case study. We show that the popular single-timescale actor-critic can attain an epsilon-optimal solution with an order of epsilon to -2 sample complexity for solving LQR on the demanding continuous state-action space. Our work provides new insights into the performance of single-timescale actor-critic, which further bridges the gap between theory and practice.

## 1 Introduction

Actor-critic (AC) methods achieved substantial success in solving many difficult reinforcement learning (RL) problems [LeCun *et al.*, 2015; Mnih *et al.*, 2016; Silver *et al.*, 2017]. In addition to a policy update, AC methods employ a parallel critic update to bootstrap the Q-value for policy gradient estimation, which often enjoys reduced variance and fast convergence in training.

Despite the empirical success, theoretical analysis of AC in the most practical form remains challenging. Existing works mostly focus on either the double-loop or the two-timescale variants. In double-loop AC, the actor is updated in the outer loop only after the critic takes sufficiently many steps to have an accurate estimation of the Q-value in the inner loop [Yang *et al.*, 2019; Kumar *et al.*, 2019; Wang *et al.*, 2019]. Hence, the convergence of the critic is decoupled from that of the actor. The analysis is separated into a policy evaluation sub-problem in the inner loop and a perturbed gradient descent in the outer loop. In two-timescale AC, the actor and the critic are updated simultaneously in each iteration using stepsizes of different timescales.

The actor stepsize (denoted by  $\alpha_t$  in the sequel) is typically smaller than that of the critic (denoted by  $\beta_t$  in the sequel), with their ratio going to zero as the iteration number goes to infinity (i.e.,  $\lim_{t \rightarrow \infty} \alpha_t / \beta_t = 0$ ). The two-timescale allows the critic to approximate the correct Q-value asymptotically. This special stepsize design essentially decouples the analysis of the actor and the critic.

The aforementioned AC variants are considered mainly for the ease of analysis, which, however, are uncommon in practical implementations. In practice, the single-timescale AC, where the actor and the critic are updated simultaneously using constantly proportional stepsizes (i.e., with  $\alpha_t / \beta_t = c > 0$ ), is more favorable due to its simplicity of implementation and empirical sample efficiency [Schulman *et al.*, 2015; Mnih *et al.*, 2016]. For online learning, the actor and the critic update only once with a single sample in each iteration using proportional stepsizes. This single-sample single-timescale AC is the most classic AC algorithm extensively discussed in the literature and introduced in [Sutton and Barto, 2018]. However, its analysis is significantly more difficult than other variants, primarily due to the more inaccurate value estimation of the critic update and the stronger coupling between critic and actor. More recent works [Chen *et al.*, 2021; Olshevsky and Ghahserifard, 2023; Chen and Zhao, 2022] investigated its local convergence and on the finite state- or action-space only. Given that most practical applications in real world are of continuous state-action space, it is demanding to ask the following challenging question:

*Can the classic single-sample single-timescale AC find a global optimal policy on continuous state-action space?*

To this end, we take a first step to consider the Linear Quadratic Regulation (LQR), a fundamental continuous state-action space control problem that is commonly employed to study the performance and the limits of RL algorithms [Fazel *et al.*, 2018; Yang *et al.*, 2019; Tu and Recht, 2018; Duan *et al.*, 2023]. We analyze the same classic single-sample single-timescale AC algorithm as those studied in the references listed in Table 1. As compared in Table 1, our result is the first to show the global optimality on continuous (infinite) state-action space, while achieving the sample complexity as the previous studies.

Specifically, we consider the time-average cost, which is a more common case for LQR formulation and more difficult to analyze than the discounted cost. The single-sample

Reference	Setting		Optimality	Sample Complexity
	State Space	action space		
[Chen <i>et al.</i> , 2021]	infinite	finite	local	$\mathcal{O}(\epsilon^{-2})$
[Olshevsky and Ghahserifard, 2023]	finite	finite	local	$\mathcal{O}(\epsilon^{-2})$
[Chen and Zhao, 2022]	infinite	finite	local	$\mathcal{O}(\epsilon^{-2})$
This Paper	infinite	infinite	global	$\mathcal{O}(\epsilon^{-2})$

Table 1: Comparison with other single-sample single-timescale actor-critic algorithms

single-timescale AC algorithm for solving LQR consists of three parallel updates in each iteration: the cost estimator, the critic, and the actor. Unlike the aforementioned double-loop or two-timescale, there is no specialized design in single-sample single-timescale AC that facilitates a decoupled analysis of its three interconnected updates. In fact, it is both conservative and difficult to bound the three iterations separately. Moreover, the existing perturbed gradient analysis can no longer be applied to establish the convergence of the actor either.

To tackle these challenges in analysis, we instead directly bound the overall interconnected iteration system altogether, without resorting to conservative decoupled analysis. In particular, despite the inaccurate estimation in all three updates, we prove the estimation errors diminish to zero if the (constant) ratio of the stepsizes between the actor and the critic is below a threshold. The identified threshold provides new insights into the practical choices of the stepsizes for single-timescale AC.

Compared with other single-sample single-timescale AC (see Table 1), the state-action space we study is infinite. We emphasize that moving from finite to infinite state-action space is highly nontrivial and requires significant analysis. Existing works [Chen *et al.*, 2021; Chen and Zhao, 2022] derived key intermediate results such as many Lipschitz constants relying on the finite size of the state-action space ( $|\mathcal{S}|, |\mathcal{A}|$ ). These results however become immaterial in the infinite state-action space scenario. Some other analysis [Olshevsky and Ghahserifard, 2023] concatenates all state-action pairs to create a finite-dimensional feature matrix. However, this will not be possible when the state-action space is infinite. Consequently, existing analyses are not applicable in our context.

We also distinguish our work from other model-free RL algorithms for solving LQR in Table 2, in addition to AC methods. The zeroth-order methods and the policy iteration method are included for completeness. In particular, we note that [Zhou and Lu, 2023] analyzed the single-timescale AC under a multi-sample setting, where the critics are updated by the least square temporal difference (LSTD) estimator. The idea is still to obtain an accurate policy gradient estimation at each iteration by using sufficient samples (in LSTD), and then follow the common perturbed gradient analysis to prove the convergence of the actor, which decouples the convergence analysis of the actor and the critic. Moreover, the analysis requires a strong assumption on the uniform boundedness of the critic parameters. In comparison, our analysis does not require this assumption and considers the more classic and challenging single-sample setting which is also considered by

the previous works as listed in Table 1.

Overall, our contributions are summarized as follows:

- Our work furthers the theoretical understanding of AC on continuous state-action space, which represents the most practical usages. We for the first time show that the single-sample single-timescale AC can provably find the  $\epsilon$ -accurate global optimum with a sample complexity of  $\mathcal{O}(\epsilon^{-2})$  for tasks with unbounded continuous state-action space. The previous works consider the more restricted finite state-action space settings with only local convergence guarantee [Chen *et al.*, 2021; Olshevsky and Ghahserifard, 2023; Chen and Zhao, 2022].

- We also contribute to the work of RL on continuous control tasks. It is novel that even with the actor updated by a roughly estimated gradient, the single-sample single-timescale AC algorithm can still find the global optimal policy for LQR, under general assumptions. Compared with all other model-free RL algorithms for solving LQR (see Table 2), our work adopts the simplest single-sample single-timescale structure, which may serve as the first step towards understanding the limits of AC methods on continuous control tasks. In addition, compared with the state-of-the-art double-loop AC for solving LQR [Yang *et al.*, 2019], we improve the sample complexity from  $\mathcal{O}(\epsilon^{-5})$  to  $\mathcal{O}(\epsilon^{-2})$ . We also show the algorithm is much more sample-efficient empirically compared to a few classic works in Experiments, which unveils the practical wisdom of AC algorithm.

