

ASYNCHRONOUS ADVANTAGE ACTOR CRITIC: NON-ASYMPTOTIC ANALYSIS AND LINEAR SPEEDUP

Han Shen

Rensselaer Polytechnic Institute
shenh5@rpi.edu

Kaiqing Zhang

University of Illinois at Urbana-Champaign
kzhang66@illinois.edu

Mingyi Hong

University of Minnesota
mhong@umn.edu

Tianyi Chen

Rensselaer Polytechnic Institute
chent18@rpi.edu

ABSTRACT

Asynchronous and parallel implementation of standard reinforcement learning (RL) algorithms is a key enabler of the tremendous success of modern RL. Among many asynchronous RL algorithms, arguably the most popular and effective one is the asynchronous advantage actor-critic (A3C) algorithm. Although A3C is becoming the workhorse of RL, its theoretical properties are still not well-understood, including the non-asymptotic analysis and the performance gain of parallelism (a.k.a. speedup). This paper revisits the A3C algorithm with TD(0) for the critic update, termed A3C-TD(0), with provable convergence guarantees. With linear value function approximation for the TD update, the convergence of A3C-TD(0) is established under both i.i.d. and Markovian sampling. Under i.i.d. sampling, A3C-TD(0) obtains sample complexity of $\mathcal{O}(\epsilon^{-2.5}/N)$ per worker to achieve ϵ accuracy, where N is the number of workers. Compared to the best-known sample complexity of $\mathcal{O}(\epsilon^{-2.5})$ for two-timescale AC, A3C-TD(0) achieves *linear speedup*, which justifies the advantage of parallelism and asynchrony in AC algorithms theoretically for the first time. Numerical tests on synthetically generated instances and OpenAI Gym environments have been provided to verify our theoretical analysis.

1 INTRODUCTION

Reinforcement learning (RL) has achieved impressive performance in many domains such as robotics [1, 2] and video games [3]. However, these empirical successes are often at the expense of significant computation. To unlock high computation capabilities, the state-of-the-art RL approaches rely on sampling data from massive parallel simulators on multiple machines [3, 4, 5, 6]. Empirically, these approaches can stabilize the learning processes and *reduce training time* when they are implemented in an asynchronous manner. One popular RL method that often achieves the best empirical performance is the asynchronous variant of the actor-critic (AC) algorithm, which is referred to as A3C [3].

A3C builds on the original AC algorithm [7]. At a high level, AC simultaneously performs policy optimization (a.k.a. the actor step) using the policy gradient method [8] and policy evaluation (a.k.a. the critic step) using the temporal difference learning (TD) algorithm [9]. To ensure scalability, both actor and critic steps can combine with various function approximation techniques. To ensure stability, AC is often implemented in a two time-scale fashion, where the actor step runs in the slow timescale and the critic step runs in the fast timescale. Similar to other on-policy RL algorithms, AC uses samples generated from the target policy. Thus, data sampling is entangled with the learning procedure, which generates significant *overhead*. To speed up the sampling process of AC, A3C introduces multiple workers with a shared policy, and each learner has its own simulator to perform data sampling. The shared policy can be then updated using samples collected from multiple learners.

Despite the tremendous empirical success achieved by A3C, to the best of our knowledge, its theoretical property is not well-understood. The following *theoretical* questions remain unclear: **Q1**) Under what assumption does A3C converge? **Q2**) What is its convergence rate? **Q3**) Can A3C obtain benefit (or linear speedup) using parallelism and asynchrony?

For **Q3**, we are interested in the *training time linear speedup* with N workers, which is the ratio between the training time using a single worker and that using N workers. Since asynchronous parallelism mitigates the effect of stragglers and keeps all workers busy, the training time speedup can be measured roughly by the sample (i.e., computational) complexity linear speedup [10], given by

$$\text{Speedup}(N) = \frac{\text{sample complexity when using one worker}}{\text{average sample complexity per worker when using } N \text{ workers}}. \quad (1)$$

If $\text{Speedup}(N) = \Theta(N)$, the speedup is linear, and the training time roughly reduces linearly as the number of workers increases. This paper aims to answer these questions, towards the goal of providing theoretical justification for the empirical successes of parallel and asynchronous RL.

1.1 RELATED WORKS

Analysis of actor critic algorithms. AC method was first proposed by [7, 11], with asymptotic convergence guarantees provided in [7, 11, 12]. It was not until recently that the *non-asymptotic* analyses of AC have been established. The finite-sample guarantee for the batch AC algorithm has been established in [13, 14, 15] with i.i.d. sampling. Later, in [16], the finite-sample analysis was established for the double-loop nested AC algorithm under the Markovian setting. An improved analysis for the Markovian setting with minibatch updates has been presented in [17] for the nested AC method. More recently, [18, 19] have provided the first finite-time analyses for the two-timescale AC algorithms under Markov sampling, with both $\tilde{O}(\epsilon^{-2.5})$ sample complexity, which is the best-known sample complexity for two-timescale AC. Through the lens of bi-level optimization, [20] has also provided finite-sample guarantees for this two-timescale Markov sampling setting, with global optimality guarantees when a *natural* policy gradient step is used in the actor. However, none of the existing works has analyzed the effect of the asynchronous and parallel updates in AC.

Empirical parallel and distributed AC. In [3], the original A3C method was proposed and became the workhorse in empirical RL. Later, [21] has provided a GPU-version of A3C which significantly decreases training time. Recently, the A3C algorithm is further optimized in modern computers by [22], where a large batch variant of A3C with improved efficiency is also proposed. In [23], an importance weighted distributed AC algorithm IMPALA has been developed to solve a collection of problems with one single set of parameters. Recently, a gossip-based distributed yet synchronous AC algorithm has been proposed in [5], which has achieved the performance competitive to A3C.

Asynchronous stochastic optimization. For solving general optimization problems, asynchronous stochastic methods have received much attention recently. The study of asynchronous stochastic methods can be traced back to 1980s [24]. With the batch size M , [25] analyzed asynchronous SGD (async-SGD) for convex functions, and derived a convergence rate of $\mathcal{O}(K^{-\frac{1}{2}} M^{-\frac{1}{2}})$ if delay K_0 is bounded by $\mathcal{O}(K^{\frac{1}{4}} M^{-\frac{3}{4}})$. This result implies linear speedup. [26] extended the analysis of [25] to smooth convex with nonsmooth regularization and derived a similar rate. Recent studies by [27] improved upper bound of K_0 to $\mathcal{O}(K^{\frac{1}{2}} M^{-\frac{1}{2}})$. However, all these works have focused on the single-timescale SGD with a single variable, which cannot capture the stochastic recursion of the AC and A3C algorithms. To best of our knowledge, non-asymptotic analysis of asynchronous two-timescale SGD has remained unaddressed, and its speedup analysis is even an uncharted territory.

1.2 THIS WORK

In this context, we revisit A3C with TD(0) for the critic update, termed A3C-TD(0). The hope is to provide *non-asymptotic* guarantee and *linear speedup* justification for this popular algorithm.

Our contributions. Compared to the existing literature on both the AC algorithms and the async-SGD, our contributions can be summarized as follows.

c1) We revisit two-timescale A3C-TD(0) and establish its convergence rates with both i.i.d. and Markovian sampling. To the best of our knowledge, this is the first non-asymptotic convergence result for *asynchronous parallel* AC algorithms.

c2) We characterize the sample complexity of A3C-TD(0). In i.i.d. setting, A3C-TD(0) achieves a sample complexity of $\mathcal{O}(\epsilon^{-2.5}/N)$ per worker, where N is the number of workers. Compared to the best-known complexity of $\mathcal{O}(\epsilon^{-2.5})$ for i.i.d. two-timescale AC [20], A3C-TD(0) achieves *linear*

speedup, thanks to the parallelism and asynchrony. In the Markovian setting, if delay is bounded, the sample complexity of A3C-TD(0) matches the order of the non-parallel AC algorithm [19].

c3) We test A3C-TD(0) on the synthetically generated environment to verify our theoretical guarantees with both i.i.d. and Markovian sampling. We also test A3C-TD(0) on the classic control tasks and Atari Games from OpenAI Gym.

Technical challenges. Compared to the recent analysis of nonparallel two-timescale AC in [18, 19, 20], several new challenges arise due to the parallelism and asynchrony.

Markovian noise coupled with asynchrony and delay. The analysis of two-timescale AC algorithm is non-trivial because of the Markovian noise coupled with both the actor and critic steps. Different from the nonparallel AC that only involves a single Markov chain, asynchronous parallel AC introduces multiple Markov chains (one per worker) that mix at different speed. This is because at a given iteration, workers collect different number of samples and thus their chains mix to different degrees. As we will show later, the worker with the slowest mixing chain will determine the convergence.

Linear speedup for SGD with two coupled sequences. Parallel async-SGD has been shown to achieve linear speedup recently [10, 28]. Different from async-SGD, asynchronous AC is a two-timescale stochastic *semi-gradient* algorithm for solving the more challenging *bilevel* optimization problem (see [20]). The errors induced by asynchrony and delay are intertwined with both actor and critic updates via a nested structure, which makes the sharp analysis more challenging. Our linear speedup analysis should be also distinguished from that of mini-batch async-SGD [29], where the speedup is a result of *variance reduction* thanks to the larger batch size generated by parallel workers.

2 PRELIMINARIES

2.1 MARKOV DECISION PROCESS AND POLICY GRADIENT THEOREM

RL problems are often modeled as an MDP described by $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{P}, R, \gamma\}$, where \mathcal{S} is the state space, \mathcal{A} is the action space, $\mathcal{P}(s'|s, a)$ is the probability of transitioning to $s' \in \mathcal{S}$ given current state $s \in \mathcal{S}$ and action $a \in \mathcal{A}$, and $R(s, a, s')$ is the reward associated with the transition (s, a, s') , and $\gamma \in [0, 1)$ is a discount factor. Throughout the paper, we assume the reward R is upper-bounded by a constant R_{\max} . We also define the normalized reward as $r(s, a, s') := (1 - \gamma)R(s, a, s')$ which is upper-bounded by $r_{\max} = (1 - \gamma)R_{\max}$. A policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ is defined as a mapping from the state space \mathcal{S} to the probability distribution over the action space \mathcal{A} .

Considering discrete time t in an infinite horizon, a policy π can generate a trajectory of state-action pairs $(s_0, a_0, s_1, a_1, \dots)$ with $a_t \sim \pi(\cdot|s_t)$ and $s_{t+1} \sim \mathcal{P}(\cdot|s_t, a_t)$. Given a policy π , we define the normalized state and state action value functions as

$$V_\pi(s) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, s_{t+1}) \mid s_0 = s \right], \quad Q_\pi(s, a) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, s_{t+1}) \mid s_0 = s, a_0 = a \right] \quad (2)$$

where \mathbb{E} is taken over the trajectory $(s_0, a_0, s_1, a_1, \dots)$ generated under policy π . With the above definitions, the normalized advantage function is $A_\pi(s, a) := Q_\pi(s, a) - V_\pi(s)$. With η denoting the initial state distribution, the discounted state visitation measure induced by policy π is defined as $d_\pi(s) := (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s \mid s_0 \sim \eta, \pi)$, and the discounted state action visitation measure is $d'_\pi(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s \mid s_0 \sim \eta, \pi) \pi(a|s)$.

The goal of RL is to find a policy that maximizes the expected accumulative reward $J(\pi) := \mathbb{E}_{s \sim \eta}[V_\pi(s)]$. When the state and action spaces are large, finding the optimal policy π becomes computationally intractable. To overcome the inherent difficulty of learning a function, the policy gradient methods search the best performing policy over a class of parameterized policies. We parameterize the policy with parameter $\theta \in \mathbb{R}^d$, and solve the optimization problem as

$$\max_{\theta \in \mathbb{R}^d} J(\theta) \quad \text{with} \quad J(\theta) := \mathbb{E}_{s \sim \eta} [V_{\pi_\theta}(s)]. \quad (3)$$

To maximize $J(\theta)$ with respect to θ , one can update θ using the policy gradient direction given by [8]

$$\nabla J(\theta) = \mathbb{E}_{s, a \sim d'_\theta} [A_{\pi_\theta}(s, a) \psi_\theta(s, a)], \quad (4)$$

where $\psi_\theta(s, a) := \nabla \log \pi_\theta(a|s)$, and $d'_\theta := (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s \mid s_0, \pi_\theta) \pi_\theta(a|s)$. Since computing \mathbb{E} in (4) is expensive if not impossible, popular policy gradient-based algorithms iteratively update θ using stochastic estimate of (4) such as REINFORCE [30] and G(PO)MDP [31].

2.2 ACTOR-CRITIC ALGORITHM WITH VALUE FUNCTION APPROXIMATION

Both REINFORCE and G(PO)MDP-based policy gradient algorithms rely on a Monte-Carlo estimate of the value function $V_{\pi_\theta}(s)$ and thus $\nabla J(\theta)$ by generating a trajectory per iteration. However, policy gradient methods based on Monte-Carlo estimate typically suffer from high variance and large sampling cost. An alternative way is to recursively refine the estimate of $V_{\pi_\theta}(s)$. For a policy π_θ , it is known that $V_{\pi_\theta}(s)$ satisfies the Bellman equation [32], that is

$$V_{\pi_\theta}(s) = \mathbb{E}_{a \sim \pi_\theta(\cdot|s), s' \sim \mathcal{P}(\cdot|s, a)} [r(s, a, s') + \gamma V_{\pi_\theta}(s')], \quad \forall s \in \mathcal{S}. \quad (5)$$

In practice, when the state space \mathcal{S} is prohibitively large, one cannot afford the computational and memory complexity of computing $V_{\pi_\theta}(s)$ and $A_{\pi_\theta}(s, a)$. To overcome this curse-of-dimensionality, a popular method is to approximate the value function using function approximation techniques. Given the state feature mapping $\phi(\cdot) : \mathcal{S} \rightarrow \mathbb{R}^{d'}$ for some $d' > 0$, we approximate the value function linearly as $V_{\pi_\theta}(s) \approx \hat{V}_\omega(s) := \phi(s)^\top \omega$, where $\omega \in \mathbb{R}^{d'}$ is the critic parameter.

Given a policy π_θ , the task of finding the best ω such that $V_{\pi_\theta}(s) \approx \hat{V}_\omega(s)$ is usually addressed by TD learning [9]. Defining the k th transition as $x_k := (s_k, a_k, s_{k+1})$ and the corresponding TD-error as $\hat{\delta}(x_k, \omega_k) := r(s_k, a_k, s_{k+1}) + \gamma \phi(s_{k+1})^\top \omega_k - \phi(s_k)^\top \omega_k$, the parameter ω is updated via

$$\omega_{k+1} = \Pi_{R_\omega}(\omega_k + \beta_k g(x_k, \omega_k)) \quad \text{with} \quad g(x_k, \omega_k) := \hat{\delta}(x_k, \omega_k) \nabla_{\omega_k} \hat{V}_{\omega_k}(s_k) \quad (6)$$

where β_k is the critic stepsize, and Π_{R_ω} is the projection with R_ω being a pre-defined constant. The projection step is often used to control the norm of the gradient. In AC, it prevents the actor and critic updates from going a too large step in the ‘wrong’ direction; see e.g., [7, 18, 19].

Using the definition of advantage function $A_{\pi_\theta}(s, a) = \mathbb{E}_{s' \sim \mathcal{P}}[r(s, a, s') + \gamma V_{\pi_\theta}(s')] - V_{\pi_\theta}(s)$, we can also rewrite (4) as $\nabla J(\theta) = \mathbb{E}_{s, a \sim d'_\theta, s' \sim \mathcal{P}} [(r(s, a, s') + \gamma V_{\pi_\theta}(s') - V_{\pi_\theta}(s)) \psi_\theta(s, a)]$. Leveraging the value function approximation, we can then approximate the policy gradient as

$$\hat{\nabla} J(\theta) = (r(s, a, s') + \gamma \hat{V}_\omega(s') - \hat{V}_\omega(s)) \psi_\theta(s, a) = \hat{\delta}(x, \omega) \psi_\theta(s, a) \quad (7)$$

which gives rise to the policy update

$$\theta_{k+1} = \theta_k + \alpha_k v(x_k, \theta_k, \omega_k) \quad \text{with} \quad v(x_k, \theta_k, \omega_k) := \hat{\delta}(x_k, \omega_k) \psi_{\theta_k}(s_k, a_k) \quad (8)$$

where α_k is the stepsize for the actor update.

To ensure convergence when simultaneously performing critic and actor updates, the stepsizes α_k and β_k often decay at two different rates, which is referred to the two-timescale AC [19, 20].

3 ASYNCHRONOUS ADVANTAGE ACTOR CRITIC WITH TD(0)

To speed up the training process, we implement AC over N workers in a shared memory setting without coordinating among workers — a setting similar to that in A3C [3]. Each worker has its own simulator to perform sampling, and then collaboratively updates the shared policy π_θ using AC updates. As there is no synchronization after each update, the policy used by workers to generate samples may be outdated, which introduces staleness.