## 1.1 Related Work

In this section, we review the existing works that are most relevant to ours.

**Actor-Critic methods.** The AC algorithm was proposed by [Konda and Tsitsiklis, 1999]. [Kakade, 2001] extended it to the natural AC algorithm. The asymptotic convergence of AC algorithms has been well established in [Kakade, 2001; Bhatnagar *et al.*, 2009; Castro and Meir, 2010; Zhang *et al.*, 2020]. Many recent works focused on the finite-time convergence of AC methods. Under the double-loop setting, [Yang *et al.*, 2019] established the global convergence of AC methods for solving LQR. [Wang *et al.*, 2019] studied the global convergence of AC methods with both the actor and the critic being parameterized by neural networks. [Kumar *et al.*, 2019] studied the finite-time local convergence of a few AC variants with linear function approximation. Under the two-timescale AC setting, [Wu *et al.*, 2020; Xu *et al.*, 2020] established the finite-time convergence to a stationary point at a sample complexity of  $\mathcal{O}(\epsilon^{-2.5})$ . Under the single-timescale setting, all the related works [Chen *et al.*, 2021; Olshevsky and Ghahserifard, 2023; Chen and Zhao, 2022]

Reference	Algorithm	Structure	
[Fazel <i>et al.</i> , 2018]	zeroth-order	double-loop	
[Malik <i>et al.</i> , 2019]	zeroth-order		
[Yang <i>et al.</i> , 2019]	actor-critic		
[Krauth <i>et al.</i> , 2019]	policy iteration	multi-sample	
[Zhou and Lu, 2023]	actor-critic	single-timescale	multi-sample
This paper	actor-critic	single-timescale	single-sample

Table 2: Comparison with other model-free RL algorithms for solving LQR.

have been reviewed in the Introduction.

**RL algorithms for LQR.** RL algorithms in the context of LQR have seen increased interest in the recent years. These works can be mainly divided into two categories: model-based methods [Dean *et al.*, 2018; Mania *et al.*, 2019; Cohen *et al.*, 2019; Dean *et al.*, 2020] and model-free methods. Our main interest lies in the model-free methods. Notably, [Fazel *et al.*, 2018] established the first global convergence result for LQR under the policy gradient method using zeroth-order optimization. [Krauth *et al.*, 2019] studied the convergence and sample complexity of the LSTD policy iteration method under the LQR setting. On the subject of adopting AC to solve LQR, [Yang *et al.*, 2019] provided the first finite-time analysis with convergence guarantee and sample complexity under the double-loop setting. [Zhou and Lu, 2023] considered the multi-sample (LSTD) and single-timescale setting. For the more practical yet challenging single-sample single-timescale AC, there is no such theoretical guarantee so far, which is the focus of this paper.

**Notation.** We use non-bold letters to denote scalars and use lower and upper case bold letters to denote vectors and matrices respectively. We also use  $\|\omega\|$  to denote the  $\ell_2$ -norm of a vector  $\omega$ ,  $\|\mathbf{A}\|$  to denote the spectral norm of a matrix  $\mathbf{A}$ , and  $\|\mathbf{A}\|_F$  to denote the Frobenius norm of a matrix  $\mathbf{A}$ . We use  $\text{Tr}(\cdot)$  to denote the trace of a matrix. For any symmetric matrix  $\mathbf{M} \in \mathbb{R}^{n \times n}$ , let  $\text{svec}(\mathbf{M}) \in \mathbb{R}^{n(n+1)/2}$  denote the vectorization of the upper triangular part of  $\mathbf{M}$  such that  $\|\mathbf{M}\|_F^2 = \langle \text{svec}(\mathbf{M}), \text{svec}(\mathbf{M}) \rangle$ . Besides, let  $\text{smat}(\cdot)$  denote the inverse of  $\text{svec}(\cdot)$  so that  $\text{smat}(\text{svec}(\mathbf{M})) = \mathbf{M}$ . Finally, we denote by  $\mathbf{A} \otimes_s \mathbf{B}$  the symmetric Kronecker product [Schacke, 2004] of two matrices  $\mathbf{A}$  and  $\mathbf{B}$ .

## 2 Preliminaries

In this section, we introduce the AC algorithm and provide the theoretical background of LQR.

### 2.1 Actor-Critic Algorithms

We consider the reinforcement learning for the standard Markov Decision Process (MDP) defined by  $(\mathcal{X}, \mathcal{U}, \mathcal{P}, c)$ , where  $\mathcal{X}$  is the state space,  $\mathcal{U}$  is the action space,  $\mathcal{P}(x_{t+1}|x_t, u_t)$  denotes the transition kernel that the agent transits to state  $x_{t+1}$  after taking action  $u_t$  at current state  $x_t$ , and  $c(x_t, u_t)$  is the running cost. A policy  $\pi_\theta(u|x)$  parameterized by  $\theta$  is defined as a mapping from a given state to a probability distribution over actions.

In this paper, we aim to find a policy  $\pi_\theta$  that minimizes the

infinite-horizon time-average cost, which is given by

$$J(\theta) := \lim_{T \rightarrow \infty} \mathbb{E}_\theta \frac{\sum_{t=0}^T c(x_t, u_t)}{T} = \mathbb{E}_{x \sim \rho_\theta, u \sim \pi_\theta} [c(x, u)], \quad (1)$$

where  $\rho_\theta$  denotes the stationary state distribution generated by policy  $\pi_\theta$ . In the time-average cost setting, the state-action value (Q-value) of policy  $\pi_\theta$  is defined as

$$Q_\theta(x, u) = \mathbb{E}_\theta \left[ \sum_{t=0}^{\infty} (c(x_t, u_t) - J(\theta)) | x_0 = x, u_0 = u \right],$$

which describes the accumulated differences between running costs and average cost for selecting  $u$  in state  $x$  and thereafter following policy  $\pi_\theta$  [Sutton and Barto, 2018]. Based on this definition, we can use the policy gradient theorem [Sutton *et al.*, 1999] to express the gradient of  $J(\theta)$  with respect to  $\theta$  as

$$\nabla_\theta J(\theta) = \mathbb{E}_{x \sim \rho_\theta, u \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(u|x) Q_\theta(x, u)]. \quad (2)$$

One can also choose to update the policy using the natural policy gradient [Kakade, 2001], which is given by

$$\nabla_\theta^N J(\theta) = F(\theta)^\dagger \nabla_\theta J(\theta). \quad (3)$$

where

$$F(\theta) = \mathbb{E}_{x \sim \rho_\theta, u \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(u|x) \nabla_\theta \log \pi_\theta(u|x)^\top]$$

is the Fisher information matrix and  $F(\theta)^\dagger$  denotes its Moore Penrose pseudoinverse.

Optimizing  $J(\theta)$  in (1) with (2) requires evaluating the Q-value of the current policy  $\pi_\theta$ , which is usually unknown. AC estimates both the Q-value and the policy. The critic update approximates Q-value towards the actual value of the current policy  $\pi_\theta$  using temporal difference (TD) learning [Sutton and Barto, 2018]. The actor improves the policy to reduce the time-average cost  $J(\theta)$  via policy gradient descent. Note that the AC with a natural policy gradient is also known as natural AC, which is a variant of AC.

### 2.2 Actor-Critic for Linear Quadratic Regulator

In this paper, we aim to demystify the convergence property of AC by focusing on the infinite-horizon time-average linear quadratic regulator (LQR) problem:

$$\begin{aligned} & \underset{\{u_t\}}{\text{minimize}} && J(\{u_t\}) := \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T x_t^\top \mathbf{Q} x_t + u_t^\top \mathbf{R} u_t \right] \\ & \text{subject to} && x_{t+1} = \mathbf{A} x_t + \mathbf{B} u_t + \epsilon_t, \end{aligned} \quad (4)$$

where  $\mathbf{x}_t \in \mathbb{R}^d$  is the state and  $\mathbf{u}_t \in \mathbb{R}^k$  is the control action at time  $t$ ;  $\mathbf{A} \in \mathbb{R}^{d \times d}$  and  $\mathbf{B} \in \mathbb{R}^{d \times k}$  are system matrices, and the  $(\mathbf{A}, \mathbf{B})$ -pair is stabilizable;  $\mathbf{Q} \in \mathbb{S}^{d \times d}$  and  $\mathbf{R} \in \mathbb{S}^{k \times k}$  are symmetric positive definite performance matrices, and hence, the  $(\mathbf{A}, \mathbf{Q}^{1/2})$ -pair is immediately observable;  $\epsilon_t \sim \mathcal{N}(0, \mathbf{D}_0)$  are i.i.d Gaussian random variables with positive definite covariance  $\mathbf{D}_0 \succ 0$ . From the optimal control theory [Anderson and Moore, 2007], the optimal policy of (4) is a linear feedback of the state

$$\mathbf{u}_t = -\mathbf{K}^* \mathbf{x}_t, \quad (5)$$

where  $\mathbf{K}^* \in \mathbb{R}^{k \times d}$  is the optimal policy which can be uniquely found by solving an Algebraic Riccati Equation (ARE) [Anderson and Moore, 2007] depending on  $\mathbf{A}, \mathbf{B}, \mathbf{Q}, \mathbf{R}$ . This means that finding  $\mathbf{K}^*$  using ARE relies on the complete model knowledge.

In the sequel, we pursue finding the optimal policy in a *model-free* way by using the AC method, without knowing or estimating  $\mathbf{A}, \mathbf{B}, \mathbf{Q}, \mathbf{R}$ . The structure of the optimal policy in (5) allows us to reformulate (4) as a static optimization problem over all feasible policy matrix  $\mathbf{K} \in \mathbb{R}^{k \times d}$ . To encourage exploration, we parameterize the policy as

$$\{\pi_{\mathbf{K}}(\cdot | \mathbf{x}) = \mathcal{N}(-\mathbf{K}\mathbf{x}, \sigma^2 \mathbf{I}_k), \mathbf{K} \in \mathbb{R}^{k \times d}\}, \quad (6)$$

where  $\mathcal{N}(\cdot, \cdot)$  denotes the Gaussian distribution and  $\sigma > 0$  is the standard deviation of the exploration noise. In other words, given a state  $\mathbf{x}_t$ , the agent will take an action  $\mathbf{u}_t$  according to  $\mathbf{u}_t = -\mathbf{K}\mathbf{x}_t + \sigma\zeta_t$ , where  $\zeta_t \sim \mathcal{N}(0, \mathbf{I}_k)$ . As a consequence, the optimization problem defined in (4) under policy (6) can be reformulated as

$$\underset{\mathbf{K}}{\text{minimize}} \quad J(\mathbf{K}) := \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T \mathbf{x}_t^\top \mathbf{Q} \mathbf{x}_t + \mathbf{u}_t^\top \mathbf{R} \mathbf{u}_t \right] \quad (7)$$

subject to

$$\begin{aligned} \mathbf{u}_t &= -\mathbf{K}\mathbf{x}_t + \sigma\zeta_t, \\ \mathbf{x}_{t+1} &= \mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t + \epsilon_t. \end{aligned} \quad (8)$$

Therefore, the closed-loop form of system (8) is given by

$$\mathbf{x}_{t+1} = (\mathbf{A} - \mathbf{B}\mathbf{K})\mathbf{x}_t + \xi_t, \quad (9)$$

where  $\xi_t = \epsilon_t + \sigma\mathbf{B}\zeta_t \sim \mathcal{N}(0, \mathbf{D}_\sigma)$  with  $\mathbf{D}_\sigma = \mathbf{D}_0 + \sigma^2 \mathbf{B}\mathbf{B}^\top$ . Note that optimizing over the set of stochastic policies (6) will lead to the same optimal  $\mathbf{K}^*$ . From (9), a policy  $\mathbf{K}$  is stabilizing if and only if  $\rho(\mathbf{A} - \mathbf{B}\mathbf{K}) < 1$ , where  $\rho(\cdot)$  denotes the spectral radius. It is well known that if  $\mathbf{K}$  is stabilizing, the Markov chain in (9) yields a stationary state distribution  $\rho_{\mathbf{K}} \sim \mathcal{N}(0, \mathbf{D}_{\mathbf{K}})$ , where  $\mathbf{D}_{\mathbf{K}}$  satisfies the following Lyapunov equation (by taking the variance of (9))

$$\mathbf{D}_{\mathbf{K}} = \mathbf{D}_\sigma + (\mathbf{A} - \mathbf{B}\mathbf{K})\mathbf{D}_{\mathbf{K}}(\mathbf{A} - \mathbf{B}\mathbf{K})^\top. \quad (10)$$

Similarly, we define  $\mathbf{P}_{\mathbf{K}}$  as the unique positive definite solution to (Bellman equation under  $\mathbf{K}$ )

$$\mathbf{P}_{\mathbf{K}} = \mathbf{Q} + \mathbf{K}^\top \mathbf{R} \mathbf{K} + (\mathbf{A} - \mathbf{B}\mathbf{K})^\top \mathbf{P}_{\mathbf{K}} (\mathbf{A} - \mathbf{B}\mathbf{K}). \quad (11)$$

Based on  $\mathbf{D}_{\mathbf{K}}$  and  $\mathbf{P}_{\mathbf{K}}$ , the following lemma characterizes  $J(\mathbf{K})$  and its gradient  $\nabla_{\mathbf{K}} J(\mathbf{K})$ .

**Lemma 1** ([Yang *et al.*, 2019]). *For any stabilizing policy  $\mathbf{K}$ , the time-average cost  $J(\mathbf{K})$  and its gradient  $\nabla_{\mathbf{K}} J(\mathbf{K})$  take the following forms*

$$J(\mathbf{K}) = \text{Tr}(\mathbf{P}_{\mathbf{K}} \mathbf{D}_\sigma) + \sigma^2 \text{Tr}(\mathbf{R}), \quad (12a)$$

$$\nabla_{\mathbf{K}} J(\mathbf{K}) = 2\mathbf{E}_{\mathbf{K}} \mathbf{D}_{\mathbf{K}}, \quad (12b)$$

where  $\mathbf{E}_{\mathbf{K}} := (\mathbf{R} + \mathbf{B}^\top \mathbf{P}_{\mathbf{K}} \mathbf{B})\mathbf{K} - \mathbf{B}^\top \mathbf{P}_{\mathbf{K}} \mathbf{A}$ .

Then, the natural gradient of  $J(\mathbf{K})$  can be calculated as [Fazel *et al.*, 2018; Yang *et al.*, 2019]

$$\nabla_{\mathbf{K}}^N J(\mathbf{K}) = \nabla_{\mathbf{K}} J(\mathbf{K}) \mathbf{D}_{\mathbf{K}}^{-1} = \mathbf{E}_{\mathbf{K}}, \quad (13)$$

which eliminates the burden of estimating  $\mathbf{D}_{\mathbf{K}}$ . Note that we omit the constant coefficient since it can be absorbed by the stepsize.