Notations on transition (s, a, s') . Since each worker will maintain a separate Markov chain, we thereafter use subscription t in (s_t, a_t, s_{t+1}) to indicate the t th transition on a Markov chain. We use k to denote the global counter (or iteration), which increases by one whenever a worker finishes the actor and critic updates in the shared memory. We use subscription (k) in $(s_{(k)}, a_{(k)}, s'_{(k)})$ to indicate the transition used in the k th update.

Specifically, we initialize θ_0, ω_0 in the shared memory. Each worker will initialize the simulator with initial state s_0 . Without coordination, workers will read θ, ω in the shared memory. From each worker’s view, it then generates sample (s_t, a_t, s_{t+1}) by either sampling s_t from $\mu_\theta(\cdot)$, where $\mu_\theta(\cdot)$ is

Algorithm 1 Asynchronous advantage AC with TD(0): each worker’s view.

```

1: Global initialize: Global counter  $k = 0$ , initial  $\theta_0, \omega_0$  in the shared memory.
2: Worker initialize: Local counter  $t = 0$ . Obtain initial state  $s_0$ .
3: for  $t = 0, 1, 2, \dots$  do
4:   Read  $\theta, \omega$  in the shared memory.
5:   Option 1 (i.i.d. sampling):
6:     Sample  $s_t \sim \mu_\theta(\cdot), a_t \sim \pi_\theta(\cdot|s), s_{t+1} \sim \mathcal{P}(\cdot|s_t, a_t)$ .
7:   Option 2 (Markovian sampling):
8:     Sample  $a_t \sim \pi_\theta(\cdot|s_t), s_{t+1} \sim \mathcal{P}(\cdot|s_t, a_t)$ .
9:   Compute  $\hat{\delta}(x_t, \omega) = r(s_t, a_t, s_{t+1}) + \gamma \hat{V}_\omega(s_{t+1}) - \hat{V}_\omega(s_t)$ .
10:  Compute  $g(x_t, \omega) = \hat{\delta}(x_t, \omega) \nabla_\omega \hat{V}_\omega(s_t)$ .
11:  Compute  $v(x_t, \theta, \omega) = \hat{\delta}(x_t, \omega) \psi_\theta(s_t, a_t)$ .
12:  In the shared memory, perform update (9).
13: end for

```

the stationary distribution of a MDP with transition probability measure $\mathcal{P}(\cdot|s_t, a_t)$ and policy π_θ , or, sampling s_t from a Markov chain under policy π_θ . In both cases, each worker obtains $a_t \sim \pi_\theta(\cdot|s_t)$ and $s_{t+1} \sim \mathcal{P}(\cdot|s_t, a_t)$. Once obtaining $x_t := (s_t, a_t, s_{t+1})$, each worker locally computes the policy gradient $v(x_t, \theta, \omega)$ and the TD(0) update $g(x_t, \omega)$, and then updates the parameters in shared memory asynchronously by

$$\omega_{k+1} = \Pi_{R_\omega} (\omega_k + \beta_k g(x_{(k)}, \omega_{k-\tau_k})), \quad (9a)$$

$$\theta_{k+1} = \theta_k + \alpha_k v(x_{(k)}, \theta_{k-\tau_k}, \omega_{k-\tau_k}), \quad (9b)$$

where τ_k is the delay in the k th actor and critic updates. See the A3C with TD(0) in Algorithm 1.

Parallel sampling. The AC update (6) and (8) uses samples generated “on-the-fly” from the target policy π_θ , which brings overhead. Compared with (6) and (8), the A3C-TD(0) update (9) allows parallel sampling from N workers, which is the key to linear speedup. We consider the case where only one worker can update parameters in the shared memory at the same time and the update cannot be interrupted. In practice, (9) can also be performed in a mini-batch fashion.

Minor differences from A3C [3]. The A3C-TD(0) algorithm resembles the popular A3C method [3]. With n_{\max} denoting the horizon of steps, for $n \in \{1, \dots, n_{\max}\}$, A3C iteratively uses n -step TD errors to compute actor and critic gradients. In A3C-TD(0), we use the TD(0) method which is the 1-step TD method for actor and critic update. When $n_{\max} = 1$, A3C method reduces to A3C-TD(0). The n -step TD method is a hybrid version of the TD(0) method and the Monte-Carlo method. The A3C method with Monte-Carlo sampling is essentially the delayed policy gradient method, and thus its convergence follows directly from the delayed SGD. Therefore, we believe that the convergence of the A3C method based on TD(0) in this paper can be easily extended to the convergence of the A3C method with n -step TD. We here focus on A3C with TD(0) just for ease of exposition.

4 CONVERGENCE ANALYSIS OF TWO-TIMESCALE A3C-TD(0)

In this section, we analyze the convergence of A3C-TD(0) in both i.i.d. and Markovian settings. Throughout this section, the notation $\mathcal{O}(\cdot)$ contains constants that are independent of N and K_0 .

To analyze the performance of A3C-TD(0), we make the following assumptions.

Assumption 1. *There exists K_0 such that the delay at each iteration is bounded by $\tau_k \leq K_0, \forall k$.*

Assumption 1 ensures the viability of analyzing the asynchronous update; see the same assumption in e.g., [5, 27]. In practice, the delay usually scales as the number of workers, that is $K_0 = \Theta(N)$.

Recall the state feature mapping $\phi(\cdot) : \mathcal{S} \rightarrow \mathbb{R}^{d'}$. We define

$$A_{\theta, \phi} := \mathbb{E}_{s \sim \mu_\theta, s' \sim \mathcal{P}_{\pi_\theta}} [\phi(s)(\gamma \phi(s') - \phi(s))^\top], \quad b_{\theta, \phi} := \mathbb{E}_{s \sim \mu_\theta, a \sim \pi_\theta, s' \sim \mathcal{P}} [r(s, a, s') \phi(s)]. \quad (10)$$

It is known that for a given θ , the stationary point ω_θ^* of the TD(0) update in Algorithm 1 satisfies

$$A_{\theta, \phi} \omega_\theta^* + b_{\theta, \phi} = 0. \quad (11)$$

Assumption 2. For all $s \in \mathcal{S}$, the feature vector $\phi(s)$ is normalized so that $\|\phi(s)\|_2 \leq 1$. For $\theta \in \mathbb{R}^d$, $A_{\theta, \phi}$ is negative definite and its maximum eigenvalue is upper bounded by constant $-\lambda$.

Assumption 2 is common in analyzing TD with linear function approximation; see e.g., [19, 33, 34, 35]. With this assumption, $A_{\theta, \phi}$ is invertible, so we have $\omega_\theta^* = -A_{\theta, \phi}^{-1} b_{\theta, \phi}$. Defining $R_\omega := r_{\max}/\lambda$, we have $\|\omega_\theta^*\|_2 \leq R_\omega$. It justifies the projection introduced in Algorithm 1. In practice, the projection radius R_ω can be estimated online by methods proposed in [33, Section 8.2] or [36, Lemma 1].

Assumption 3. For any $\theta, \theta' \in \mathbb{R}^d$, $s \in \mathcal{S}$ and $a \in \mathcal{A}$, there exist constants such that: i) $\|\psi_\theta(s, a)\|_2 \leq C_\psi$; ii) $\|\psi_\theta(s, a) - \psi_{\theta'}(s, a)\|_2 \leq L_\psi \|\theta - \theta'\|_2$; iii) $|\pi_\theta(a|s) - \pi_{\theta'}(a|s)| \leq L_\pi \|\theta - \theta'\|_2$.

Assumption 3 is common in analyzing policy gradient-type algorithms which has also been made by e.g., [36, 37, 38]. This assumption holds for many policy parameterization methods such as tabular softmax policy [38], Gaussian policy [39] and Boltzmann policy [40].

Assumption 4. For any $\theta \in \mathbb{R}^d$, the Markov chain under policy π_θ and transition kernel $\mathcal{P}(\cdot|s, a)$ is irreducible and aperiodic. Then there exist constants $\kappa > 0$ and $\rho \in (0, 1)$ such that

$$\sup_{s \in \mathcal{S}} d_{TV}(\mathbb{P}(s_t \in \cdot | s_0 = s, \pi_\theta), \mu_\theta) \leq \kappa \rho^t, \quad \forall t \quad (12)$$

where μ_θ is the stationary state distribution under π_θ , and s_t is the state of Markov chain at time t .

Assumption 4 assumes the Markov chain mixes at a geometric rate; see also [33, 35, 41].

4.1 LINEAR SPEEDUP RESULT WITH I.I.D. SAMPLING

In this section, we consider A3C-TD(0) under the i.i.d. sampling setting, which is widely used for analyzing RL algorithms; see e.g., [14, 20, 42].

We first give the convergence result of critic update as follows.

Theorem 1 (Critic convergence). *Suppose Assumptions 1–4 hold. Consider Algorithm 1 with i.i.d. sampling and $\hat{V}_\omega(s) = \phi(s)^\top \omega$. Select step size $\alpha_k = \frac{c_1}{(1+k)^{\sigma_1}}$, $\beta_k = \frac{c_2}{(1+k)^{\sigma_2}}$, where $0 < \sigma_2 < \sigma_1 < 1$ and c_1, c_2 are positive constants. Then it holds that*

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\omega_k - \omega_{\theta_k}^*\|_2^2 = \mathcal{O}\left(\frac{1}{K^{1-\sigma_2}}\right) + \mathcal{O}\left(\frac{1}{K^{2(\sigma_1-\sigma_2)}}\right) + \mathcal{O}\left(\frac{K_0^2}{K^{2\sigma_2}}\right) + \mathcal{O}\left(\frac{K_0}{K^{\sigma_1}}\right) + \mathcal{O}\left(\frac{1}{K^{\sigma_2}}\right). \quad (13)$$

Different from async-SGD (e.g., [10]), the optimal critic parameter ω_θ^* is constantly drifting as θ changes at each iteration. This necessitates setting $\sigma_1 > \sigma_2$ to make the policy change slower than the critic, which can be observed in the second term in (13). If $\sigma_1 > \sigma_2$, then the policy is static relative to the critic in an asymptotic sense.

To introduce the convergence of actor update, we first define the critic approximation error as

$$\epsilon_{app} := \max_{\theta \in \mathbb{R}^d} \sqrt{\mathbb{E}_{s \sim \mu_\theta} |V_{\pi_\theta}(s) - \hat{V}_{\omega_\theta^*}(s)|^2}, \quad (14)$$

where μ_θ is the stationary distribution under π_θ and \mathcal{P} . This error captures the quality of the critic function approximation; see also [16, 17, 19]. It becomes zero in the tabular case where the value function V_{π_θ} belongs to the linear function space for any θ .

Because A3C samples from the transition kernel \mathcal{P} , the marginal state distribution for each worker will converge to the stationary distribution μ_θ by Assumption 4. However, the policy gradient in (4) requires samples from the discounted visitation measure d_θ' . The mismatch between the visitation measure and the stationary distribution will inevitably introduce the sampling error defined as:

$$\epsilon_{sp} := 4R_{\max}^2 C_\psi^2 \left(\log_\rho \kappa^{-1} + \frac{1}{1-\rho} \right) (1-\gamma). \quad (15)$$

The sampling error is small when γ is close to 1. This is because the discounted visitation measure approaches to the stationary distribution when γ approaches 1. This fact is also consistent with the practice where a large γ is commonly used in two-timescale AC algorithms [3].

Theorem 2 (Actor convergence). *Under the same assumptions of Theorem 1, select step size $\alpha_k = \frac{c_1}{(1+k)^{\sigma_1}}$, $\beta_k = \frac{c_2}{(1+k)^{\sigma_2}}$, where $0 < \sigma_2 < \sigma_1 < 1$ and c_1, c_2 are positive constants. Then it holds that*

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 = \mathcal{O}\left(\frac{1}{K^{1-\sigma_1}}\right) + \mathcal{O}\left(\frac{K_0}{K^{\sigma_1}}\right) + \mathcal{O}\left(\frac{K_0^2}{K^{2\sigma_2}}\right) + \mathcal{O}\left(\frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\omega_k - \omega_{\theta_k}^*\|_2^2\right) + \mathcal{O}(\epsilon_{app} + \epsilon_{sp}). \quad (16)$$

Different from the analysis of async-SGD, in actor update, the stochastic gradient $v(x, \theta, \omega)$ is biased because of inexact value function approximation. The bias introduced by the critic optimality gap and the function approximation error correspond to the last two terms in (16).

Select $\sigma_1 = \frac{3}{5}$ and $\sigma_2 = \frac{2}{5}$ in Theorems 1 and 2. If $N = \mathcal{O}(K^{\frac{1}{5}})$, then it holds that

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 = \mathcal{O}(K^{-\frac{2}{5}}) + \mathcal{O}(\epsilon_{app} + \epsilon_{sp}) \quad (17)$$

where $\mathcal{O}(\cdot)$ contains constants that are independent of N and K_0 .

Corollary 1 (Linear speedup). *To reach ϵ -accuracy in (17), the required number of iterations is $\mathcal{O}(\epsilon^{-2.5})$. Since each iteration of A3C-TD(0) only uses one sample (one transition), the sample complexity is $\mathcal{O}(\epsilon^{-2.5})$, which matches the state-of-the-art sample complexity of two-timescale AC running on one worker. Then under A3C-TD(0), the average sample complexity per worker is $\mathcal{O}(\epsilon^{-2.5}/N)$ which indicates linear speedup in (1). Intuitively, the negative effect of parameter staleness introduced by parallel asynchrony vanishes asymptotically, which implies linear speedup.*

4.2 CONVERGENCE RESULT WITH MARKOVIAN SAMPLING

Theorem 3 (Critic convergence). *Suppose Assumptions 1–4 hold. Consider Algorithm 1 with Markovian sampling and $\hat{V}_\omega(s) = \phi(s)^\top \omega$. Select step size $\alpha_k = \frac{c_1}{(1+k)^{\sigma_1}}$, $\beta_k = \frac{c_2}{(1+k)^{\sigma_2}}$, where $0 < \sigma_2 < \sigma_1 < 1$ and c_1, c_2 are positive constants. Then it holds that*

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\omega_k - \omega_{\theta_k}^*\|_2^2 = \mathcal{O}\left(\frac{1}{K^{1-\sigma_2}}\right) + \mathcal{O}\left(\frac{1}{K^{2(\sigma_1-\sigma_2)}}\right) + \mathcal{O}\left(\frac{K_0^2}{K^{2\sigma_2}}\right) + \mathcal{O}\left(\frac{K_0^2 \log^2 K}{K^{\sigma_1}}\right) + \mathcal{O}\left(\frac{K_0 \log K}{K^{\sigma_2}}\right). \quad (18)$$

The following theorem gives the convergence rate of actor update in Algorithm 1.

Theorem 4 (Actor convergence). *Under the same assumptions of Theorem 3, select step size $\alpha_k = \frac{c_1}{(1+k)^{\sigma_1}}$, $\beta_k = \frac{c_2}{(1+k)^{\sigma_2}}$, where $0 < \sigma_2 < \sigma_1 < 1$ and c_1, c_2 are positive constants. Then it holds that*

$$\begin{aligned} & \frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 \\ &= \mathcal{O}\left(\frac{1}{K^{1-\sigma_1}}\right) + \mathcal{O}\left(\frac{K_0^2 \log^2 K}{K^{\sigma_1}}\right) + \mathcal{O}\left(\frac{K_0^2}{K^{2\sigma_2}}\right) + \mathcal{O}\left(\frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\omega_k - \omega_{\theta_k}^*\|_2^2\right) + \mathcal{O}(\epsilon_{app} + \epsilon_{sp}). \end{aligned} \quad (19)$$

Assume $K_0 = \mathcal{O}(K^{\frac{1}{5}})$. Given Theorem 3, select $\sigma_1 = \frac{3}{5}$ and $\sigma_2 = \frac{2}{5}$, then it holds that

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 = \tilde{\mathcal{O}}\left(K_0 K^{-\frac{2}{5}}\right) + \mathcal{O}(\epsilon_{app} + \epsilon_{sp}), \quad (20)$$

where $\tilde{\mathcal{O}}(\cdot)$ hides constants and the logarithmic order of K .

Different from i.i.d. sampling, the stochastic gradients $g(x, \omega)$ and $v(x, \theta, \omega)$ are biased for Markovian sampling, and the bias decreases as the chain mixes. The mixing time corresponds to the logarithmic terms $\log K$ in (18) and (19). Because of asynchrony, at a given iteration, workers collect different number of samples and their chains mix to different degrees. The worker with the slowest mixing chain will determine the rate of convergence. The product of K_0 and $\log K$ in (18) and (19) appears due to the slowest mixing chain. As the last term in (18) dominates other terms asymptotically, the convergence rate degrades as the number of workers increases. While the theoretical linear speedup is difficult to establish in the Markovian setting, we will empirically demonstrate it in Section 5.2.

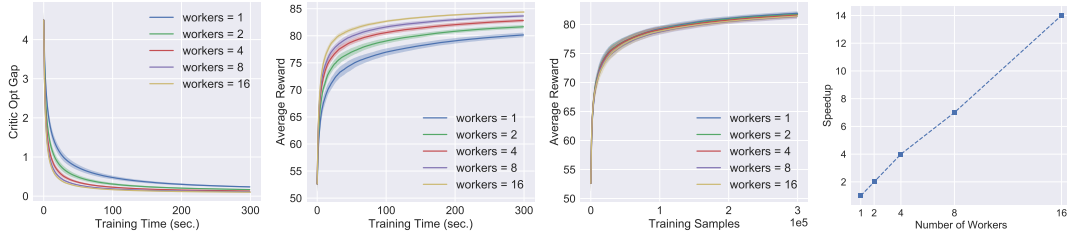


Figure 1: Convergence results of A3C-TD(0) with i.i.d. sampling in synthetic environment.

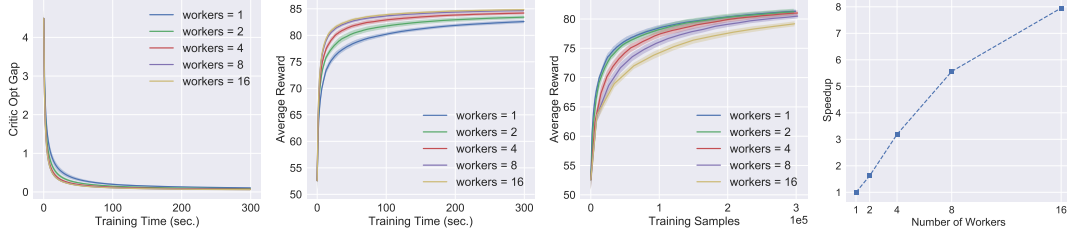


Figure 2: Convergence results of A3C-TD(0) with Markovian sampling in synthetic environment.

5 NUMERICAL EXPERIMENTS

We test the speedup performance of A3C-TD(0) on both synthetically generated and Gym environments. The settings, parameters, and codes are provided in supplementary material.

5.1 A3C-TD(0) IN SYNTHETIC ENVIRONMENT

To verify the theoretical result, we tested A3C-TD(0) with linear value function approximation in a synthetic environment. We use tabular softmax policy parameterization [38], which satisfies Assumption 3. The MDP has a state space $|\mathcal{S}| = 100$, an discrete action space of $|\mathcal{A}| = 5$. Each state feature has a dimension of 10. Elements of the transition matrix, the reward and the state features are randomly sampled from a uniform distribution over $(0, 1)$. We evaluate the convergence of actor and critic respectively with the running average of test reward and critic optimality gap $\|\omega_k - \omega_k^*\|_2$.

Figures 1 and 2 show the training time and sample complexity of running A3C-TD(0) with i.i.d. sampling and Markovian sampling respectively. The speedup plot is measured by the number of samples needed to achieve a target running average reward under different number of workers. All the results are average over 10 Monte-Carlo runs. Figure 1 shows that the sample complexity of A3C-TD(0) stays about the same with different number of workers under i.i.d. sampling. Also, it can be observed from the speedup plot of Figure 1 that the A3C-TD(0) achieves roughly linear speedup with i.i.d. sampling, which is consistent with Corollary 1. The speedup of A3C-TD(0) with Markovian sampling shown in Figure 2 is roughly linear when number of workers is small.

5.2 A3C-TD(0) IN OPENAI GYM ENVIRONMENTS

We have also tested A3C-TD(0) with neural network parametrization in the classic control (Carpole) environment and the Atari game (Seaquest and Beamrider) environments. In Figures 3-5, each curve is generated by averaging over 5 Monte-Carlo runs with 95% confidence interval. Figures 3-5 show the speedup of A3C-TD(0) under different number of workers, where the average reward is computed by taking the running average of test rewards. The speedup and runtime speedup plots are respectively measured by the number of samples and training time needed to achieve a target running average reward under different number of workers. Although not justified theoretically, Figures 3-5 suggest that the sample complexity speedup is roughly linear, and the runtime speedup slightly degrades when the number of workers increases. This is partially due to our hardware limit. Similar observation has also been obtained in async-SGD [10].

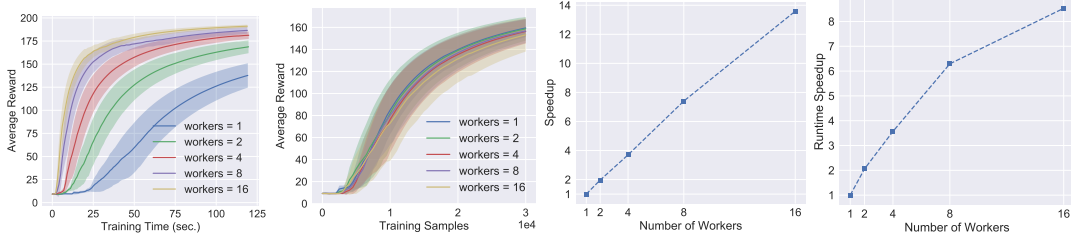


Figure 3: Speedup of A3C-TD(0) in OpenAI gym classic control task (Carpole).



Figure 4: Speedup of A3C-TD(0) in OpenAI Gym Atari game (Seaquest).

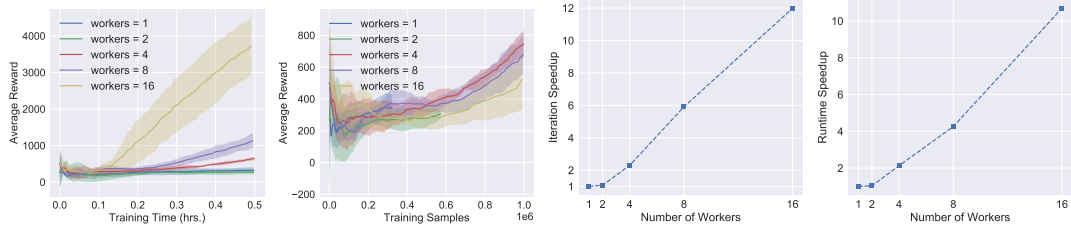


Figure 5: Speedup of A3C-TD(0) in OpenAI Gym Atari game (Beamrider).

6 CONCLUSIONS

This paper revisits the A3C algorithm with TD(0) for the critic update, termed A3C-TD(0). With linear value function approximation, the convergence of the A3C-TD(0) algorithm has been established under both i.i.d. and Markovian sampling settings. Under i.i.d. sampling, A3C-TD(0) achieves linear speedup compared to the best-known sample complexity of two-timescale AC, theoretically justifying the benefit of parallelism and asynchrony for the first time. Under Markov sampling, such a linear speedup can be observed in most classic benchmark tasks.

REFERENCES

- [1] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, “Continuous control with deep reinforcement learning,” in *Proc. of International Conference on Learning Representations*, 2016.
- [2] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, p. 529, 2015.
- [3] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, “Asynchronous methods for deep reinforcement learning,” in *Proc. of International Conference on Machine Learning*, 2016.
- [4] A. Nair, P. Srinivasan, S. Blackwell, C. Alcicek, R. Fearon, A. De Maria, V. Panneershelvam, M. Suleyman, C. Beattie, S. Petersen *et al.*, “Massively parallel methods for deep reinforcement learning,” *arXiv preprint:1507.04296*, 2015.

-
- [5] M. Assran, J. Romoff, N. Ballas, J. Pineau, and M. Rabbat, “Gossip-based actor-learner architectures for deep reinforcement learning,” in *Proc. of Advances in Neural Information Processing Systems*, 2019.
 - [6] T. Chen, K. Zhang, G. B. Giannakis, and T. Başar, “Communication-efficient distributed reinforcement learning,” *arXiv preprint: 1812.03239*, Dec. 2018.
 - [7] V. Konda, *Actor-critic algorithms*. PhD thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, 2002.
 - [8] R. Sutton, D. McAllester, S. Singh, and Y. Mansour, “Policy gradient methods for reinforcement learning with function approximation,” in *Proc. of Advances in Neural Information Processing Systems*, 2000.
 - [9] R. Sutton, “Learning to predict by the methods of temporal differences,” *Machine Learning*, vol. 3, pp. 9–44, 1988.
 - [10] X. Lian, H. Zhang, C. Hsieh, Y. Yijun Huang, and J. Liu, “A comprehensive linear speedup analysis for asynchronous stochastic parallel optimization from zeroth-order to first-order,” in *Proc. of the Advances in Neural Information Processing Systems*, 2016.
 - [11] V. Borkar and V. Konda, “The actor-critic algorithm as multi-time-scale stochastic approximation,” *Sadhana*, vol. 22, no. 4, pp. 525–543, 1997.
 - [12] S. Bhatnagar, R. Sutton, M. Ghavamzadeh, and M. Lee, “Natural actor critic algorithms,” *Automatica*, vol. 45, pp. 2471–2482, 2009.
 - [13] Z. Yang, K. Zhang, M. Hong, and T. Başar, “A finite sample analysis of the actor-critic algorithm,” in *Proc. of IEEE Conference on Decision and Control*, 2018, pp. 2759–2764.
 - [14] H. Kumar, A. Koppel, and A. Ribeiro, “On the sample complexity of actor-critic method for reinforcement learning with function approximation,” *arXiv preprint:1910.08412*, 2019.
 - [15] Z. Fu, Z. Yang, and Z. Wang, “Single-timescale actor-critic provably finds globally optimal policy,” *arXiv preprint:2008.00483*, 2020.
 - [16] S. Qiu, Z. Yang, J. Ye, and Z. Wang, “On the finite-time convergence of actor-critic algorithm,” in *Optimization Foundations for Reinforcement Learning Workshop at Advances in Neural Information Processing Systems*, 2019.
 - [17] T. Xu, Z. Wang, and Y. Liang, “Improving sample complexity bounds for (natural) actor-critic algorithms,” in *Proc. of Advances in Neural Information Processing Systems*, 2020.
 - [18] —, “Non-asymptotic convergence analysis of two time-scale (natural) actor-critic algorithms,” *arXiv preprint:2005.03557*, 2020.
 - [19] Y. Wu, W. Zhang, P. Xu, and Q. Gu, “A finite time analysis of two time-scale actor critic methods,” in *Proc. of Advances in Neural Information Processing Systems*, 2020.
 - [20] M. Hong, H.-T. Wai, Z. Wang, and Z. Yang, “A two-timescale framework for bilevel optimization: Complexity analysis and application to actor-critic,” *arXiv preprint:2007.05170*, 2020.
 - [21] M. Babaeizadeh, I. Frosio, S. Tyree, J. Clemons, and J. Kautz, “Reinforcement learning through asynchronous advantage actor-critic on a gpu,” in *Proc. of International Conference on Learning Representations*, 2017.
 - [22] A. Stooke and P. Abbeel, “Accelerated methods for deep reinforcement learning,” *arXiv preprint:1803.02811*, 2019.
 - [23] L. Espeholt, H. Soyer, R. Munos, K. Simonyan, V. Mnih, T. Ward, Y. Doron, V. Firoiu, T. Harley, I. Dunning, S. Legg, and K. Kavukcuoglu, “Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures,” *arXiv preprint:1802.01561*, 2018.

-
- [24] D. Bertsekas and J. Tsitsiklis, *Parallel and distributed computation: numerical methods*. Prentice-Hall, 1989.
 - [25] A. Agarwal and J. Duchi, “Distributed delayed stochastic optimization,” in *Proc. of Advances in Neural Information Processing Systems*, 2011.
 - [26] H. Feyzmahdavian, A. Aytekin, and M. Johansson, “An asynchronous mini-batch algorithm for regularized stochastic optimization,” *arXiv preprint:1505.04824*, 2015.
 - [27] X. Lian, Y. Huang, Y. Li, and J. Liu, “Asynchronous parallel stochastic gradient for nonconvex optimization,” in *Proc. of Advances in Neural Information Processing Systems*, 2015.
 - [28] T. Sun, R. Hannah, and W. Yin, “Asynchronous coordinate descent under more realistic assumptions,” in *Proc. of Advances in Neural Information Processing Systems*, 2017.
 - [29] X. Lian, C. Zhang, H. Zhang, C.-J. Hsieh, W. Zhang, and J. Liu, “Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent,” in *Proc. of Advances in Neural Information Processing Systems*, 2017.
 - [30] R. J. Williams, “Simple statistical gradient-following algorithms for connectionist reinforcement learning,” *Machine Learning*, vol. 8, no. 3-4, pp. 229–256, May 1992.
 - [31] J. Baxter and P. L. Bartlett, “Infinite-horizon policy-gradient estimation,” *J. Artificial Intelligence Res.*, vol. 15, pp. 319–350, 2001.
 - [32] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT Press, 2018.
 - [33] J. Bhandari, D. Russo, and R. Singal, “A finite-time analysis of temporal difference learning with linear function approximation,” in *Proc. of Conference on Learning Theory*, 2018.
 - [34] T. Sun, H. Shen, T. Chen, and D. Li, “Adaptive temporal difference learning with linear function approximation,” *arXiv preprint:2002.08537*, 2020.
 - [35] T. Xu, Z. Wang, Y. Zhou, and Y. Liang, “Reanalysis of variance reduced temporal difference learning,” in *Proc. of International Conference on Learning Representations*, 2020.
 - [36] S. Zou, T. Xu, and Y. Liang, “Finite-sample analysis for SARSA with linear function approximation,” in *Proc. of Advances in Neural Information Processing Systems*, 2019.
 - [37] K. Zhang, A. Koppel, H. Zhu, and T. Başar, “Global convergence of policy gradient methods to (almost) locally optimal policies,” *arXiv preprint:1906.08383*, 2019.
 - [38] A. Agarwal, S. M. Kakade, J. D. Lee, and G. Mahajan, “Optimality and approximation with policy gradient methods in markov decision processes,” in *Proc. of Thirty Third Conference on Learning Theory*, 2020.
 - [39] K. Doya, “Reinforcement learning in continuous time and space,” *Neural Computation*, vol. 12, no. 1, pp. 219–245, 2000.
 - [40] V. Konda and V. Borkar, “Actor-critic-type learning algorithms for markov decision processes,” *SIAM Journal on Control and Optimization*, vol. 38, no. 1, pp. 94–123, 1999.
 - [41] T. Sun, Y. Sun, and W. Yin, “On markov chain gradient descent,” in *Proc. of Advances in Neural Information Processing Systems*, 2018.
 - [42] R. Sutton, H. Maei, D. Precup, S. Bhatnagar, D. Silver, and E. Szepesvári, C. and Wiewiora, “Fast gradient-descent methods for temporal-difference learning with linear function approximation,” in *Proc. of International Conference on Machine Learning*, 2009.
 - [43] A. Y. Mitrophanov, “Sensitivity and convergence of uniformly ergodic markov chains,” *Journal of Applied Probability*, vol. 42, no. 4, pp. 1003–1014, 2005.
 - [44] Dgriff, “Pytorch implementation of a3c,” https://github.com/dgriff777/rl_a3c_pytorch, 2018.

Supplementary Material

A PRELIMINARY LEMMAS

A.1 GEOMETRIC MIXING

The operation $p \otimes q$ denotes the tensor product between two distributions $p(x)$ and $q(y)$, i.e. $(p \otimes q)(x, y) = p(x) \cdot q(y)$.

Lemma 1. *Suppose Assumption 4 holds. For any $\theta \in \mathbb{R}^d$, we have*

$$\sup_{s_0 \in \mathcal{S}} d_{TV}(\mathbb{P}((s_t, a_t, s_{t+1}) \in \cdot | s_0, \pi_\theta), \mu_\theta \otimes \pi_\theta \otimes \mathcal{P}) \leq \kappa \rho^t. \quad (21)$$

where $\mu_\theta(\cdot)$ is the stationary distribution with policy π_θ and transition kernel $\mathcal{P}(\cdot | s, a)$.

Proof. We start with

$$\begin{aligned} & \sup_{s_0 \in \mathcal{S}} d_{TV}(\mathbb{P}((s_t, a_t, s_{t+1}) \in \cdot | s_0, \pi_\theta), \mu_\theta \otimes \pi_\theta \otimes \mathcal{P}) \\ &= \sup_{s_0 \in \mathcal{S}} d_{TV}(\mathbb{P}(s_t = \cdot | s_0, \pi_\theta) \otimes \pi_\theta \otimes \mathcal{P}, \mu_\theta \otimes \pi_\theta \otimes \mathcal{P}) \\ &= \sup_{s_0 \in \mathcal{S}} \frac{1}{2} \int_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \int_{s' \in \mathcal{S}} |\mathbb{P}(s_t = ds | s_0, \pi_\theta) \pi_\theta(a | s) \mathcal{P}(ds' | s, a) - \mu_\theta(ds) \pi_\theta(a | s) \mathcal{P}(ds' | s, a)| \\ &= \sup_{s_0 \in \mathcal{S}} \frac{1}{2} \int_{s \in \mathcal{S}} |\mathbb{P}(s_t = ds | s_0, \pi_\theta) - \mu_\theta(ds)| \sum_{a \in \mathcal{A}} \pi_\theta(a | s) \int_{s' \in \mathcal{S}} \mathcal{P}(ds' | s, a) \\ &= \sup_{s_0 \in \mathcal{S}} d_{TV}(\mathbb{P}(s_t \in \cdot | s_0, \pi_\theta), \mu_\theta) \\ &\leq \kappa \rho^t, \end{aligned}$$

which completes the proof. \square

For the use in the later proof, given $K > 0$, we first define m_K as:

$$m_K := \min \{ m \in \mathbb{N}^+ \mid \kappa \rho^{m-1} \leq \min\{\alpha_k, \beta_k\} \}, \quad (22)$$

where κ and ρ are constants defined in (4). m_K is the minimum number of samples needed for the Markov chain to approach the stationary distribution so that the bias incurred by the Markovian sampling is small enough.