Calculating the natural gradient  $\nabla_{\mathbf{K}}^N J(\mathbf{K})$  requires estimating  $\mathbf{P}_{\mathbf{K}}$ , which depends on  $\mathbf{A}, \mathbf{B}, \mathbf{Q}, \mathbf{R}$ . To estimate the gradient without the knowledge of the model, we instead directly utilize the Q-value.

**Lemma 2** ([Bradtke *et al.*, 1994; Yang *et al.*, 2019]). *For any stabilizing policy  $\mathbf{K}$ , the Q-value  $Q_{\mathbf{K}}(\mathbf{x}, \mathbf{u})$  takes the following form*

$$\begin{aligned} Q_{\mathbf{K}}(\mathbf{x}, \mathbf{u}) &= (\mathbf{x}^\top, \mathbf{u}^\top) \Omega_{\mathbf{K}} \begin{pmatrix} \mathbf{x} \\ \mathbf{u} \end{pmatrix} - \text{Tr}(\mathbf{P}_{\mathbf{K}} \mathbf{D}_{\mathbf{K}}) \\ &\quad - \sigma^2 \text{Tr}(\mathbf{R} + \mathbf{P}_{\mathbf{K}} \mathbf{B} \mathbf{B}^\top), \end{aligned} \quad (14)$$

where

$$\Omega_{\mathbf{K}} := \begin{bmatrix} \Omega_{\mathbf{K}}^{11} & \Omega_{\mathbf{K}}^{12} \\ \Omega_{\mathbf{K}}^{21} & \Omega_{\mathbf{K}}^{22} \end{bmatrix} := \begin{bmatrix} \mathbf{Q} + \mathbf{A}^\top \mathbf{P}_{\mathbf{K}} \mathbf{A} & \mathbf{A}^\top \mathbf{P}_{\mathbf{K}} \mathbf{B} \\ \mathbf{B}^\top \mathbf{P}_{\mathbf{K}} \mathbf{A} & \mathbf{R} + \mathbf{B}^\top \mathbf{P}_{\mathbf{K}} \mathbf{B} \end{bmatrix}. \quad (15)$$

Clearly, if we can estimate  $\Omega_{\mathbf{K}}$ , then  $\mathbf{E}_{\mathbf{K}}$  in (13) can be readily estimated by using  $\Omega_{\mathbf{K}}^{21}$  and  $\Omega_{\mathbf{K}}^{22}$ , which represent the bottom left corner block and bottom right corner block of matrix  $\Omega_{\mathbf{K}}$ , respectively.

### 3 Single-sample Single-timescale Actor-Critic

In this section, we describe the single-sample single-timescale AC algorithm for solving LQR. In view of the structure of the Q-value given in (14) and the fact that [Schacke, 2004]

$$(\mathbf{x}^\top, \mathbf{u}^\top) \Omega_{\mathbf{K}} \begin{pmatrix} \mathbf{x} \\ \mathbf{u} \end{pmatrix} = \phi(\mathbf{x}, \mathbf{u})^\top \text{svec}(\Omega_{\mathbf{K}}), \quad (16)$$

where

$$\phi(\mathbf{x}, \mathbf{u}) := \text{svec} \left[ \begin{pmatrix} \mathbf{x} \\ \mathbf{u} \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{u} \end{pmatrix}^\top \right] \quad (17)$$

and  $\text{svec}(\cdot)$  denotes the vectorization of the upper triangular part of a symmetric matrix as defined in [Schacke, 2004]. We can then parameterize the Q-estimator (critic) by

$$\hat{Q}_{\mathbf{K}}(\mathbf{x}, \mathbf{u}; \boldsymbol{\omega}, b) = \phi(\mathbf{x}, \mathbf{u})^\top \boldsymbol{\omega} + b,$$

where  $\phi(\mathbf{x}, \mathbf{u})$  defined in (17) is the feature function and  $\boldsymbol{\omega}$  is the critic. Using the TD(0) learning, the critic update is followed by

$$\begin{aligned} \boldsymbol{\omega}_{t+1} &= \boldsymbol{\omega}_t + \beta_t [(c_t - J(\mathbf{K}) + \phi(\mathbf{x}_{t+1}, \mathbf{u}_{t+1})^\top \boldsymbol{\omega}_t \\ &\quad + b - \phi(\mathbf{x}_t, \mathbf{u}_t)^\top \boldsymbol{\omega}_t - b)] \phi(\mathbf{x}_t, \mathbf{u}_t), \end{aligned} \quad (18)$$

where  $\beta_t$  is the stepsize of the critic and  $\mathbf{K}$  denotes the policy under which the state-action pairs are sampled. Note that the constant  $b$  is not required for updating the linear coefficient  $\omega$ .

Taking the expectation of  $\omega_{t+1}$  in (18) with respect to the stationary distribution, conditioned on  $\omega_t$ , the expected subsequent critic can be written as

$$\mathbb{E}[\omega_{t+1}|\omega_t] = \omega_t + \beta_t(\mathbf{b}_K - \mathbf{A}_K\omega_t), \quad (19)$$

where

$$\begin{aligned} \mathbf{A}_K &= \mathbb{E}_{(x,u)}[\phi(x,u)(\phi(x,u) - \phi(x',u'))^\top], \\ \mathbf{b}_K &= \mathbb{E}_{(x,u)}[(c(x,u) - J(\mathbf{K}))\phi(x,u)]. \end{aligned} \quad (20)$$

Note that for ease of exposition, we denote  $(x',u')$  as the next state-action pair after  $(x,u)$  and abbreviate  $\mathbb{E}_{x \sim \rho_K, u \sim \pi_K(\cdot|x)}$  as  $\mathbb{E}_{(x,u)}$ .

**Assumption 1.** We consider the policy class  $\mathbb{K}$  such that  $\forall \mathbf{K} \in \mathbb{K}$ ,  $\mathbf{K}$  is norm bounded and the spectral radius satisfies  $\rho(\mathbf{A} - \mathbf{BK}) \leq \lambda$  for some constant  $\lambda \in (0, 1)$ .

The above assumes the uniform boundedness of the policy (actor) parameter  $\mathbf{K}$ , which is common in the literature of actor-critic algorithms [Karmakar and Bhatnagar, 2018; Barakat *et al.*, 2022; Zhou and Lu, 2023]. One potential approach to address the boundedness assumption involves formulating a projection map capable of diminishing the magnitude of  $\|\mathbf{K}\|$  when it exceeds the specified boundary [Konda and Tsitsiklis, 1999; Bhatnagar *et al.*, 2009], which is deferred to future research endeavors.

As previously discussed, a policy  $\mathbf{K}$  is considered stabilizing if and only if  $\rho(\mathbf{A} - \mathbf{BK}) < 1$ . Therefore, Assumption 1 also implies the stability of policy  $\mathbf{K}$ , which is equivalent to assuming the existence of  $\mathbf{A}_K$  due to the expectation being taken over the stationary distribution. Such assumption is standard in the literature [Wu *et al.*, 2020; Chen *et al.*, 2021; Olshevsky and Ghahserifard, 2023]. Without loss of generality, we slightly strengthen the requirement to  $\rho(\mathbf{A} - \mathbf{BK}) \leq \lambda$  for some constant  $\lambda \in (0, 1)$ . This is made to avoid tedious computation of the probability of bounded learning trajectories. It is worth noting that one could alternatively assume  $\rho(\mathbf{A} - \mathbf{BK}) < 1$  and deduce that the same results presented in the sequel with additional high probability characterization.