A.2 DISTRIBUTION MISMATCH

In this section, we will bound the mismatch between the visitation measure d_θ and stationary distribution μ_θ in terms of divergence. It will later be used in the analysis of actor convergence.

Lemma 2. *Suppose Assumption 4 holds. For any $\theta \in \mathbb{R}^d$, it holds that*

$$d_{TV}(d_\theta, \mu_\theta) \leq 2 \left(\log_\rho \kappa^{-1} + \frac{1}{1-\rho} \right) (1-\gamma). \quad (23)$$

Proof. By [7], the discounted visitation measure d_θ can be viewed as the stationary distribution of an artificial Markov chain with transition kernel $\tilde{\mathcal{P}} = (1-\gamma)\eta + \gamma\mathcal{P}$. Denote the stationary distribution of the artificial Markov chain with policy π_θ and transition kernel $\tilde{\mathcal{P}}$ as $\tilde{\mu}_\theta$, then we have

$$d_{TV}(d_\theta, \mu_\theta) = \int_{s \in \mathcal{S}} |d_\theta(s) - \mu_\theta(s)| = \int_{s \in \mathcal{S}} |\tilde{\mu}_\theta(s) - \mu_\theta(s)|. \quad (24)$$

Following [43, Theorem 3.1], the second term in (24) can be bounded as

$$\begin{aligned} \int_{s \in \mathcal{S}} |\tilde{\mu}_\theta(s) - \mu_\theta(s)| &\leq \left(\log_\rho \kappa^{-1} + \frac{1}{1-\rho} \right) \sup_s \int_{s' \in \mathcal{S}} \left| \sum_a \pi_\theta(a | s) \left(\tilde{\mathcal{P}}(s' | s, a) - \mathcal{P}(s' | s, a) \right) \right| \\ &\leq \left(\log_\rho \kappa^{-1} + \frac{1}{1-\rho} \right) \sup_s \sum_a \pi_\theta(a | s) \int_{s' \in \mathcal{S}} \left| \tilde{\mathcal{P}}(s' | s, a) - \mathcal{P}(s' | s, a) \right| \end{aligned} \quad (25)$$

in which we have

$$\int_{s' \in \mathcal{S}} \left| \tilde{\mathcal{P}}(s'|s, a) - \mathcal{P}(s'|s, a) \right| = (1 - \gamma) \int_{s' \in \mathcal{S}} |\mathcal{P}(s'|s, a) - \eta(s')| \leq 2(1 - \gamma). \quad (26)$$

Substituting the last inequality into (25) completes the proof. \square

A.3 AUXILIARY MARKOV CHAIN

The auxiliary Markov chain is a virtual Markov chain with no policy drifting — a technique developed in [36] to analyze stochastic approximation algorithms in non-stationary settings.

Lemma 3. *Under Assumption 1 and Assumption 3, consider the update (9) in Algorithm 1 with Markovian sampling. For a given number of samples m , consider the Markov chain of the worker that contributes to the k th update:*

$$s_{t-m} \xrightarrow{\theta_{k-d_m}} a_{t-m} \xrightarrow{\mathcal{P}} s_{t-m+1} \xrightarrow{\theta_{k-d_{m-1}}} a_{t-m+1} \cdots s_{t-1} \xrightarrow{\theta_{k-d_1}} a_{t-1} \xrightarrow{\mathcal{P}} s_t \xrightarrow{\theta_{k-d_0}} a_t \xrightarrow{\mathcal{P}} s_{t+1},$$

where $(s_t, a_t, s_{t+1}) = (s_{(k)}, a_{(k)}, s'_{(k)})$, and $\{d_j\}_{j=0}^m$ is some increasing sequence with $d_0 := \tau_k$.

Given $(s_{t-m}, a_{t-m}, s_{t-m+1})$ and θ_{k-d_m} , we construct its auxiliary Markov chain by repeatedly applying $\pi_{\theta_{k-d_m}}$:

$$s_{t-m} \xrightarrow{\theta_{k-d_m}} a_{t-m} \xrightarrow{\mathcal{P}} s_{t-m+1} \xrightarrow{\theta_{k-d_m}} \tilde{a}_{t-m+1} \cdots \tilde{s}_{t-1} \xrightarrow{\theta_{k-d_m}} \tilde{a}_{t-1} \xrightarrow{\mathcal{P}} \tilde{s}_t \xrightarrow{\theta_{k-d_m}} \tilde{a}_t \xrightarrow{\mathcal{P}} \tilde{s}_{t+1}.$$

Define $x_t := (s_t, a_t, s_{t+1})$, then we have:

$$\begin{aligned} & d_{TV}(\mathbb{P}(x_t \in \cdot | \theta_{k-d_m}, s_{t-m+1}), \mathbb{P}(\tilde{x}_t \in \cdot | \theta_{k-d_m}, s_{t-m+1})) \\ & \leq \frac{1}{2} |\mathcal{A}| L_\pi \sum_{i=\tau_k}^{d_m} \mathbb{E}[\|\theta_{k-i} - \theta_{k-d_m}\|_2 | \theta_{k-d_m}, s_{t-m+1}]. \end{aligned} \quad (27)$$

Proof. Throughout the lemma, all expectations and probabilities are conditioned on θ_{k-d_m} and s_{t-m+1} . We omit this condition for convenience.

With $\tilde{x}_t := (\tilde{s}_t, \tilde{a}_t, \tilde{s}_{t+1})$, first we have

$$\begin{aligned} & d_{TV}(\mathbb{P}(s_{t+1} \in \cdot), \mathbb{P}(\tilde{s}_{t+1} \in \cdot)) \\ & = \frac{1}{2} \int_{s' \in \mathcal{S}} |\mathbb{P}(s_{t+1} = ds') - \mathbb{P}(\tilde{s}_{t+1} = ds')| \\ & = \frac{1}{2} \int_{s' \in \mathcal{S}} \left| \int_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \mathbb{P}(s_t = ds, a_t = a, s_{t+1} = ds') - \mathbb{P}(\tilde{s}_t = ds, \tilde{a}_t = a, \tilde{s}_{t+1} = ds') \right| \\ & \leq \frac{1}{2} \int_{s' \in \mathcal{S}} \int_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} |\mathbb{P}(s_t = ds, a_t = a, s_{t+1} = ds') - \mathbb{P}(\tilde{s}_t = ds, \tilde{a}_t = a, \tilde{s}_{t+1} = ds')| \\ & = \frac{1}{2} \int_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \int_{s' \in \mathcal{S}} |\mathbb{P}(s_t = ds, a_t = a, s_{t+1} = ds') - \mathbb{P}(\tilde{s}_t = ds, \tilde{a}_t = a, \tilde{s}_{t+1} = ds')| \\ & = d_{TV}(\mathbb{P}(x_t \in \cdot), \mathbb{P}(\tilde{x}_t \in \cdot)), \end{aligned} \quad (28)$$

where the second last equality is due to Tonelli's theorem. Next we have

$$\begin{aligned} & d_{TV}(\mathbb{P}(x_t \in \cdot), \mathbb{P}(\tilde{x}_t \in \cdot)) \\ & = \frac{1}{2} \int_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \int_{s' \in \mathcal{S}} |\mathbb{P}(s_t = ds, a_t = a, s_{t+1} = ds') - \mathbb{P}(\tilde{s}_t = ds, \tilde{a}_t = a, \tilde{s}_{t+1} = ds')| \\ & = \frac{1}{2} \int_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} |\mathbb{P}(s_t = ds, a_t = a) - \mathbb{P}(\tilde{s}_t = ds, \tilde{a}_t = a)| \int_{s' \in \mathcal{S}} \mathcal{P}(s_{t+1} = ds' | s_t = ds, a_t = a) \\ & = \frac{1}{2} \int_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} |\mathbb{P}(s_t = ds, a_t = a) - \mathbb{P}(\tilde{s}_t = ds, \tilde{a}_t = a)| \\ & = d_{TV}(\mathbb{P}((s_t, a_t) \in \cdot), \mathbb{P}((\tilde{s}_t, \tilde{a}_t) \in \cdot)). \end{aligned} \quad (29)$$

Due to the fact that $\theta_{k-\tau_k}$ is dependent on s_t , we need to write $\mathbb{P}(s_t, a_t)$ as

$$\begin{aligned}\mathbb{P}(s_t, a_t) &= \int_{\theta_{k-\tau_k} \in \mathbb{R}^d} \mathbb{P}(s_t, \theta_{k-\tau_k}, a_t) \\ &= \int_{\theta \in \mathbb{R}^d} \mathbb{P}(s_t) \mathbb{P}(\theta_{k-\tau_k} = d\theta | s_t) \pi_{\theta_{k-\tau_k}}(a_t | s_t) \\ &= \mathbb{P}(s_t) \mathbb{E}[\pi_{\theta_{k-\tau_k}}(a_t | s_t) | s_t].\end{aligned}$$

Then we have

$$\begin{aligned}d_{TV}(\mathbb{P}((s_t, a_t) \in \cdot), \mathbb{P}((\tilde{s}_t, \tilde{a}_t) \in \cdot)) &= \frac{1}{2} \int_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \left| \mathbb{P}(s_t = ds) \mathbb{E}[\pi_{\theta_{k-\tau_k}}(a_t = a | s_t = ds) | s_t = ds] - \mathbb{P}(\tilde{s}_t = ds) \pi_{\theta_{k-d_m}}(\tilde{a}_t = a | \tilde{s}_t = ds) \right| \\ &\leq \frac{1}{2} \int_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \left| \mathbb{P}(s_t = ds) \mathbb{E}[\pi_{\theta_{k-\tau_k}}(a_t = a | s_t = ds) | s_t = ds] - \mathbb{P}(s_t = ds) \pi_{\theta_{k-d_m}}(a_t = a | s_t = ds) \right| \\ &\quad + \frac{1}{2} \int_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \left| \mathbb{P}(s_t = ds) \pi_{\theta_{k-d_m}}(\tilde{a}_t = a | \tilde{s}_t = ds) - \mathbb{P}(\tilde{s}_t = ds) \pi_{\theta_{k-d_m}}(\tilde{a}_t = a | \tilde{s}_t = ds) \right| \\ &= \frac{1}{2} \int_{s \in \mathcal{S}} \mathbb{P}(s_t = ds) \sum_{a \in \mathcal{A}} \left| \mathbb{E}[\pi_{\theta_{k-\tau_k}}(a_t = a | s_t = ds) | s_t = ds] - \pi_{\theta_{k-d_m}}(a_t = a | s_t = ds) \right| \\ &\quad + \frac{1}{2} \int_{s \in \mathcal{S}} |\mathbb{P}(s_t = ds) - \mathbb{P}(\tilde{s}_t = ds)|. \tag{30}\end{aligned}$$

Using Jensen's inequality, we have

$$\begin{aligned}d_{TV}(\mathbb{P}((s_t, a_t) \in \cdot), \mathbb{P}((\tilde{s}_t, \tilde{a}_t) \in \cdot)) &\leq \frac{1}{2} \int_{s \in \mathcal{S}} \mathbb{P}(s_t = ds) \sum_{a \in \mathcal{A}} \mathbb{E} \left[\left| \pi_{\theta_{k-\tau_k}}(a_t = a | s_t = ds) - \pi_{\theta_{k-d_m}}(a_t = a | s_t = ds) \right| \middle| s_t = ds \right] \\ &\quad + \frac{1}{2} \int_{s \in \mathcal{S}} |\mathbb{P}(s_t = ds) - \mathbb{P}(\tilde{s}_t = ds)| \\ &\leq \frac{1}{2} \int_{s \in \mathcal{S}} \mathbb{P}(s_t = ds) \sum_{a \in \mathcal{A}} \mathbb{E} [\|\theta_{k-\tau_k} - \theta_{k-d_m}\|_2 | s_t = ds] + \frac{1}{2} \int_{s \in \mathcal{S}} |\mathbb{P}(s_t = ds) - \mathbb{P}(\tilde{s}_t = ds)| \\ &= \frac{1}{2} |\mathcal{A}| L_\pi \mathbb{E} \|\theta_{k-\tau_k} - \theta_{k-d_m}\|_2 + d_{TV}(\mathbb{P}(s_t \in \cdot), \mathbb{P}(\tilde{s}_t \in \cdot)) \tag{31}\end{aligned}$$

where the last inequality follows Assumption 3.

Now we start to prove (27). First we have

$$\begin{aligned}d_{TV}(\mathbb{P}(x_t \in \cdot), \mathbb{P}(\tilde{x}_t \in \cdot)) &\stackrel{(29)}{=} d_{TV}(\mathbb{P}((s_t, a_t) \in \cdot), \mathbb{P}((\tilde{s}_t, \tilde{a}_t) \in \cdot)) \\ &\stackrel{(31)}{\leq} d_{TV}(\mathbb{P}(s_t \in \cdot), \mathbb{P}(\tilde{s}_t \in \cdot)) + \frac{1}{2} |\mathcal{A}| L_\pi \mathbb{E} \|\theta_{k-\tau_k} - \theta_{k-d_m}\|_2 \\ &\stackrel{(28)}{\leq} d_{TV}(\mathbb{P}(x_{t-1} \in \cdot), \mathbb{P}(\tilde{x}_{t-1} \in \cdot)) + \frac{1}{2} |\mathcal{A}| L_\pi \mathbb{E} \|\theta_{k-\tau_k} - \theta_{k-d_m}\|_2. \tag{32}\end{aligned}$$

Since $d_{TV}(\mathbb{P}(x_{t-m} \in \cdot), \mathbb{P}(x_{t-m} \in \cdot)) = 0$, recursively applying (32) for $\{t-1, \dots, t-m\}$ gives

$$\begin{aligned}d_{TV}(\mathbb{P}(x_t \in \cdot), \mathbb{P}(\tilde{x}_t \in \cdot)) &\leq \frac{1}{2} |\mathcal{A}| L_\pi \sum_{j=0}^m \mathbb{E} \|\theta_{k-d_j} - \theta_{k-d_m}\|_2 \\ &\leq \frac{1}{2} |\mathcal{A}| L_\pi \sum_{i=\tau_k}^{d_m} \mathbb{E} \|\theta_{k-i} - \theta_{k-d_m}\|_2,\end{aligned}$$

which completes the proof. \square

A.4 LIPSCHITZ CONTINUITY OF VALUE FUNCTION

Lemma 4. Suppose Assumption 3 holds. For any $\theta_1, \theta_2 \in \mathbb{R}^d$ and $s \in \mathcal{S}$, we have

$$\|\nabla V_{\pi_{\theta_1}}(s)\|_2 \leq L_V, \quad (33a)$$

$$|V_{\pi_{\theta_1}}(s) - V_{\pi_{\theta_2}}(s)| \leq L_V \|\theta_1 - \theta_2\|_2, \quad (33b)$$

where the constant is $L_V := C_\psi R_{\max}$ with C_ψ defined as in Assumption 3.

Proof. First we have

$$\begin{aligned} Q_\pi(s, a) &= \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, s_{t+1}) | s_0 = s, a_0 = a \right] \\ &\leq \sum_{t=0}^{\infty} \gamma^t r_{\max} = R_{\max}. \end{aligned}$$

By the policy gradient theorem [8], we have

$$\begin{aligned} \|\nabla V_{\pi_{\theta_1}}(s)\|_2 &= \|\mathbb{E} [Q_{\pi_{\theta_1}}(s, a) \psi_{\theta_1}(s, a)]\|_2 \\ &\leq \mathbb{E} \|Q_{\pi_{\theta_1}}(s, a) \psi_{\theta_1}(s, a)\|_2 \\ &\leq \mathbb{E} [|Q_{\pi_{\theta_1}}(s, a)| \|\psi_{\theta_1}(s, a)\|_2] \\ &\leq R_{\max} C_\psi, \end{aligned}$$

where the first inequality is due to Jensen's inequality, and the last inequality follows Assumption 3 and the fact that $Q_\pi(s, a) \leq R_{\max}$. By the mean value theorem, we immediately have

$$|V_{\pi_{\theta_1}}(s) - V_{\pi_{\theta_2}}(s)| \leq \sup_{\theta_1 \in \mathbb{R}^d} \|\nabla V_{\pi_{\theta_1}}(s)\|_2 \|\theta_1 - \theta_2\|_2 = L_V \|\theta_1 - \theta_2\|_2,$$

which completes the proof. \square

A.5 LIPSCHITZ CONTINUITY OF POLICY GRADIENT

We give a proposition regarding the L_J -Lipschitz of the policy gradient under proper assumptions, which has been shown by [37].

Proposition 1. Suppose Assumption 3 and 4 hold. For any $\theta, \theta' \in \mathbb{R}^d$, we have $\|\nabla J(\theta) - \nabla J(\theta')\|_2 \leq L_J \|\theta - \theta'\|_2$, where L_J is a positive constant.

A.6 LIPSCHITZ CONTINUITY OF OPTIMAL CRITIC PARAMETER

We provide a justification for Lipschitz continuity of ω_θ^* in the next proposition.

Proposition 2. Suppose Assumption 3 and 4 hold. For any $\theta_1, \theta_2 \in \mathbb{R}^d$, we have

$$\|\omega_{\theta_1}^* - \omega_{\theta_2}^*\|_2 \leq L_\omega \|\theta_1 - \theta_2\|_2,$$

where $L_\omega := 2r_{\max} |\mathcal{A}| L_\pi (\lambda^{-1} + \lambda^{-2}(1 + \gamma))(1 + \log_\rho \kappa^{-1} + (1 - \rho)^{-1})$.

Proof. We use A_1, A_2, b_1 and b_2 as shorthand notations of $A_{\pi_{\theta_1}}, A_{\pi_{\theta_2}}, b_{\pi_{\theta_1}}$ and $b_{\pi_{\theta_2}}$ respectively. By Assumption 2, $A_{\theta, \phi}$ is invertible for any $\theta \in \mathbb{R}^d$, so we can write $\omega_\theta^* = -A_{\theta, \phi}^{-1} b_{\theta, \phi}$. Then we have

$$\begin{aligned} \|\omega_1^* - \omega_2^*\|_2 &= \|-A_1^{-1} b_1 + A_2^{-1} b_2\|_2 \\ &= \|-A_1^{-1} b_1 - A_1^{-1} b_2 + A_1^{-1} b_2 + A_2^{-1} b_2\|_2 \\ &= \|-A_1^{-1} (b_1 - b_2) - (A_1^{-1} - A_2^{-1}) b_2\|_2 \\ &\leq \|A_1^{-1} (b_1 - b_2)\|_2 + \|(A_1^{-1} - A_2^{-1}) b_2\|_2 \\ &\leq \|A_1^{-1}\|_2 \|b_1 - b_2\|_2 + \|A_1^{-1} - A_2^{-1}\|_2 \|b_2\|_2 \\ &= \|A_1^{-1}\|_2 \|b_1 - b_2\|_2 + \|A_1^{-1} (A_2 - A_1) A_2^{-1}\|_2 \|b_2\|_2 \\ &\leq \|A_1^{-1}\|_2 \|b_1 - b_2\|_2 + \|A_1^{-1}\|_2 \|A_2^{-1}\|_2 \|b_2\|_2 \|A_2 - A_1\|_2 \\ &\leq \lambda^{-1} \|b_1 - b_2\|_2 + \lambda^{-2} r_{\max} \|A_1 - A_2\|_2, \end{aligned} \quad (34)$$

where the last inequality follows Assumption 2, and the fact that

$$\|b_2\|_2 = \|\mathbb{E}[r(s, a, s')\phi(s)]\|_2 \leq \mathbb{E}\|r(s, a, s')\phi(s)\|_2 \leq \mathbb{E}[\|r(s, a, s')\| \|\phi(s)\|_2] \leq r_{\max}.$$

Denote (s^1, a^1, s'^1) and (s^2, a^2, s'^2) as samples drawn with θ_1 and θ_2 respectively, i.e. $s^1 \sim \mu_{\theta_1}$, $a^1 \sim \pi_{\theta_1}$, $s'^1 \sim \mathcal{P}$ and $s^2 \sim \mu_{\theta_2}$, $a^2 \sim \pi_{\theta_2}$, $s'^2 \sim \mathcal{P}$. Then we have

$$\begin{aligned} \|b_1 - b_2\|_2 &= \|\mathbb{E}[r(s^1, a^1, s'^1)\phi(s^1)] - \mathbb{E}[r(s^2, a^2, s'^2)\phi(s^2)]\|_2 \\ &\leq \sup_{s, a, s'} \|r(s, a, s')\phi(s)\|_2 \|\mathbb{P}((s^1, a^1, s'^1) \in \cdot) - \mathbb{P}((s^2, a^2, s'^2) \in \cdot)\|_{TV} \\ &\leq r_{\max} \|\mathbb{P}((s^1, a^1, s'^1) \in \cdot) - \mathbb{P}((s^2, a^2, s'^2) \in \cdot)\|_{TV} \\ &= 2r_{\max} d_{TV}(\mu_{\theta_1} \otimes \pi_{\theta_1} \otimes \mathcal{P}, \mu_{\theta_2} \otimes \pi_{\theta_2} \otimes \mathcal{P}) \\ &\leq 2r_{\max} |\mathcal{A}| L_{\pi} (1 + \log_{\rho} \kappa^{-1} + (1 - \rho)^{-1}) \|\theta_1 - \theta_2\|_2, \end{aligned} \quad (35)$$

where the first inequality follows the definition of total variation (TV) norm, and the last inequality follows Lemma A.1. in [19]. Similarly we have:

$$\begin{aligned} \|A_1 - A_2\|_2 &\leq 2(1 + \gamma) d_{TV}(\mu_{\theta_1} \otimes \pi_{\theta_1}, \mu_{\theta_2} \otimes \pi_{\theta_2}) \\ &= (1 + \gamma) |\mathcal{A}| L_{\pi} (1 + \log_{\rho} \kappa^{-1} + (1 - \rho)^{-1}) \|\theta_1 - \theta_2\|_2. \end{aligned} \quad (36)$$

Substituting (35) and (36) into (34) completes the proof. \square

B PROOF OF MAIN THEOREMS

B.1 PROOF OF THEOREM 1

For brevity, we first define the following notations:

$$\begin{aligned} x &:= (s, a, s'), \\ \hat{\delta}(x, \omega) &:= r(s, a, s') + \gamma \phi(s')^\top \omega - \phi(s)^\top \omega, \\ g(x, \omega) &:= \hat{\delta}(x, \omega) \phi(s), \\ \bar{g}(\theta, \omega) &:= \mathbb{E}_{s \sim \mu_{\theta}, a \sim \pi_{\theta}, s' \sim \mathcal{P}} [g(x, \omega)]. \end{aligned}$$

We also define constant $C_{\delta} := r_{\max} + (1 + \gamma) \max\{R_{\max}, R_{\omega}\}$, and we immediately have

$$\|g(x, \omega)\|_2 \leq |r(x) + \gamma \phi(s')^\top \omega - \phi(s)^\top \omega| \leq r_{\max} + (1 + \gamma) R_{\omega} \leq C_{\delta} \quad (37)$$

and likewise, we have $\|\bar{g}(x, \omega)\|_2 \leq C_{\delta}$.

The critic update in Algorithm 1 can be written compactly as:

$$\omega_{k+1} = \Pi_{R_{\omega}}(\omega_k + \beta_k g(x_{(k)}, \omega_{k-\tau_k})), \quad (38)$$

where τ_k is the delay of the parameters used in evaluating the k th stochastic gradient, and $x_{(k)} := (s_{(k)}, a_{(k)}, s'_{(k)})$ is the sample used to evaluate the stochastic gradient at k th update.

Proof. Using ω_k^* as shorthand notation of $\omega_{\theta_k}^*$, we start with the optimality gap

$$\begin{aligned} &\|\omega_{k+1} - \omega_{k+1}^*\|_2^2 \\ &= \|\Pi_{R_{\omega}}(\omega_k + \beta_k g(x_{(k)}, \omega_{k-\tau_k})) - \omega_{k+1}^*\|_2^2 \\ &\leq \|\omega_k + \beta_k g(x_{(k)}, \omega_{k-\tau_k}) - \omega_{k+1}^*\|_2^2 \\ &= \|\omega_k - \omega_k^*\|_2^2 + 2\beta_k \langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_{k-\tau_k}) \rangle + 2\langle \omega_k - \omega_k^*, \omega_k^* - \omega_{k+1}^* \rangle \\ &\quad + \|\omega_k^* - \omega_{k+1}^* + \beta_k g(x_{(k)}, \omega_{k-\tau_k})\|_2^2 \\ &= \|\omega_k - \omega_k^*\|_2^2 + 2\beta_k \langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_{k-\tau_k}) - g(x_{(k)}, \omega_k) \rangle + 2\beta_k \langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_k) - \bar{g}(\theta_k, \omega_k) \rangle \\ &\quad + 2\beta_k \langle \omega_k - \omega_k^*, \bar{g}(\theta_k, \omega_k) \rangle + 2\langle \omega_k - \omega_k^*, \omega_k^* - \omega_{k+1}^* \rangle + \|\omega_k^* - \omega_{k+1}^* + \beta_k g(x_{(k)}, \omega_{k-\tau_k})\|_2^2 \\ &\leq \|\omega_k - \omega_k^*\|_2^2 + 2\beta_k \langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_{k-\tau_k}) - g(x_{(k)}, \omega_k) \rangle + 2\beta_k \langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_k) - \bar{g}(\theta_k, \omega_k) \rangle \\ &\quad + 2\beta_k \langle \omega_k - \omega_k^*, \bar{g}(\theta_k, \omega_k) \rangle + 2\langle \omega_k - \omega_k^*, \omega_k^* - \omega_{k+1}^* \rangle + 2\|\omega_k^* - \omega_{k+1}^*\|_2^2 + 2C_{\delta}^2 \beta_k^2. \end{aligned} \quad (39)$$

We first bound $\langle \omega_k - \omega_k^*, \bar{g}(\theta_k, \omega_k) \rangle$ in (39) as

$$\begin{aligned}
\langle \omega_k - \omega_k^*, \bar{g}(\theta_k, \omega_k) \rangle &= \langle \omega_k - \omega_k^*, \bar{g}(\theta_k, \omega_k) - \bar{g}(\theta_k, \omega_k^*) \rangle \\
&= \left\langle \omega_k - \omega_k^*, \mathbb{E} \left[(\gamma \phi(s') - \phi(s))^\top (\omega_k - \omega_k^*) \phi(s) \right] \right\rangle \\
&= \left\langle \omega_k - \omega_k^*, \mathbb{E} \left[\phi(s) (\gamma \phi(s') - \phi(s))^\top \right] (\omega_k - \omega_k^*) \right\rangle \\
&= \left\langle \omega_k - \omega_k^*, A_{\pi_{\theta_k}} (\omega_k - \omega_k^*) \right\rangle \\
&\leq -\lambda \|\omega_k - \omega_k^*\|_2^2,
\end{aligned} \tag{40}$$

where the first equality is due to $\bar{g}(\theta, \omega_\theta^*) = A_{\theta, \phi} \omega_\theta^* + b = 0$, and the last inequality follows Assumption 2. Substituting (40) into (39), then taking expectation on both sides of (39) yield

$$\begin{aligned}
\mathbb{E} \|\omega_{k+1} - \omega_{k+1}^*\|_2^2 &\leq (1 - 2\lambda\beta_k) \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 + 2\beta_k \mathbb{E} \langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_{k-\tau_k}) - g(x_{(k)}, \omega_k) \rangle \\
&\quad + 2\beta_k \mathbb{E} \langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_k) - \bar{g}(\theta_k, \omega_k) \rangle + 2 \mathbb{E} \langle \omega_k - \omega_k^*, \omega_k^* - \omega_{k+1}^* \rangle \\
&\quad + 2 \mathbb{E} \|\omega_k^* - \omega_{k+1}^*\|_2^2 + 2C_\delta^2 \beta_k^2.
\end{aligned} \tag{41}$$

We then bound the term $\mathbb{E} \langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_{k-\tau_k}) - g(x_{(k)}, \omega_k) \rangle$ in (41) as

$$\begin{aligned}
\mathbb{E} \langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_{k-\tau_k}) - g(x_{(k)}, \omega_k) \rangle &= \mathbb{E} \left\langle \omega_k - \omega_k^*, \left(\gamma \phi(s'_{(k)}) - \phi(s_{(k)}) \right)^\top (\omega_{k-\tau_k} - \omega_k) \phi(s_{(k)}) \right\rangle \\
&\leq (1 + \gamma) \mathbb{E} [\|\omega_k - \omega_k^*\|_2 \|\omega_{k-\tau_k} - \omega_k\|_2] \\
&\leq (1 + \gamma) \mathbb{E} \left[\|\omega_k - \omega_k^*\|_2 \left\| \sum_{i=k-\tau_k}^{k-1} (\omega_{i+1} - \omega_i) \right\|_2 \right] \\
&\leq (1 + \gamma) \mathbb{E} \left[\|\omega_k - \omega_k^*\|_2 \sum_{i=k-\tau_k}^{k-1} \beta_i \|g(x_i, \omega_{i-\tau_i})\|_2 \right] \\
&\leq (1 + \gamma) \mathbb{E} \left[\|\omega_k - \omega_k^*\|_2 \sum_{i=k-\tau_k}^{k-1} \beta_{k-K_0} \|g(x_i, \omega_{i-\tau_i})\|_2 \right] \\
&\leq C_\delta (1 + \gamma) K_0 \beta_{k-K_0} \mathbb{E} \|\omega_k - \omega_k^*\|_2,
\end{aligned} \tag{42}$$

where the second last inequality is due to the monotonicity of step size, and the last inequality follows the definition of C_δ in (37).

Next we jointly bound the fourth and fifth term in (41) as

$$\begin{aligned}
&2 \mathbb{E} \langle \omega_k - \omega_k^*, \omega_k^* - \omega_{k+1}^* \rangle + 2 \mathbb{E} \|\omega_k^* - \omega_{k+1}^*\|_2^2 \\
&\leq 2 \mathbb{E} [\|\omega_k - \omega_k^*\|_2 \|\omega_k^* - \omega_{k+1}^*\|_2] + 2 \mathbb{E} \|\omega_k^* - \omega_{k+1}^*\|_2^2 \\
&\leq 2L_\omega \mathbb{E} [\|\omega_k - \omega_k^*\|_2 \|\theta_k - \theta_{k+1}\|_2] + 2L_\omega^2 \mathbb{E} \|\theta_k - \theta_{k+1}\|_2^2 \\
&= 2L_\omega \alpha_k \mathbb{E} \left[\|\omega_k - \omega_k^*\|_2 \left\| \hat{\delta}(x_{(k)}, \omega_{k-\tau_k}) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\|_2 \right] \\
&\quad + 2L_\omega^2 \alpha_k^2 \mathbb{E} \left\| \hat{\delta}(x_{(k)}, \omega_{k-\tau_k}) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\|_2^2 \\
&\leq 2L_\omega C_p \alpha_k \mathbb{E} \|\omega_k - \omega_k^*\|_2 + 2L_\omega^2 C_p^2 \alpha_k^2,
\end{aligned} \tag{43}$$

where constant $C_p := C_\delta C_\psi$. The second inequality is due to the L_ω -Lipschitz of ω_θ^* shown in Proposition 2, and the last inequality follows the fact that

$$\left\| \hat{\delta}(x_{(k)}, \omega_{k-\tau_k}) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\|_2 \leq C_\delta C_\psi = C_p. \tag{44}$$

Substituting (42) and (43) into (41) yields

$$\begin{aligned}
\mathbb{E} \|\omega_{k+1} - \omega_{k+1}^*\|_2^2 &\leq (1 - 2\lambda\beta_k) \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 + 2\beta_k (C_1 \frac{\alpha_k}{\beta_k} + C_2 K_0 \beta_{k-K_0}) \mathbb{E} \|\omega_k - \omega_k^*\|_2 \\
&\quad + 2\beta_k \mathbb{E} \langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_k) - \bar{g}(\theta_k, \omega_k) \rangle + C_q \beta_k^2,
\end{aligned} \tag{45}$$

where $C_1 := L_\omega C_p$, $C_2 := C_\delta(1 + \gamma)$ and $C_q := 2C_\delta^2 + 2L_\omega^2 C_p^2 \max_{(k)} \frac{\alpha_k^2}{\beta_k^2} = 2C_\delta^2 + 2L_\omega^2 C_p^2 \frac{c_1^2}{c_2^2}$.

For brevity, we use $x \sim \theta$ to denote $s \sim \mu_\theta$, $a \sim \pi_\theta$ and $s' \sim \mathcal{P}$ in this proof. Consider the third term in (45) conditioned on $\theta_k, \omega_k, \theta_{k-\tau_k}$. We bound it as

$$\begin{aligned}
& \mathbb{E} [\langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_k) - \bar{g}(\theta_k, \omega_k) \rangle | \theta_k, \omega_k, \theta_{k-\tau_k}] \\
&= \left\langle \omega_k - \omega_k^*, \mathbb{E}_{x_{(k)} \sim \theta_{k-\tau_k}} [g(x_{(k)}, \omega_k) | \omega_k] - \bar{g}(\theta_k, \omega_k) \right\rangle \\
&= \left\langle \omega_k - \omega_k^*, \bar{g}(\theta_{k-\tau_k}, \omega_k) - \bar{g}(\theta_k, \omega_k) \right\rangle \\
&\leq \|\omega_k - \omega_k^*\|_2 \|\bar{g}(\theta_{k-\tau_k}, \omega_k) - \bar{g}(\theta_k, \omega_k)\|_2 \\
&\leq 2R_\omega \left\| \mathbb{E}_{x \sim \theta_{k-\tau_k}} [g(x, \omega_k)] - \mathbb{E}_{x \sim \theta_k} [g(x, \omega_k)] \right\|_2 \\
&\leq 2R_\omega \sup_x \|g(x, \omega_k)\|_2 \left\| \mu_{\theta_{k-\tau_k}} \otimes \pi_{\theta_{k-\tau_k}} \otimes \mathcal{P} - \mu_{\theta_k} \otimes \pi_{\theta_k} \otimes \mathcal{P} \right\|_{TV} \\
&\leq 4R_\omega C_\delta d_{TV}(\mu_{\theta_{k-\tau_k}} \otimes \pi_{\theta_{k-\tau_k}} \otimes \mathcal{P}, \mu_{\theta_k} \otimes \pi_{\theta_k} \otimes \mathcal{P}), \tag{46}
\end{aligned}$$

where second last inequality follows the definition of TV norm and the last inequality uses the definition of C_δ in (37).