We then provide the coercive property of cost function  $J(\mathbf{K})$ , illustrating that  $J(\mathbf{K})$  tends towards infinity as  $\|\mathbf{K}\|$  approaches infinity or when  $\rho(\mathbf{A} - \mathbf{BK})$  approaches 1.

**Lemma 3 (Coercive Property).** *The cost function  $J(\mathbf{K})$  defined in (7) is coercive, that is, for any sequence  $\{\mathbf{K}_i\}_{i=1}^\infty$  of stabilizing policies, we have*

$$J(\mathbf{K}_i) \rightarrow +\infty, \quad \text{if } \|\mathbf{K}_i\| \rightarrow +\infty \text{ or } \rho(\mathbf{A} - \mathbf{BK}_i) \rightarrow 1.$$

Lemma 3 demonstrates the safety of boundary cutting ( $\|\mathbf{K}_i\| \rightarrow +\infty, \rho(\mathbf{A} - \mathbf{BK}_i) \rightarrow 1$ ), ensuring that the optimal  $\mathbf{K}^*$  that minimizes  $J(\mathbf{K})$  resides within the class  $\mathbb{K}$ , thereby justifying Assumption 1. Additionally, we present some numerical examples in Section 5 to support this assumption.

As the existence of  $\mathbf{A}_K$  and  $\mathbf{b}_K$  are ensured by Assumption 1, given a policy  $\pi_K$ , it is not hard to show that if the update in (19) has converged to some limiting point  $\omega_K^*$ , i.e.,  $\lim_{t \rightarrow \infty} \omega_t = \omega_K^*$ ,  $\omega_K^*$  must be the solution of  $\mathbf{A}_K\omega = \mathbf{b}_K$ .

**Lemma 4.** *Suppose  $K \in \mathbb{K}$ . Then the matrix  $\mathbf{A}_K$  defined in (20) is invertible and  $\mathbf{A}_K\omega = \mathbf{b}_K$  has a unique solution  $\omega_K^*$  that satisfies*

$$\omega_K^* = \text{svec}(\Omega_K). \quad (21)$$

where  $\Omega_K$  is defined in (15).

Since  $\text{smat}(\cdot)$  represents the inverse of  $\text{svec}(\cdot)$ , it follows that  $\Omega_K$  can be expressed as  $\text{smat}(\omega_K^*)$ , thereby completing the estimation of  $\Omega_K$ .

Combining (13), (15), and (21), we can express the natural gradient of  $J(\mathbf{K})$  using  $\omega_K^*$ :

$$\nabla_{\mathbf{K}}^N J(\mathbf{K}) = \Omega_K^{22}\mathbf{K} - \Omega_K^{21} = \text{smat}(\omega_K^*)^{22}\mathbf{K} - \text{smat}(\omega_K^*)^{21},$$

where  $\text{smat}(\omega_K^*)^{21}$  and  $\text{smat}(\omega_K^*)^{22}$  represent the bottom left corner block and bottom right corner block of matrix  $\text{smat}(\omega_K^*)$ , respectively.

This allows us to estimate the natural policy gradient using the critic parameters  $\omega_t$ , and then update the actor in a model-free manner

$$\mathbf{K}_{t+1} = \mathbf{K}_t - \alpha_t \widehat{\nabla_{\mathbf{K}_t}^N J(\mathbf{K}_t)}, \quad (22)$$

where  $\alpha_t$  is the actor stepsize and  $\widehat{\nabla_{\mathbf{K}_t}^N J(\mathbf{K}_t)}$  is the natural gradient estimation depending on  $\omega_t$ :

$$\widehat{\nabla_{\mathbf{K}_t}^N J(\mathbf{K}_t)} = \text{smat}(\omega_t)^{22}\mathbf{K}_t - \text{smat}(\omega_t)^{21}. \quad (23)$$

Furthermore, we introduce a cost estimator  $\eta_t$  to estimate the time-average cost  $J(\mathbf{K}_t)$ . Combining the critic update (18) and the actor update (22)-(23), the single-sample single-timescale AC for solving LQR is listed below.

---

**Algorithm 1** Single-Sample Single-timescale Actor-Critic for Linear Quadratic Regulator

---

- 1: **Input** initialize actor parameter  $\mathbf{K}_0 \in \mathbb{K}$ , critic parameter  $\omega_0$ , time-average cost  $\eta_0$ , stepsizes  $\alpha_t$  for actor,  $\beta_t$  for critic, and  $\gamma_t$  for cost estimator.
  - 2: **for**  $t = 0, 1, 2, \dots, T - 1$  **do**
  - 3:   Sample  $\mathbf{x}_t$  from the stationary distribution  $\rho_{\mathbf{K}_t}$ .
  - 4:   Take action  $\mathbf{u}_t \sim \pi_{\mathbf{K}_t}(\cdot|\mathbf{x}_t)$  and receive cost  $c_t = c(\mathbf{x}_t, \mathbf{u}_t)$  and the next state  $\mathbf{x}'_t$ .
  - 5:   Obtain  $\mathbf{u}'_t \sim \pi_{\mathbf{K}_t}(\cdot|\mathbf{x}'_t)$ .
  - 6:    $\delta_t = c_t - \eta_t + \phi(\mathbf{x}'_t, \mathbf{u}'_t)^\top \omega_t - \phi(\mathbf{x}_t, \mathbf{u}_t)^\top \omega_t$
  - 7:    $\eta_{t+1} = \text{proj}_{\mathcal{B}_{\bar{\eta}}}(\eta_t + \gamma_t(c_t - \eta_t))$
  - 8:    $\omega_{t+1} = \text{proj}_{\mathcal{B}_{\bar{\omega}}}(\omega_t + \beta_t \delta_t \phi(\mathbf{x}_t, \mathbf{u}_t))$
  - 9:    $\mathbf{K}_{t+1} = \mathbf{K}_t - \alpha_t(\text{smat}(\omega_t)^{22}\mathbf{K}_t - \text{smat}(\omega_t)^{21})$
  - 10: **end for**
- 

Note that *single-sample* refers to the fact that only one sample is used to update the critic per actor step. Line 3 of Algorithm 1 samples from the stationary distribution induced by the policy  $\pi_{\mathbf{K}_t}$ , which is a mild requirement in the analysis of uniformly ergodic Markov chain, such as in the LQR problem [Yang *et al.*, 2019]. It is only made to simplify the theoretical analysis. Indeed, as shown in [Tu and Recht, 2018], when  $\mathbf{K} \in \mathbb{K}$ , (9) is geometrically  $\beta$ -mixing and thus its

distribution converges to the stationary distribution exponentially. In practice, one can run the Markov chain in (9) a sufficient number of steps and sample one state from the last step to approximate the stationary distribution. In addition, *single-timescale* refers to the fact that the stepsizes for the critic and the actor updates are constantly proportional.

Since the update of the critic parameter in (18) requires the time-average cost  $J(\mathbf{K}_t)$ , Line 7 provides an estimation of it. Besides, on top of (18), we additionally introduce a projection in Line 8 and Line 9 to keep the critic norm-bounded. The projection follows the standard definition, i.e.,  $\text{proj}_{\mathcal{B}_y}(\mathbf{x})$  means project  $\mathbf{x}$  to the set  $\mathcal{B}_y := \{\mathbf{x} \mid \|\mathbf{x}\| \leq y\}$ . This is common in the literature [Wu *et al.*, 2020; Yang *et al.*, 2019; Chen and Zhao, 2022]. In our analysis, the projection is relaxed using its nonexpansive property.