Define constant $C_3 := 2R_\omega C_\delta |\mathcal{A}| L_\pi (1 + \log_\rho \kappa^{-1} + (1 - \rho)^{-1})$. Then by following the third item in [19, Lemma A.1], we can write (46) as

$$\begin{aligned}
& \mathbb{E} [\langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_k) - \bar{g}(\theta_k, \omega_k) \rangle | \theta_k, \omega_k, \theta_{k-\tau_k}] \\
&\leq 4R_\omega C_\delta d_{TV}(\mu_{\theta_{k-\tau_k}} \otimes \pi_{\theta_{k-\tau_k}} \otimes \mathcal{P}, \mu_{\theta_k} \otimes \pi_{\theta_k} \otimes \mathcal{P}) \\
&\leq C_3 \|\theta_{k-\tau_k} - \theta_k\|_2 \\
&\leq C_3 \sum_{i=k-\tau_k}^{k-1} \alpha_i \|g(x_i, \omega_{i-\tau_i})\|_2 \\
&\leq C_3 C_\delta K_0 \alpha_{k-K_0}, \tag{47}
\end{aligned}$$

where we used the monotonicity of α_k and Assumption 1.

Taking total expectation on both sides of (47) and substituting it into (45) yield

$$\begin{aligned}
\mathbb{E} \|\omega_{k+1} - \omega_{k+1}^*\|_2^2 &\leq (1 - 2\lambda\beta_k) \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 + 2\beta_k \left(C_1 \frac{\alpha_k}{\beta_k} + C_2 K_0 \beta_{k-K_0} \right) \mathbb{E} \|\omega_k - \omega_k^*\|_2 \\
&\quad + 2C_3 C_\delta K_0 \beta_k \alpha_{k-K_0} + C_q \beta_k^2. \tag{48}
\end{aligned}$$

Taking summation on both sides of (48) and rearranging yield

$$\begin{aligned}
2\lambda \sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 &\leq \underbrace{\sum_{k=K_0}^K \frac{1}{\beta_k} \left(\mathbb{E} \|\omega_k - \omega_k^*\|_2^2 - \mathbb{E} \|\omega_{k+1} - \omega_{k+1}^*\|_2^2 \right)}_{I_1} + \underbrace{C_q \sum_{k=K_0}^K \beta_k}_{I_2} \\
&\quad + 2 \underbrace{\sum_{k=K_0}^K 2C_3 C_\delta K_0 \alpha_{k-K_0}}_{I_3} + 2 \underbrace{\sum_{k=K_0}^K \left(C_1 \frac{\alpha_k}{\beta_k} + C_2 K_0 \beta_{k-K_0} \right) \mathbb{E} \|\omega_k - \omega_k^*\|_2}_{I_4}. \tag{49}
\end{aligned}$$

We bound I_1 as

$$\begin{aligned}
I_1 &= \sum_{k=M_K}^K \frac{1}{\beta_k} \left(\mathbb{E} \|\omega_k - \omega_k^*\|_2^2 - \mathbb{E} \|\omega_{k+1} - \omega_{k+1}^*\|_2^2 \right) \\
&= \sum_{k=M_K}^K \left(\frac{1}{\beta_k} - \frac{1}{\beta_{k-1}} \right) \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 + \frac{1}{\beta_{M_K-1}} \mathbb{E} \|\omega_{M_K} - \omega_{M_K}^*\|_2^2 - \frac{1}{\beta_K} \mathbb{E} \|\omega_{K+1} - \omega_{K+1}^*\|_2^2 \\
&\leq \sum_{k=M_K}^K \left(\frac{1}{\beta_k} - \frac{1}{\beta_{k-1}} \right) \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 + \frac{1}{\beta_{M_K-1}} \mathbb{E} \|\omega_{M_K} - \omega_{M_K}^*\|_2^2 \\
&\leq 4R_\omega^2 \left(\sum_{k=M_K}^K \left(\frac{1}{\beta_k} - \frac{1}{\beta_{k-1}} \right) + \frac{1}{\beta_{M_K-1}} \right) = \frac{4R_\omega^2}{\beta_K} = \mathcal{O}(K^{\sigma_2}), \tag{50}
\end{aligned}$$

where the last inequality is due to the fact that

$$\|\omega_k - \omega_\theta^*\|_2 \leq \|\omega_k\|_2 + \|\omega_\theta^*\|_2 \leq 2R_\omega.$$

We bound I_2 as

$$\sum_{k=M_K}^K \beta_k = \sum_{k=M_K}^K \frac{c_2}{(1+k)^{\sigma_2}} = \mathcal{O}(K^{1-\sigma_2}) \tag{51}$$

where the inequality follows from the integration rule $\sum_{k=a}^b k^{-\sigma} \leq \frac{b^{1-\sigma}}{1-\sigma}$.

We bound I_3 as

$$I_3 = \sum_{k=K_0}^K 2C_3 C_\delta K_0 \alpha_{k-K_0} = 2C_3 C_\delta c_1 K_0 \sum_{k=0}^{K-K_0} (1+k)^{-\sigma_1} = \mathcal{O}(K_0 K^{1-\sigma_1}). \tag{52}$$

For the last term I_4 , we have

$$\begin{aligned}
I_4 &= \sum_{k=K_0}^K \left(C_1 \frac{\alpha_k}{\beta_k} + C_2 K_0 \beta_{k-K_0} \right) \mathbb{E} \|\omega_k - \omega_k^*\|_2 \\
&\leq \sqrt{\sum_{k=K_0}^K \left(C_1 \frac{\alpha_k}{\beta_k} + C_2 K_0 \beta_{k-K_0} \right)^2} \sqrt{\sum_{k=K_0}^K (\mathbb{E} \|\omega_k - \omega_k^*\|_2)^2} \\
&\leq \sqrt{\sum_{k=K_0}^K \left(C_1 \frac{\alpha_k}{\beta_k} + C_2 K_0 \beta_{k-K_0} \right)^2} \sqrt{\sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2}, \tag{53}
\end{aligned}$$

where the first inequality follows Cauchy–Schwartz inequality, and the second inequality follows Jensen’s inequality. In (53), we have

$$\begin{aligned}
\sum_{k=K_0}^K \left(C_1 \frac{\alpha_k}{\beta_k} + C_2 K_0 \beta_{k-K_0} \right)^2 &\leq \sum_{k=0}^{K-K_0} \left(C_1 \frac{\alpha_k}{\beta_k} + C_2 K_0 \beta_k \right)^2 \\
&= C_1^2 \sum_{k=0}^{K-K_0} \frac{\alpha_k^2}{\beta_k^2} + 2C_1 C_2 K_0 \sum_{k=0}^{K-K_0} \alpha_k + C_2^2 K_0^2 \sum_{k=0}^{K-K_0} \beta_k^2 \\
&= \mathcal{O} \left(K^{2(\sigma_2-\sigma_1)+1} \right) + \mathcal{O} \left(K_0 K^{-\sigma_1+1} \right) + \mathcal{O} \left(K_0^2 K^{1-2\sigma_2} \right) \tag{54}
\end{aligned}$$

where the first inequality is due to the fact that $\frac{\alpha_k}{\beta_k}$ and β_{k-K_0} are monotonically decreasing.

Substituting (54) into (53) gives

$$I_4 \leq \sqrt{\mathcal{O} \left(K^{2(\sigma_2-\sigma_1)+1} \right) + \mathcal{O} \left(K_0 K^{-\sigma_1+1} \right) + \mathcal{O} \left(K_0^2 K^{1-2\sigma_2} \right)} \sqrt{\sum_{k=M_K}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2}. \tag{55}$$

Substituting (50), (51), (52) and (55) into (49), and dividing both sides of (49) by $K - K_0 + 1$ give

$$\begin{aligned}
& 2\lambda \frac{1}{K - K_0 + 1} \sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 \\
& \leq \frac{\sqrt{\mathcal{O}(K^{2(\sigma_2 - \sigma_1) + 1}) + \mathcal{O}(K_0 K^{-\sigma_1 + 1}) + \mathcal{O}(K_0^2 K^{1 - 2\sigma_2})}}{K - K_0 + 1} \sqrt{\sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2} \\
& \quad + \mathcal{O}\left(\frac{1}{K^{1 - \sigma_2}}\right) + \mathcal{O}\left(\frac{1}{K^{\sigma_2}}\right) + \mathcal{O}\left(\frac{K_0}{K^{\sigma_1}}\right).
\end{aligned} \tag{56}$$

We define the following functions:

$$\begin{aligned}
T_1(K) &:= \frac{1}{K - K_0 + 1} \sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2, \\
T_2(K) &:= \mathcal{O}\left(\frac{1}{K^{1 - \sigma_2}}\right) + \mathcal{O}\left(\frac{1}{K^{\sigma_2}}\right) + \mathcal{O}\left(\frac{K_0}{K^{\sigma_1}}\right), \\
T_3(K) &:= \frac{\mathcal{O}(K^{2(\sigma_2 - \sigma_1) + 1}) + \mathcal{O}(K_0 K^{-\sigma_1 + 1}) + \mathcal{O}(K_0^2 K^{1 - 2\sigma_2})}{K - K_0 + 1}.
\end{aligned}$$

Then (56) can be written as:

$$T_1(K) - \frac{1}{2\lambda} \sqrt{T_1(K)} \sqrt{T_3(K)} \leq \frac{1}{2\lambda} T_2(K).$$

Solving this quadratic inequality in terms of $T_1(K)$, we obtain

$$T_1(K) \leq \frac{1}{\lambda} T_2(K) + \frac{1}{2\lambda^2} T_3(K), \tag{57}$$

which implies

$$\begin{aligned}
& \frac{1}{K - K_0 + 1} \sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 \\
& = \mathcal{O}\left(\frac{1}{K^{1 - \sigma_2}}\right) + \mathcal{O}\left(\frac{1}{K^{2(\sigma_1 - \sigma_2)}}\right) + \mathcal{O}\left(\frac{K_0^2}{K^{2\sigma_2}}\right) + \mathcal{O}\left(\frac{K_0}{K^{\sigma_1}}\right) + \mathcal{O}\left(\frac{1}{K^{\sigma_2}}\right).
\end{aligned}$$

We further have

$$\begin{aligned}
\frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 & \leq \frac{1}{K} \left(\sum_{k=1}^{K_0-1} 4R_\omega^2 + \sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 \right) \\
& = \frac{K_0 - 1}{K} 4R_\omega^2 + \frac{K - K_0 + 1}{K} \frac{1}{K - K_0 + 1} \sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 \\
& = \mathcal{O}\left(\frac{K_0}{K}\right) + \mathcal{O}\left(\frac{1}{K - K_0 + 1} \sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2\right) \\
& = \mathcal{O}\left(\frac{1}{K - K_0 + 1} \sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2\right)
\end{aligned} \tag{58}$$

which completes the proof. \square

B.2 PROOF OF THEOREM 2

We first clarify the notations:

$$\begin{aligned}
x &:= (s, a, s'), \\
\hat{\delta}(x, \omega) &:= r(s, a, s') + \gamma \phi(s')^\top \omega - \phi(s)^\top \omega, \\
\delta(x, \theta) &:= r(s, a, s') + \gamma V_{\pi_\theta}(s') - V_{\pi_\theta}(s),
\end{aligned}$$

The update in Algorithm 1 can be written compactly as:

$$\theta_{k+1} = \theta_k + \alpha_k \hat{\delta}(x_{(k)}, \omega_{k-\tau_k}) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}). \quad (59)$$

For brevity, we use ω_k^* as shorthand notation of $\omega_{\theta_k}^*$ in this proof. Then we are ready to give the convergence proof.

Proof. From L_J -Lipschitz of policy gradient shown in Proposition 1, we have:

$$\begin{aligned} J(\theta_{k+1}) &\geq J(\theta_k) + \langle \nabla J(\theta_k), \theta_{k+1} - \theta_k \rangle - \frac{L_J}{2} \|\theta_{k+1} - \theta_k\|_2^2 \\ &= J(\theta_k) + \alpha_k \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_{k-\tau_k}) - \hat{\delta}(x_{(k)}, \omega_k^*) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\ &\quad + \alpha_k \left\langle \nabla J(\theta_k), \hat{\delta}(x_{(k)}, \omega_k^*) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle - \frac{L_J}{2} \alpha_k^2 \|\hat{\delta}(x_{(k)}, \omega_{k-\tau_k}) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)})\|_2^2 \\ &\geq J(\theta_k) + \alpha_k \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_{k-\tau_k}) - \hat{\delta}(x_{(k)}, \omega_k^*) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\ &\quad + \alpha_k \left\langle \nabla J(\theta_k), \hat{\delta}(x_{(k)}, \omega_k^*) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle - \frac{L_J}{2} C_p^2 \alpha_k^2, \end{aligned}$$

where the last inequality follows the definition of C_p in (44).

Taking expectation on both sides of the last inequality yields

$$\begin{aligned} \mathbb{E}[J(\theta_{k+1})] &\geq \mathbb{E}[J(\theta_k)] + \underbrace{\alpha_k \mathbb{E} \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_{k-\tau_k}) - \hat{\delta}(x_{(k)}, \omega_k^*) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle}_{I_1} \\ &\quad + \underbrace{\alpha_k \mathbb{E} \left\langle \nabla J(\theta_k), \hat{\delta}(x_{(k)}, \omega_k^*) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle}_{I_2} - \frac{L_J}{2} C_p^2 \alpha_k^2. \end{aligned} \quad (60)$$

We first decompose I_1 as

$$\begin{aligned} I_1 &= \mathbb{E} \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_{k-\tau_k}) - \hat{\delta}(x_{(k)}, \omega_k^*) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\ &= \mathbb{E} \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_{k-\tau_k}) - \hat{\delta}(x_{(k)}, \omega_k) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\ &\quad + \mathbb{E} \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_k) - \hat{\delta}(x_{(k)}, \omega_k^*) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle. \end{aligned}$$

$I_1^{(1)}$ $I_1^{(2)}$

We bound $I_1^{(1)}$ as

$$\begin{aligned} I_1^{(1)} &= \mathbb{E} \left\langle \nabla J(\theta_k), \left(\gamma \phi(s'_{(k)}) - \phi(s_{(k)}) \right)^\top (\omega_{k-\tau_k} - \omega_k) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\ &\geq -\mathbb{E} \left[\|\nabla J(\theta_k)\|_2 \|\gamma \phi(s'_{(k)}) - \phi(s_{(k)})\|_2 \|\omega_k - \omega_{k-\tau_k}\|_2 \|\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)})\|_2 \right] \\ &\geq -(1 + \gamma) C_\psi \mathbb{E} [\|\nabla J(\theta_k)\|_2 \|\omega_k - \omega_{k-\tau_k}\|_2] \\ &\geq -(1 + \gamma) C_\psi C_\delta K_0 \beta_{k-1} \mathbb{E} \|\nabla J(\theta_k)\|_2, \end{aligned}$$

where the last inequality follows

$$\begin{aligned} \|\omega_k - \omega_{k-\tau_k}\|_2 &= \left\| \sum_{i=k-\tau_k}^{k-1} (\omega_{i+1} - \omega_i) \right\|_2 \\ &\leq \sum_{i=k-\tau_k}^{k-1} \|\beta_i g(x_i, \omega_{i-\tau_i})\|_2 \\ &\leq \beta_{k-1} \sum_{i=k-\tau_k}^{k-1} \|g(x_i, \omega_{i-\tau_i})\|_2 \\ &\leq \beta_{k-1} K_0 C_\delta, \end{aligned}$$

where the second inequality is due to the monotonicity of step size, and the third one follows (37).