## 4 Main Theory

In this section, we establish the global optimality and analyze the finite-time performance of Algorithm 1. All the proofs can be found in the Supplementary Material.

**Theorem 1.** *Suppose that Assumptions 1 hold and choose  $\alpha_t = \frac{c}{\sqrt{T}}$ ,  $\beta_t = \gamma_t = \frac{1}{\sqrt{T}}$ , where  $c$  is a small positive constant. It holds that*

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}(\eta_t - J(\mathbf{K}_t))^2 &= \mathcal{O}\left(\frac{1}{\sqrt{T}}\right), \\ \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\|\omega_t - \omega_{\mathbf{K}_t}^*\|^2 &= \mathcal{O}\left(\frac{1}{\sqrt{T}}\right), \\ \min_{0 \leq t < T} \mathbb{E}[J(\mathbf{K}_t) - J(\mathbf{K}^*)] &= \mathcal{O}\left(\frac{1}{\sqrt{T}}\right). \end{aligned}$$

The theorem shows that the cost estimator, the critic, and the actor all converge at a sub-linear rate of  $\mathcal{O}(T^{-\frac{1}{2}})$ . The  $\mathcal{O}$  notation hides the polynomials of the dependence parameters. Note that we have explicitly characterized all the necessary problem parameters in the proofs before the last step of the analysis of the interconnected system. One can easily keep all the problem parameters in the interconnected system analysis and get the order for all parameters. To focus on the key factors and for ease of comprehension, we only show the convergence rate in terms of the iteration number.

Correspondingly, to obtain an  $\epsilon$ -optimal policy, the required sample complexity is  $\mathcal{O}(\epsilon^{-2})$ . This order is consistent with the existing results on single-sample single-timescale AC [Chen *et al.*, 2021; Olshevsky and Ghahserifard, 2023; Chen and Zhao, 2022]. Nevertheless, our result is the first finite-time analysis of the single-sample single-timescale AC with a global optimality guarantee and considers the challenging continuous state-action space.

### 4.1 Proof Sketch

The main challenge in the finite-time analysis lies in that the estimation errors of the time-average cost, the critic, and the natural policy gradient are strongly coupled. To overcome this issue, we view the propagation of these errors as an interconnected system and analyze them comprehensively. To

see the merit of our analysis framework, we sketch the main proof steps of Theorem 1 in the following. The supporting lemmas and theorems mentioned below can be found in the Supplementary Material.

We define three measures  $A_T, B_T, C_T$  which denote average values of the cost estimation error, the critic error, and the square norm of natural policy gradient, respectively:

$$A_T := \frac{\sum_{t=0}^{T-1} \mathbb{E}y_t^2}{T}, B_T := \frac{\sum_{t=0}^{T-1} \mathbb{E}\|z_t\|^2}{T}, C_T := \frac{\sum_{t=0}^{T-1} \mathbb{E}\|\mathbf{E}_{\mathbf{K}_t}\|^2}{T},$$

where  $y_t := \eta_t - J(\mathbf{K}_t)$  is the cost estimation error and  $z_t := \omega_t - \omega_{\mathbf{K}_t}^*$  with  $\omega_{\mathbf{K}_t}^* := \omega_{\mathbf{K}_t}^*$  is the critic error. Note that  $\mathbf{E}_{\mathbf{K}_t} = \nabla_{\mathbf{K}_t}^N J(\mathbf{K}_t)$  is the natural policy gradient according to (13).

We first derive implicit (coupled) upper bounds for the cost estimation error  $y_t$ , the critic error  $z_t$ , and the natural gradient  $\mathbf{E}_{\mathbf{K}_t}$ , respectively. After that, we solve an interconnected system of inequalities in terms of  $A_T, B_T, C_T$  to establish the finite-time convergence.

**Step 1: Cost estimation error analysis.** From the cost estimator update rule (Line 7 of Algorithm 1), we decompose the cost estimation error into (neglecting the projection for the time being):

$$\begin{aligned} y_{t+1}^2 &= (1 - 2\gamma_t)y_t^2 + 2\gamma_t y_t (c_t - J(\mathbf{K}_t)) \\ &\quad + 2y_t (J(\mathbf{K}_t) - J(\mathbf{K}_{t+1})) \\ &\quad + [J(\mathbf{K}_t) - J(\mathbf{K}_{t+1}) + \gamma_t (c_t - \eta_t)]^2. \end{aligned} \quad (24)$$

The second term on the right hand side of (24) is a noise term introduced by random sampling of state-action pairs, which reduces to 0 after taking the expectations. The third term is the variation of the moving targets  $J(\mathbf{K}_t)$  tracked by cost estimator. It is bounded by  $y_t, z_t, \mathbf{E}_{\mathbf{K}_t}$  utilizing the Lipschitz continuity of  $J(\mathbf{K}_t)$  (Lemma 9), the actor update rule (23), and the Cauchy-Schwartz inequality. The last term reflects the variance in cost estimation, which is bounded by  $\mathcal{O}(\gamma_t)$ .

**Step 2: Critic error analysis.** By the critic update rule (Line 8 of Algorithm 1), we decompose the squared error by (neglecting the projection for the time being)

$$\begin{aligned} \|z_{t+1}\|^2 &= \|z_t\|^2 + 2\beta_t \langle z_t, \bar{h}(\omega_t, \mathbf{K}_t) \rangle + 2\beta_t \Lambda(\mathbf{O}_t, \omega_t, \mathbf{K}_t) \\ &\quad + 2\beta_t \langle z_t, \Delta h(\mathbf{O}_t, \eta_t, \mathbf{K}_t) \rangle + 2\langle z_t, \omega_t^* - \omega_{t+1}^* \rangle \\ &\quad + \|\beta_t (\mathbf{h}(\mathbf{O}_t, \omega_t, \mathbf{K}_t) + \Delta h(\mathbf{O}_t, \eta_t, \mathbf{K}_t)) \\ &\quad + (\omega_t^* - \omega_{t+1}^*)\|^2, \end{aligned} \quad (25)$$

where the definitions of  $\mathbf{h}, \bar{h}, \Delta h, \Lambda$ , and  $\mathbf{O}_t$  can be found in (28) in the Supplementary Material. The second term on the right hand side of (25) is bounded by  $-\mu \|z_t\|^2$ , where  $\mu$  is a lower bound of  $\sigma_{\min}(A_{\mathbf{K}_t})$  proved in Lemma 10. The third term is a random noise introduced by sampling, which reduces to 0 after taking expectation. The fourth term is caused by inaccurate cost and critic estimations, which can be bounded by the norm of  $y_t$  and  $z_t$ . The fifth term tracks the difference between the drifting critic targets. We control it by the Lipschitz continuity of the critic target established in Lemma 11. The last term reflects the variances of various estimations, which is bounded by  $\mathcal{O}(\beta_t)$ .