Then we bound $I_1^{(2)}$ as

$$\begin{aligned}
I_1^{(2)} &= \mathbb{E} \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_k) - \hat{\delta}(x_{(k)}, \omega_k^*) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\
&= -\mathbb{E} \left\langle \nabla J(\theta_k), \left(\gamma \phi(s'_{(k)}) - \phi(s_{(k)}) \right)^\top (\omega_k^* - \omega_k) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\
&\geq -\mathbb{E} \left[\|\nabla J(\theta_k)\|_2 \|\gamma \phi(s'_{(k)}) - \phi(s_{(k)})\|_2 \|\omega_k - \omega_k^*\|_2 \|\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)})\|_2 \right] \\
&\geq -(1 + \gamma) C_\psi \mathbb{E} [\|\nabla J(\theta_k)\|_2 \|\omega_k - \omega_k^*\|_2].
\end{aligned}$$

Collecting lower bounds of $I_1^{(1)}$ and $I_1^{(2)}$ gives

$$\begin{aligned}
I_1 &\geq -(1 + \gamma) C_\psi \mathbb{E} [\|\nabla J(\theta_k)\|_2 (C_\delta K_0 \beta_{k-1} + \|\omega_k - \omega_k^*\|_2)] \\
&= -(1 + \gamma) C_\psi \sqrt{(\mathbb{E} [\|\nabla J(\theta_k)\|_2 (C_\delta K_0 \beta_{k-1} + \|\omega_k - \omega_k^*\|_2)]^2} \\
&\geq -(1 + \gamma) C_\psi \sqrt{\mathbb{E} \|\nabla J(\theta_k)\|_2^2 \mathbb{E} [(C_\delta K_0 \beta_{k-1} + \|\omega_k - \omega_k^*\|_2)^2]} \\
&\geq -(1 + \gamma) C_\psi \sqrt{\mathbb{E} \|\nabla J(\theta_k)\|_2^2} \sqrt{\mathbb{E} [2C_\delta^2 K_0^2 \beta_{k-1}^2 + 2\|\omega_k - \omega_k^*\|_2^2]} \\
&= -\sqrt{2}(1 + \gamma) C_\psi \sqrt{\mathbb{E} \|\nabla J(\theta_k)\|_2^2} \sqrt{C_\delta^2 K_0^2 \beta_{k-1}^2 + \mathbb{E} \|\omega_k - \omega_k^*\|_2^2}, \tag{61}
\end{aligned}$$

where the second inequality follows Cauchy-Schwartz inequality, and the third inequality follows Young's inequality.

Now we consider I_2 . We first decompose I_2 as

$$\begin{aligned}
I_2 &= \mathbb{E} \left\langle \nabla J(\theta_k), \hat{\delta}(x_{(k)}, \omega_k^*) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\
&= \underbrace{\mathbb{E} \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_k^*) - \hat{\delta}(x_{(k)}, \omega_{k-\tau_k}^*) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle}_{I_2^{(1)}} \\
&\quad + \underbrace{\mathbb{E} \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_{k-\tau_k}^*) - \delta(x_{(k)}, \theta_{k-\tau_k}) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle}_{I_2^{(2)}} \\
&\quad + \underbrace{\mathbb{E} \left\langle \nabla J(\theta_k), \delta(x_{(k)}, \theta_{k-\tau_k}) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \nabla J(\theta_k) \right\rangle}_{I_2^{(3)}} + \|\nabla J(\theta_k)\|_2^2.
\end{aligned}$$

We bound $I_2^{(1)}$ as

$$\begin{aligned}
I_2^{(1)} &= \mathbb{E} \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_k^*) - \hat{\delta}(x_{(k)}, \omega_{k-\tau_k}^*) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\
&= \mathbb{E} \left\langle \nabla J(\theta_k), \left(\gamma \phi(s'_{(k)}) - \phi(s_{(k)}) \right)^\top (\omega_k^* - \omega_{k-\tau_k}^*) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\
&\geq -\mathbb{E} \left[\|\nabla J(\theta_k)\|_2 \|\gamma \phi(s'_{(k)}) - \phi(s_{(k)})\|_2^\top \|\omega_k^* - \omega_{k-\tau_k}^*\|_2 \|\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)})\|_2 \right] \\
&\geq -L_V C_\psi (1 + \gamma) \mathbb{E} \|\omega_k^* - \omega_{k-\tau_k}^*\|_2 \\
&\geq -L_V L_\omega C_\psi (1 + \gamma) \mathbb{E} \|\theta_k - \theta_{k-\tau_k}\|_2 \\
&\geq -L_V L_\omega C_\psi C_p (1 + \gamma) K_0 \alpha_{k-K_0},
\end{aligned}$$

where the second last inequality follows from Proposition 2 and the last inequality uses (44) as

$$\begin{aligned}
\|\theta_k - \theta_{k-\tau_k}\|_2 &\leq \sum_{i=k-\tau_k}^{k-1} \|\theta_{i+1} - \theta_i\|_2 \\
&= \sum_{i=k-\tau_k}^{k-1} \alpha_i \|\hat{\delta}(x_i, \omega_{i-\tau_i}) \psi_{\theta_{i-\tau_i}}(s_i, a_i)\|_2 \\
&\leq \sum_{i=k-\tau_k}^{k-1} \alpha_{k-\tau_k} C_p \leq C_p K_0 \alpha_{k-K_0}.
\end{aligned} \tag{62}$$

We bound $I_2^{(2)}$ as

$$\begin{aligned}
I_2^{(2)} &= \mathbb{E} \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_{k-\tau_k}^*) - \delta(x_{(k)}, \theta_{k-\tau_k}) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\
&\geq -\mathbb{E} \left[\|\nabla J(\theta_k)\|_2 \left| \hat{\delta}(x_{(k)}, \omega_{k-\tau_k}^*) - \delta(x_{(k)}, \theta_{k-\tau_k}) \right| \|\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)})\|_2 \right] \\
&\geq -L_V C_\psi \mathbb{E} \left| \hat{\delta}(x_{(k)}, \omega_{k-\tau_k}^*) - \delta(x_{(k)}, \theta_{k-\tau_k}) \right| \\
&= -L_V C_\psi \mathbb{E} \left| \gamma \left(\phi(s'_{(k)})^\top \omega_{k-\tau_k}^* - V_{\pi_{\theta_{k-\tau_k}}}(s'_{(k)}) \right) + V_{\pi_{\theta_{k-\tau_k}}}(s_{(k)}) - \phi(s_{(k)})^\top \omega_{k-\tau_k}^* \right| \\
&\geq -L_V C_\psi \left(\gamma \mathbb{E} \left| \phi(s'_{(k)})^\top \omega_{k-\tau_k}^* - V_{\pi_{\theta_{k-\tau_k}}}(s'_{(k)}) \right| + \mathbb{E} \left| V_{\pi_{\theta_{k-\tau_k}}}(s_{(k)}) - \phi(s_{(k)})^\top \omega_{k-\tau_k}^* \right| \right) \\
&\geq -L_V C_\psi \left(\gamma \sqrt{\mathbb{E} \left| \phi(s'_{(k)})^\top \omega_{k-\tau_k}^* - V_{\pi_{\theta_{k-\tau_k}}}(s'_{(k)}) \right|^2} + \sqrt{\mathbb{E} \left| V_{\pi_{\theta_{k-\tau_k}}}(s_{(k)}) - \phi(s_{(k)})^\top \omega_{k-\tau_k}^* \right|^2} \right) \\
&\geq -L_V C_\psi (1 + \gamma) \epsilon_{app}.
\end{aligned}$$

We bound $I_2^{(3)}$ as

$$\begin{aligned}
I_2^{(3)} &= \mathbb{E} \left\langle \nabla J(\theta_k), \delta(x_{(k)}, \theta_{k-\tau_k}) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \nabla J(\theta_k) \right\rangle \\
&= \mathbb{E} \left[\mathbb{E} \left[\left\langle \nabla J(\theta_k), \delta(x_{(k)}, \theta_{k-\tau_k}) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \nabla J(\theta_k) \right\rangle \middle| \theta_{k-\tau_k}, \theta_k \right] \right] \\
&= \mathbb{E} \left\langle \nabla J(\theta_k), \mathbb{E} \left[\delta(x_{(k)}, \theta_{k-\tau_k}) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \middle| \theta_{k-\tau_k}, \theta_k \right] - \nabla J(\theta_k) \right\rangle \\
&= \mathbb{E} \left\langle \nabla J(\theta_k), \mathbb{E}_{\substack{s_{(k)} \sim \mu_{\theta_{k-\tau_k}} \\ a_{(k)} \sim \pi_{\theta_{k-\tau_k}}}} \left[A_{\pi_{\theta_{k-\tau_k}}}(s_{(k)}, a_{(k)}) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right] - \nabla J(\theta_k) \right\rangle \tag{63}
\end{aligned}$$

where we used the fact that

$$\begin{aligned}
&\mathbb{E} \left[\delta(x_{(k)}, \theta_{k-\tau_k}) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \middle| \theta_{k-\tau_k}, \theta_k \right] \\
&= \mathbb{E}_{\substack{s_{(k)} \sim \mu_{\theta_{k-\tau_k}} \\ a_{(k)} \sim \pi_{\theta_{k-\tau_k}} \\ s'_{(k)} \sim \mathcal{P}}} \left[\left(r(s_{(k)}, a_{(k)}, s'_{(k)}) + \gamma V_{\pi_{\theta_{k-\tau_k}}}(s'_{(k)}) - V_{\pi_{\theta_{k-\tau_k}}}(s_{(k)}) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \middle| \theta_{k-\tau_k}, \theta_k \right] \\
&= \mathbb{E}_{\substack{s_{(k)} \sim \mu_{\theta_{k-\tau_k}} \\ a_{(k)} \sim \pi_{\theta_{k-\tau_k}}}} \left[\left(\mathbb{E}_{s'_{(k)} \sim \mathcal{P}} \left[r(s_{(k)}, a_{(k)}, s'_{(k)}) + \gamma V_{\pi_{\theta_{k-\tau_k}}}(s'_{(k)}) \right] - V_{\pi_{\theta_{k-\tau_k}}}(s_{(k)}) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \middle| \theta_{k-\tau_k}, \theta_k \right] \\
&= \mathbb{E}_{\substack{s_{(k)} \sim \mu_{\theta_{k-\tau_k}} \\ a_{(k)} \sim \pi_{\theta_{k-\tau_k}}}} \left[\left(Q_{\pi_{\theta_{k-\tau_k}}}(s_{(k)}, a_{(k)}) - V_{\pi_{\theta_{k-\tau_k}}}(s_{(k)}) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \middle| \theta_{k-\tau_k}, \theta_k \right] \\
&= \mathbb{E}_{\substack{s_{(k)} \sim \mu_{\theta_{k-\tau_k}} \\ a_{(k)} \sim \pi_{\theta_{k-\tau_k}}}} \left[A_{\pi_{\theta_{k-\tau_k}}}(s_{(k)}, a_{(k)}) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \middle| \theta_{k-\tau_k}, \theta_k \right].
\end{aligned}$$

We can further decompose (63) as

$$I_2^{(3)} = \mathbb{E} \left\langle \nabla J(\theta_k), \mathbb{E}_{\substack{s(k) \sim \mu_{\theta_{k-\tau_k}} \\ a(k) \sim \pi_{\theta_{k-\tau_k}}} \left[A_{\pi_{\theta_{k-\tau_k}}}(s(k), a(k)) \psi_{\theta_{k-\tau_k}}(s(k), a(k)) \right] - \nabla J(\theta_{k-\tau_k}) \right\rangle \\ + \mathbb{E} \left\langle \nabla J(\theta_k), \nabla J(\theta_{k-\tau_k}) - \nabla J(\theta_k) \right\rangle. \quad (64)$$

The first term in (64) corresponds to the bias introduced by the mismatch between the stationary distribution $\mu_{\theta_{k-\tau_k}}$ and the visitation measure $d_{\theta_{k-\tau_k}}$ in discounted MDP. It can be bounded as

$$\begin{aligned} & \left\langle \nabla J(\theta_k), \mathbb{E}_{\substack{s(k) \sim \mu_{\theta_{k-\tau_k}} \\ a(k) \sim \pi_{\theta_{k-\tau_k}}} \left[A_{\pi_{\theta_{k-\tau_k}}}(s(k), a(k)) \psi_{\theta_{k-\tau_k}}(s(k), a(k)) \right] - \nabla J(\theta_{k-\tau_k}) \right\rangle \\ & \geq -L_V \left\| \mathbb{E}_{\substack{s(k) \sim \mu_{\theta_{k-\tau_k}} \\ a(k) \sim \pi_{\theta_{k-\tau_k}}} \left[A_{\pi_{\theta_{k-\tau_k}}}(s(k), a(k)) \psi_{\theta_{k-\tau_k}}(s(k), a(k)) \right] - \nabla J(\theta_{k-\tau_k}) \right\|_2 \\ & \geq -L_V \sup_{s,a} \|A_{\pi_{\theta_{k-\tau_k}}}(s, a) \psi_{\theta_{k-\tau_k}}(s, a)\|_2 \|\mu_{\theta_{k-\tau_k}} - d_{\theta_{k-\tau_k}}\|_{TV} \\ & \geq -4R_{\max} C_\psi L_V \left(\log_\rho \kappa^{-1} + \frac{1}{1-\rho} \right) (1-\gamma), \end{aligned}$$

where the last inequality follows Lemma 2 and the fact that $|A_\pi(s, a)| \leq r_{\max}/(1-\gamma) = R_{\max}$.

The second term in (64) corresponds to the bias introduced by delay. It can be bounded as

$$\begin{aligned} \left\langle \nabla J(\theta_k), \nabla J(\theta_{k-\tau_k}) - \nabla J(\theta_k) \right\rangle & \geq -\|\nabla J(\theta_k)\|_2 \|\nabla J(\theta_{k-\tau_k}) - \nabla J(\theta_k)\|_2 \\ & \geq -L_V L_J \|\theta_{k-\tau_k} - \theta_k\|_2 \\ & \geq -L_V L_J C_p K_0 \alpha_{k-K_0}, \end{aligned}$$

where the second last inequality is due to L_J -Lipschitz of policy gradient shown in Proposition 1, and the last inequality follows (62).

Collecting the lower bounds gives

$$I_2^{(3)} \geq -\epsilon_{sp} - L_V L_J C_p K_0 \alpha_{k-K_0},$$

where $\epsilon_{sp} = 4R_{\max} C_\psi L_V \left(\log_\rho \kappa^{-1} + \frac{1}{1-\rho} \right) (1-\gamma)$ is the sampling error.

Collecting lower bounds of $I_2^{(1)}$, $I_2^{(2)}$ and $I_2^{(3)}$ gives

$$I_2 \geq -D_1 K_0 \alpha_{k-K_0} - L_V C_\psi (1+\gamma) \epsilon_{app} - \epsilon_{sp}, \quad (65)$$

where constant $D_1 := L_V L_\omega C_\psi C_p (1+\gamma) + L_V L_J C_p$.

Substituting (61) and (65) into (60) yields

$$\begin{aligned} \mathbb{E}[J(\theta_{k+1})] & \geq \mathbb{E}[J(\theta_k)] - \alpha_k \sqrt{2}(1+\gamma) C_\psi \sqrt{\mathbb{E} \|\nabla J(\theta_k)\|_2^2} \sqrt{C_\delta^2 K_0^2 \beta_{k-1}^2 + \mathbb{E} \|\omega_k - \omega_k^*\|_2^2} \\ & \quad - \alpha_k D_1 K_0 \alpha_{k-K_0} - \alpha_k L_V C_\psi (1+\gamma) \epsilon_{app} - \alpha_k \epsilon_{sp} + \alpha_k \|\nabla J(\theta_k)\|_2^2 - \frac{L_J}{2} C_p^2 \alpha_k^2. \end{aligned} \quad (66)$$

Dividing both sides of (66) by α_k , then rearranging and taking summation on both sides give

$$\begin{aligned} \sum_{k=K_0}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 & \leq \underbrace{\sum_{k=K_0}^K \frac{1}{\alpha_k} (\mathbb{E}[J(\theta_{k+1})] - \mathbb{E}[J(\theta_k)])}_{I_3} + \underbrace{\sum_{k=K_0}^K \left(D_1 K_0 \alpha_{k-K_0} + \frac{L_J}{2} C_p^2 \alpha_k \right)}_{I_4} \\ & \quad + \underbrace{\sqrt{2}(1+\gamma) C_\psi \sum_{k=K_0}^K \sqrt{\mathbb{E} \|\nabla J(\theta_k)\|_2^2} \sqrt{C_\delta^2 K_0^2 \beta_{k-1}^2 + \mathbb{E} \|\omega_k - \omega_k^*\|_2^2}}_{I_5} \\ & \quad + (K - K_0 + 1)(1+\gamma) L_V C_\psi \epsilon_{app} + (K - K_0 + 1) \epsilon_{sp}, \end{aligned} \quad (67)$$

We bound I_3 as

$$\begin{aligned}
I_3 &= \sum_{k=M_K}^K \frac{1}{\alpha_k} (\mathbb{E}[J(\theta_{k+1})] - \mathbb{E}[J(\theta_k)]) \\
&= \sum_{k=M_K}^K \left(\frac{1}{\alpha_{k-1}} - \frac{1}{\alpha_k} \right) \mathbb{E}[J(\theta_k)] - \frac{1}{\alpha_{M_K-1}} \mathbb{E}[J(\theta_{M_K})] + \frac{1}{\alpha_K} \mathbb{E}[J(\theta_{K+1})] \\
&\leq \frac{1}{\alpha_K} \mathbb{E}[J(\theta_{K+1})] \\
&\leq R_{\max} \frac{1}{\alpha_K} = \mathcal{O}(K^{\sigma_1}),
\end{aligned} \tag{68}$$

where the first inequality is due to the α_k is monotonic decreasing and positive, and last inequality is due to $V_{\pi_\theta}(s) \leq R_{\max}$ for any $s \in \mathcal{S}$ and π_θ .