**Step 3: Natural gradient norm analysis.** From the actor update rule (Line 9 of Algorithm 1) and the almost smoothness property of LQR (Lemma 12), we derive

$$\begin{aligned} 2\text{Tr}(\mathbf{D}_{\mathbf{K}_{t+1}}\mathbf{E}_{\mathbf{K}_t}^\top\mathbf{E}_{\mathbf{K}_t}) &= \frac{1}{\alpha_t}[J(\mathbf{K}_t) - J(\mathbf{K}_{t+1})] \\ &- 2\text{Tr}(\mathbf{D}_{\mathbf{K}_{t+1}}(\hat{\mathbf{E}}_{\mathbf{K}_t} - \mathbf{E}_{\mathbf{K}_t})^\top\mathbf{E}_{\mathbf{K}_t}) \\ &+ \alpha_t\text{Tr}(\mathbf{D}_{\mathbf{K}_{t+1}}\hat{\mathbf{E}}_{\mathbf{K}_t}^\top(\mathbf{R} + \mathbf{B}^\top\mathbf{P}_{\mathbf{K}_t}\mathbf{B})\hat{\mathbf{E}}_{\mathbf{K}_t}), \end{aligned} \quad (26)$$

where  $\hat{\mathbf{E}}_{\mathbf{K}_t}$  denotes the estimation of the natural gradient  $\mathbf{E}_{\mathbf{K}_t}$ . The first term on the left hand side of (26) can be considered as the scaled square norm of the natural gradient. The first term on the right hand side compares the actor's performances between consecutive updates, which is bounded via Abel summation by parts. The second term evaluates the inaccurate natural gradient estimation, which is then bounded by the critic error  $z_t$  and the natural gradient  $\mathbf{E}_{\mathbf{K}_t}$ . The last term can be considered as the variance of the perturbed natural gradient update, which is bounded by  $\mathcal{O}(\alpha_t)$ .

**Step 4: Interconnected iteration system analysis.** Taking expectation and summing (24), (25), (26) from 0 to  $T-1$ , we obtain the following interconnected iteration system:

$$\begin{aligned} A_T &\leq \mathcal{O}\left(\frac{1}{\sqrt{T}}\right) + h_2B_T + h_2C_T, \\ B_T &\leq \mathcal{O}\left(\frac{1}{\sqrt{T}}\right) + h_4\sqrt{A_TB_T} + h_5C_T, \\ C_T &\leq \mathcal{O}\left(\frac{1}{\sqrt{T}}\right) + h_7\sqrt{B_TC_T}, \end{aligned} \quad (27)$$

where  $h_2, h_4, h_5$ , and  $h_7$  are positive constants defined in (47). By solving the above inequalities, we further prove that if  $h_2h_4^2 + h_2h_4^2h_7^2 + 2h_5h_7^2 < 1$ , then  $A_T, B_T, C_T$  converge at a rate of  $\mathcal{O}(T^{-\frac{1}{2}})$ . This condition can be easily satisfied by choosing the stepsize ratio  $c$  to be smaller than a threshold defined in (51).

**Step 5: Global convergence analysis.** To prove the global optimality, we utilize the gradient domination condition of LQR (Lemma 13):

$$J(\mathbf{K}) - J(\mathbf{K}^*) \leq \frac{1}{\sigma_{\min}(\mathbf{R})}\|\mathbf{D}_{\mathbf{K}^*}\|\text{Tr}(\mathbf{E}_{\mathbf{K}}^\top\mathbf{E}_{\mathbf{K}}).$$

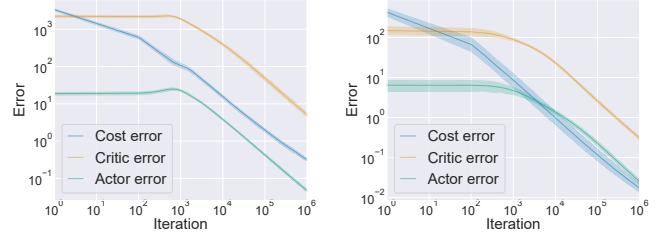
This property shows that the actor performance error can be bounded by the norm of the natural gradient ( $\text{Tr}(\mathbf{E}_{\mathbf{K}}^\top\mathbf{E}_{\mathbf{K}})$ ). Since we have proved the average natural gradient norm  $C_T$  converges to zero, summation over both sides of the above inequality yields

$$\min_{0 \leq t < T} \mathbb{E}[J(\mathbf{K}_t) - J(\mathbf{K}^*)] = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right),$$

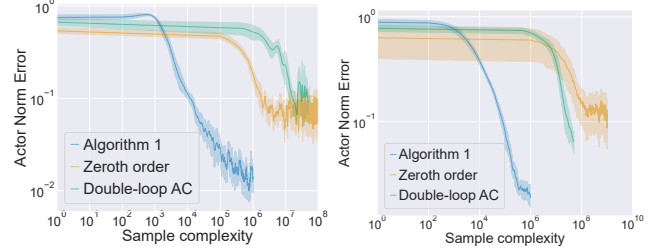
which is the convergence of the actor performance error. We thus complete the proof of Theorem 1.

## 5 Experiments

While our main contribution lies in the theoretical analysis, we also present several examples to validate the efficiency of Algorithm 1. We provide two examples to illustrate our theoretical results. The first example (first column in Figure 1)



(a) Learning results of Algorithm 1



(b) Comparison of Algorithm 1 with two other algorithms

Figure 1: (a) Learning results of Algorithm 1. In the figure, the cost error refers to  $\frac{1}{T}\sum_{t=0}^{T-1}(\eta_t - J(\mathbf{K}_t))^2$ , Critic error refers to  $\frac{1}{T}\sum_{t=0}^{T-1}\|\omega_t - \omega_{\mathbf{K}_t}^*\|^2$ , and the Actor error refers to  $\frac{1}{T}\sum_{t=0}^{T-1}[J(\mathbf{K}_t) - J(\mathbf{K}^*)]$ , corresponding to the conclusion in Theorem 1 empirically.

(b) Comparison of Algorithm 1 with two other algorithms. The actor norm error refers to  $\|\mathbf{K} - \mathbf{K}^*\|_F$ . In this figure, the solid lines correspond to the mean and the shaded regions correspond to 95% confidence interval over 10 independent runs.

is a two-dimensional system and the second example (second column in Figure 1) is a four-dimensional system. The detailed parameters are shown in Supplementary Material.

The performance of Algorithm 1 is shown in Figure 1, where the left column corresponds to the two-dimensional system and the right column to the four-dimensional system. The solid lines plot the mean values and the shaded regions denote the 95% confidence interval over 10 independent runs. Consistent with our theorem, Figure 1(a) shows that the cost estimation error, the critic error, and the actor performance error all diminish at a rate of at least  $\mathcal{O}(T^{-\frac{1}{2}})$ . The convergence also suggests that the intermediate closed-loop linear systems during iteration are uniformly stable.

We compare Algorithm 1 with the zeroth-order method [Fazel *et al.*, 2018] and the double-loop AC algorithm [Yang *et al.*, 2019] (listed in Algorithm 2 and Algorithm 3 respectively, in Supplementary Material). We plotted the relative errors of the actor parameters for all three methods in Figure 1(b). As it can be seen that Algorithm 1 demonstrates superior sample efficiency compared to the other two algorithms.

## 6 Conclusion and Discussion

In this paper, we establish the finite-time analysis for the single-sample single-timescale AC method under the LQR setting. We for the first time show that this method can find a global optimal policy under the general continuous state-action space, which contributes to understanding the limits of the AC on continuous control tasks.

## Acknowledgements

This work was supported by the Singapore Ministry of Education Tier 1 Academic Research Fund (22-5460-A0001) and Tier 2 AcRF (T2EP20123-0037). (Corresponding author: Lin Zhao.) The work of J. Duan was supported in part by the NSF China under Grant 52202487 and in part by the Young Elite Scientists Sponsorship Program by CAST under Grant 2023QNRC001.

## References

- [Anderson and Moore, 2007] Brian DO Anderson and John B Moore. *Optimal control: linear quadratic methods*. Courier Corporation, 2007.
- [Barakat *et al.*, 2022] Anas Barakat, Pascal Bianchi, and Julien Lehmann. Analysis of a target-based actor-critic algorithm with linear function approximation. In *International Conference on Artificial Intelligence and Statistics*, pages 991–1040. PMLR, 2022.
- [Bhatnagar *et al.*, 2009] Shalabh Bhatnagar, Richard S Sutton, Mohammad Ghavamzadeh, and Mark Lee. Natural actor-critic algorithms. *Automatica*, 45(11):2471–2482, 2009.
- [Bradtke *et al.*, 1994] Steven J Bradtke, B Erik Ydstie, and Andrew G Barto. Adaptive linear quadratic control using policy iteration. In *Proceedings of 1994 American Control Conference-ACC'94*, volume 3, pages 3475–3479. IEEE, 1994.
- [Castro and Meir, 2010] Dotan Di Castro and Ron Meir. A convergent online single time scale actor critic algorithm. *The Journal of Machine Learning Research*, 11:367–410, 2010.
- [Chen and Zhao, 2022] Xuyang Chen and Lin Zhao. Finite-time analysis of single-timescale actor-critic. *arXiv preprint arXiv:2210.09921*, 2022.
- [Chen *et al.*, 2021] Tianyi Chen, Yuejiao Sun, and Wotao Yin. Closing the gap: Tighter analysis of alternating stochastic gradient methods for bilevel problems. *Advances in Neural Information Processing Systems*, 34:25294–25307, 2021.
- [Cohen *et al.*, 2019] Alon Cohen, Tomer Koren, and Yishay Mansour. Learning linear-quadratic regulators efficiently with only  $\sqrt{T}$  regret. In *International Conference on Machine Learning*, pages 1300–1309. PMLR, 2019.
- [Dean *et al.*, 2018] Sarah Dean, Horia Mania, Nikolai Matni, Benjamin Recht, and Stephen Tu. Regret bounds for robust adaptive control of the linear quadratic regulator. *Advances in Neural Information Processing Systems*, 31, 2018.
- [Dean *et al.*, 2020] Sarah Dean, Horia Mania, Nikolai Matni, Benjamin Recht, and Stephen Tu. On the sample complexity of the linear quadratic regulator. *Foundations of Computational Mathematics*, 20(4):633–679, 2020.
- [Duan *et al.*, 2023] Jingliang Duan, Jie Li, Xuyang Chen, Kai Zhao, Shengbo Eben Li, and Lin Zhao. Optimization landscape of policy gradient methods for discrete-time static output feedback. *IEEE Transactions on Cybernetics*, 2023.
- [Fazel *et al.*, 2018] Maryam Fazel, Rong Ge, Sham Kakade, and Mehran Mesbahi. Global convergence of policy gradient methods for the linear quadratic regulator. In *International Conference on Machine Learning*, pages 1467–1476. PMLR, 2018.
- [Joarder and Omar, 2011] Anwar H Joarder and M Hafidz Omar. On statistical characteristics of the product of two correlated chi-square variables. *Journal of Applied Statistical Science*, 19(4):89–101, 2011.
- [Kakade, 2001] Sham M Kakade. A natural policy gradient. *Advances in neural information processing systems*, 14, 2001.
- [Karmakar and Bhatnagar, 2018] Prasenjit Karmakar and Shalabh Bhatnagar. Two time-scale stochastic approximation with controlled markov noise and off-policy temporal-difference learning. *Mathematics of Operations Research*, 43(1):130–151, 2018.
- [Konda and Tsitsiklis, 1999] Vijay Konda and John Tsitsiklis. Actor-critic algorithms. *Advances in neural information processing systems*, 12, 1999.
- [Krauth *et al.*, 2019] Karl Krauth, Stephen Tu, and Benjamin Recht. Finite-time analysis of approximate policy iteration for the linear quadratic regulator. *Advances in Neural Information Processing Systems*, 32, 2019.
- [Kumar *et al.*, 2019] Harshat Kumar, Alec Koppel, and Alejandro Ribeiro. On the sample complexity of actor-critic method for reinforcement learning with function approximation. *arXiv preprint arXiv:1910.08412*, 2019.
- [LeCun *et al.*, 2015] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- [Magnus, 1978] Jan R Magnus. The moments of products of quadratic forms in normal variables. *Statistica Neerlandica*, 32(4):201–210, 1978.
- [Malik *et al.*, 2019] Dhruv Malik, Ashwin Pananjady, Kush Bhatia, Koulik Khamaru, Peter Bartlett, and Martin Wainwright. Derivative-free methods for policy optimization: Guarantees for linear quadratic systems. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2916–2925. PMLR, 2019.
- [Mania *et al.*, 2019] Horia Mania, Stephen Tu, and Benjamin Recht. Certainty equivalence is efficient for linear quadratic control. *Advances in Neural Information Processing Systems*, 32, 2019.
- [Mnih *et al.*, 2016] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.
- [Nagar, 1959] Anirudh L Nagar. The bias and moment matrix of the general k-class estimators of the parameters in



- simultaneous equations. *Econometrica: Journal of the Econometric Society*, pages 575–595, 1959.
- [Olshevsky and Ghahserifard, 2023] Alex Olshevsky and Bahman Ghahserifard. A small gain analysis of single timescale actor critic. *SIAM Journal on Control and Optimization*, 61(2):980–1007, 2023.
- [Rencher and Schaalje, 2008] Alvin C Rencher and G Bruce Schaalje. *Linear models in statistics*. John Wiley & Sons, 2008.
- [Schacke, 2004] Kathrin Schacke. On the kronecker product. *Master's thesis, University of Waterloo*, 2004.
- [Schulman *et al.*, 2015] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR, 2015.
- [Silver *et al.*, 2017] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- [Sutton and Barto, 2018] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [Sutton *et al.*, 1999] Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.
- [Tu and Recht, 2018] Stephen Tu and Benjamin Recht. Least-squares temporal difference learning for the linear quadratic regulator. In *International Conference on Machine Learning*, pages 5005–5014. PMLR, 2018.
- [Wang *et al.*, 2019] Lingxiao Wang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural policy gradient methods: Global optimality and rates of convergence. *arXiv preprint arXiv:1909.01150*, 2019.
- [Wu *et al.*, 2020] Yue Frank Wu, Weitong Zhang, Pan Xu, and Quanquan Gu. A finite-time analysis of two time-scale actor-critic methods. *Advances in Neural Information Processing Systems*, 33:17617–17628, 2020.
- [Xu *et al.*, 2020] Tengyu Xu, Zhe Wang, and Yingbin Liang. Non-asymptotic convergence analysis of two time-scale (natural) actor-critic algorithms. *arXiv preprint arXiv:2005.03557*, 2020.
- [Yang *et al.*, 2019] Zhuoran Yang, Yongxin Chen, Mingyi Hong, and Zhaoran Wang. Provably global convergence of actor-critic: A case for linear quadratic regulator with ergodic cost. *Advances in neural information processing systems*, 32, 2019.
- [Zhang *et al.*, 2020] Shangdong Zhang, Bo Liu, Hengshuai Yao, and Shimon Whiteson. Provably convergent two-timescale off-policy actor-critic with function approximation. In *International Conference on Machine Learning*, pages 11204–11213. PMLR, 2020.
- [Zhou and Lu, 2023] Mo Zhou and Jianfeng Lu. Single timescale actor-critic method to solve the linear quadratic regulator with convergence guarantees. *Journal of Machine Learning Research*, 24(222):1–34, 2023.