We bound I_4 as

$$I_4 = \sum_{k=K_0}^K \left(D_1 K_0 \alpha_{k-K_0} + \frac{L_J}{2} C_p^2 \alpha_k \right) \leq \sum_{k=0}^{K-K_0} \left(D_1 K_0 \alpha_k + \frac{L_J}{2} C_p^2 \alpha_k \right) = \mathcal{O}(K_0 K^{1-\sigma_1}).$$

We bound I_5 as

$$\begin{aligned}
I_5 &= \sum_{k=K_0}^K \sqrt{\mathbb{E} \|\nabla J(\theta_k)\|_2^2} \sqrt{C_\delta^2 K_0^2 \beta_{k-1}^2 + \mathbb{E} \|\omega_k - \omega_k^*\|_2^2} \\
&\leq \sqrt{\sum_{k=K_0}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2} \sqrt{\sum_{k=K_0}^K (C_\delta^2 K_0^2 \beta_{k-1}^2 + \mathbb{E} \|\omega_k - \omega_k^*\|_2^2)} \\
&= \sqrt{\sum_{k=K_0}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2} \sqrt{C_\delta^2 K_0^2 \sum_{k=K_0}^K \beta_{k-1}^2 + \sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2},
\end{aligned} \tag{69}$$

where the first inequality follows Cauchy-Schwartz inequality. In (69), we have

$$\sum_{k=K_0}^K \beta_{k-1}^2 \leq \sum_{k=0}^{K-K_0} \beta_k^2 = \sum_{k=0}^{K-K_0} c_2^2 (1+k)^{-2\sigma_2} = \mathcal{O}(K^{1-2\sigma_2}).$$

Substituting the last equality into (69) gives

$$I_5 \leq \sqrt{\sum_{k=K_0}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2} \sqrt{\mathcal{O}(K_0^2 K^{1-2\sigma_2}) + \sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2}. \tag{70}$$

Dividing both sides of (66) by $K - K_0 + 1$ and collecting upper bounds of I_3 , I_4 and I_5 give

$$\begin{aligned}
&\frac{1}{K - K_0 + 1} \sum_{k=K_0}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 \\
&\leq \frac{\sqrt{2}(1+\gamma)C_\psi}{K - K_0 + 1} \sqrt{\sum_{k=K_0}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2} \sqrt{\mathcal{O}(K_0^2 K^{1-2\sigma_2}) + \sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2} \\
&\quad + \mathcal{O}\left(\frac{1}{K^{1-\sigma_1}}\right) + \mathcal{O}\left(\frac{K_0}{K^{\sigma_1}}\right) + \mathcal{O}(\epsilon_{app} + \epsilon_{sp}).
\end{aligned} \tag{71}$$

Define the following functions

$$\begin{aligned} T_4(K) &:= \frac{1}{K - K_0 + 1} \sum_{k=K_0}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2, \\ T_5(K) &:= \frac{1}{K - K_0 + 1} \left(\mathcal{O}(K_0^2 K^{1-2\sigma_2}) + \sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 \right), \\ T_6(K) &:= \mathcal{O}\left(\frac{1}{K^{1-\sigma_1}}\right) + \mathcal{O}\left(\frac{K_0}{K^{\sigma_1}}\right) + \mathcal{O}(\epsilon_{app} + \epsilon_{sp}). \end{aligned}$$

Then (90) can be rewritten as

$$T_4(K) \leq T_6(K) + \sqrt{2}(1 + \gamma)C_\psi \sqrt{T_4(K)} \sqrt{T_5(K)}.$$

Solving this quadratic inequality in terms of $T_4(K)$, we obtain

$$T_4(K) \leq 2T_6(K) + 4(1 + \gamma)^2 C_\psi^2 T_5(K), \quad (72)$$

which implies

$$\begin{aligned} &\frac{1}{K - K_0 + 1} \sum_{k=K_0}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 \\ &= \mathcal{O}\left(\frac{1}{K^{1-\sigma_1}}\right) + \mathcal{O}\left(\frac{K_0}{K^{\sigma_1}}\right) + \mathcal{O}\left(\frac{K_0^2}{K^{2\sigma_2}}\right) + \mathcal{O}\left(\frac{1}{K - K_0 + 1} \sum_{k=K_0}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2\right) + \mathcal{O}(\epsilon_{app} + \epsilon_{sp}) \end{aligned}$$

We further have

$$\begin{aligned} \frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 &\leq \frac{1}{K} \left(\sum_{k=1}^{K_0-1} L_V^2 + \sum_{k=K_0}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 \right) \\ &= \frac{K_0 - 1}{K} L_V^2 + \frac{K - K_0 + 1}{K} \frac{1}{K - K_0 + 1} \sum_{k=K_0}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 \\ &= \mathcal{O}\left(\frac{K_0}{K}\right) + \mathcal{O}\left(\frac{1}{K - K_0 + 1} \sum_{k=K_0}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2\right) \\ &= \mathcal{O}\left(\frac{1}{K - K_0 + 1} \sum_{k=K_0}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2\right) \end{aligned} \quad (73)$$

which completes the proof. \square

B.3 PROOF OF THEOREM 3

Given the definition in Section B.1, we now give the convergence proof of critic update in Algorithm 1 with linear function approximation and Markovian sampling.

By following the derivation of (45), we have

$$\begin{aligned} \mathbb{E} \|\omega_{k+1} - \omega_{k+1}^*\|_2^2 &\leq (1 - 2\lambda\beta_k) \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 + 2\beta_k \left(C_1 \frac{\alpha_k}{\beta_k} + C_2 K_0 \beta_{k-K_0} \right) \mathbb{E} \|\omega_k - \omega_k^*\|_2 \\ &\quad + 2\beta_k \mathbb{E} \langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_k) - \bar{g}(\theta_k, \omega_k) \rangle + C_q \beta_k^2, \end{aligned} \quad (74)$$

where $C_1 := C_p L_\omega$, $C_2 := C_\delta(1 + \gamma)$ and $C_q := 2C_\delta^2 + 2L_\omega^2 C_p^2 \max_{(k)} \frac{\alpha_k^2}{\beta_k^2} = 2C_\delta^2 + 2L_\omega^2 C_p^2 \frac{c_1^2}{c_2^2}$.

Now we consider the third item in the last inequality. For some $m \in \mathbb{N}^+$, we define $M := (K_0 + 1)m + K_0$. Following Lemma 5 (to be presented in Section C.1), for some $d_m \leq M$ and

positive constants C_4, C_5, C_6, C_7 , we have

$$\begin{aligned}
& \mathbb{E} \langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_k) - \bar{g}(\theta_k, \omega_k) \rangle \\
& \leq C_4 \mathbb{E} \|\theta_k - \theta_{k-d_m}\|_2 + C_5 \sum_{i=\tau_k}^{d_m} \mathbb{E} \|\theta_{k-i} - \theta_{k-d_m}\|_2 + C_6 \mathbb{E} \|\omega_k - \omega_{k-d_m}\|_2 + C_7 \kappa \rho^{m-1} \\
& \leq C_4 \sum_{i=k-d_m}^{k-1} \mathbb{E} \|\theta_{i+1} - \theta_i\|_2 + C_5 \sum_{i=\tau_k}^{d_m-1} \sum_{j=k-d_m}^{k-i-1} \mathbb{E} \|\theta_{j+1} - \theta_j\|_2 + C_6 \sum_{i=k-d_m}^{k-1} \mathbb{E} \|\omega_{i+1} - \omega_i\|_2 + C_7 \kappa \rho^{m-1} \\
& \leq C_4 \sum_{i=k-d_m}^{k-1} \alpha_i C_p + C_5 \sum_{i=\tau_k}^{d_m-1} \sum_{j=k-d_m}^{k-i-1} \alpha_j C_p + C_6 \sum_{i=k-d_m}^{k-1} \beta_i C_\delta + C_7 \kappa \rho^{m-1} \\
& \leq C_4 \alpha_{k-d_m} \sum_{i=k-d_m}^{k-1} C_p + C_5 \alpha_{k-d_m} \sum_{i=\tau_k}^{d_m-1} \sum_{j=k-d_m}^{k-i-1} C_p + C_6 \beta_{k-d_m} \sum_{i=k-d_m}^{k-1} C_\delta + C_7 \kappa \rho^{m-1} \\
& \leq C_4 d_m C_p \alpha_{k-d_m} + C_5 (d_m - \tau_k)^2 C_p \alpha_{k-d_m} + C_6 d_m C_\delta \beta_{k-d_m} + C_7 \kappa \rho^{m-1} \\
& \leq (C_4 M + C_5 M^2) C_p \alpha_{k-M} + C_6 M C_\delta \beta_{k-M} + C_7 \kappa \rho^{m-1}, \tag{75}
\end{aligned}$$

where the third last inequality is due to the monotonicity of step size, and the last inequality is due to $\tau_k \geq 0$ and $d_m \leq M$.

Further letting $m = m_K$ which is defined in (22) yields

$$\begin{aligned}
& \mathbb{E} \langle \omega_k - \omega_k^*, g(x_{(k)}, \omega_k) - \bar{g}(\theta_k, \omega_k) \rangle \\
& = (C_4 M_K + C_5 M_K^2) C_p \alpha_{k-M_K} + C_6 C_\delta M_K \beta_{k-M_K} + C_7 \kappa \rho^{m_K-1} \\
& \leq (C_4 M_K + C_5 M_K^2) C_p \alpha_{k-M_K} + C_6 C_\delta M_K \beta_{k-M_K} + C_7 \alpha_K, \tag{76}
\end{aligned}$$

where $M_K = (K_0 + 1)m_K + K_0$, and the last inequality follows the definition of $m_K = \log K$.

Substituting (76) into (74), then rearranging and summing up both sides over $k = M_K, \dots, K$ yield

$$\begin{aligned}
2\lambda \sum_{k=M_K}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 & \leq \underbrace{\sum_{k=M_K}^K \frac{1}{\beta_k} \left(\mathbb{E} \|\omega_k - \omega_k^*\|_2^2 - \mathbb{E} \|\omega_{k+1} - \omega_{k+1}^*\|_2^2 \right)}_{I_1} + C_q \underbrace{\sum_{k=M_K}^K \beta_k}_{I_2} \\
& \quad + 2 \underbrace{\sum_{k=M_K}^K ((C_4 M_K + C_5 M_K^2) C_p \alpha_{k-M_K} + C_6 C_\delta M_K \beta_{k-M_K} + C_7 \alpha_K)}_{I_3} \\
& \quad + 2 \underbrace{\sum_{k=M_K}^K \left(C_1 \frac{\alpha_k}{\beta_k} + C_2 K_0 \beta_{k-K_0} \right) \mathbb{E} \|\omega_k - \omega_k^*\|_2}_{I_4}, \tag{77}
\end{aligned}$$

where I_1, I_2 and I_4 have already been bounded in (50), (51) and (55) respectively.

We bound I_3 as

$$\begin{aligned}
I_3 & = (C_4 M_K + C_5 M_K^2) C_p \sum_{k=M_K}^K \alpha_k + C_6 C_\delta M_K \sum_{k=M_K}^K \beta_k + C_7 \alpha_K \sum_{k=M_K}^K 1 \\
& \leq (C_4 M_K + C_5 M_K^2) C_p c_1 \frac{K^{1-\sigma_1}}{1-\sigma_1} + C_6 C_\delta M_K c_2 \frac{K^{1-\sigma_2}}{1-\sigma_2} + C_7 c_1 K(1+K)^{-\sigma_1} \\
& = \mathcal{O}((K_0^2 \log^2 K) K^{1-\sigma_1}) + \mathcal{O}((K_0 \log K) K^{1-\sigma_2}), \tag{78}
\end{aligned}$$

where the last inequality follows from the integration rule $\sum_{k=a}^b k^{-\sigma} \leq \frac{b^{1-\sigma}}{1-\sigma}$, and the last equality is due to $\mathcal{O}(M_K) = \mathcal{O}(K_0 m_K) = \mathcal{O}(K_0 \log K)$.

Collecting the upper bounds of I_1 , I_2 , I_3 and I_4 , and dividing both sides of (77) by $K - M_K + 1$ yield

$$\begin{aligned}
& 2\lambda \frac{1}{K - M_K + 1} \sum_{k=M_K}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 \\
& \leq \frac{\sqrt{\mathcal{O}(K^{2(\sigma_2 - \sigma_1 + 1)}) + \mathcal{O}(K_0 K^{-\sigma_1 + 1}) + \mathcal{O}(K_0^2 K^{1 - 2\sigma_2})}}{K - M_K + 1} \sqrt{\sum_{k=M_K}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2} \\
& \quad + \mathcal{O}\left(\frac{1}{K^{1 - \sigma_2}}\right) + \mathcal{O}\left(\frac{K_0^2 \log^2 K}{K^{\sigma_1}}\right) + \mathcal{O}\left(\frac{K_0 \log K}{K^{\sigma_2}}\right). \tag{79}
\end{aligned}$$

Similar to the derivation of (57), (79) implies

$$\begin{aligned}
& \frac{1}{K - M_K + 1} \sum_{k=M_K}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 \\
& = \mathcal{O}\left(\frac{1}{K^{1 - \sigma_2}}\right) + \mathcal{O}\left(\frac{1}{K^{2(\sigma_1 - \sigma_2)}}\right) + \mathcal{O}\left(\frac{K_0^2}{K^{2\sigma_2}}\right) + \mathcal{O}\left(\frac{K_0^2 \log^2 K}{K^{\sigma_1}}\right) + \mathcal{O}\left(\frac{K_0 \log K}{K^{\sigma_2}}\right).
\end{aligned}$$

Similar to (58), we have

$$\begin{aligned}
\frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2 & = \mathcal{O}\left(\frac{K_0 \log K}{K}\right) + \mathcal{O}\left(\frac{1}{K - M_K + 1} \sum_{k=M_K}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2\right) \\
& = \mathcal{O}\left(\frac{1}{K - M_K + 1} \sum_{k=M_K}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2\right) \tag{80}
\end{aligned}$$

which completes the proof. \square

B.4 PROOF OF THEOREM 4

Given the definition in section B.2, we now give the convergence proof of actor update in Algorithm 1 with linear value function approximation and Markovian sampling method.

By following the derivation of (60), we have

$$\begin{aligned}
\mathbb{E}[J(\theta_{k+1})] & \geq \mathbb{E}[J(\theta_k)] + \underbrace{\alpha_k \mathbb{E} \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_{k-\tau_k}) - \hat{\delta}(x_{(k)}, \omega_k^*) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle}_{I_1} \\
& \quad + \underbrace{\alpha_k \mathbb{E} \left\langle \nabla J(\theta_k), \hat{\delta}(x_{(k)}, \omega_k^*) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle}_{I_2} - \frac{L_J}{2} C_p^2 \alpha_k^2. \tag{81}
\end{aligned}$$

The item I_1 can be bounded by following (61) as

$$I_1 \geq -\sqrt{2}(1 + \gamma) C_\psi \sqrt{\mathbb{E} \|\nabla J(\theta_k)\|_2^2} \sqrt{C_\delta^2 K_0^2 \beta_{k-1}^2 + \mathbb{E} \|\omega_k - \omega_k^*\|_2^2}. \tag{82}$$

Next we consider I_2 . We first decompose it as

$$\begin{aligned}
I_2 & = \mathbb{E} \left\langle \nabla J(\theta_k), \hat{\delta}(x_{(k)}, \omega_k^*) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\
& = \mathbb{E} \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_k^*) - \delta(x_{(k)}, \theta_k) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\
& \quad + \underbrace{\mathbb{E} \left\langle \nabla J(\theta_k), \delta(x_{(k)}, \theta_k) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle}_{I_2^{(1)}} + \mathbb{E} \|\nabla J(\theta_k)\|_2^2. \tag{83} \\
& \quad \underbrace{\hspace{10em}}_{I_2^{(2)}}
\end{aligned}$$

For some $m \in \mathbb{N}^+$, define $M := (K_0 + 1)m + K_0$. Following Lemma 6, for some $d_m \leq M$ and positive constants $D_2, D_3, D_4, D_5, I_2^{(1)}$ can be bounded as

$$\begin{aligned}
I_2^{(1)} &= \mathbb{E} \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_k^*) - \delta(x_{(k)}, \theta_k) \right) \psi_{\theta_k - \tau_k}(s_{(k)}, a_{(k)}) \right\rangle \\
&\geq -D_2 \mathbb{E} \|\theta_{k-\tau_k} - \theta_{k-d_m}\|_2 - D_3 \mathbb{E} \|\theta_k - \theta_{k-d_m}\|_2 - D_4 \sum_{i=k-d_m}^{k-\tau_k} \mathbb{E} \|\theta_i - \theta_{k-d_m}\|_2 \\
&\quad - D_5 \kappa \rho^{m-1} - L_V C_\psi (1 + \gamma) \epsilon_{app} \\
&\geq -D_2 (d_m - \tau_k) C_p \alpha_{k-d_m} - D_3 d_m C_p \alpha_{k-d_m} - D_4 (d_m - \tau_k)^2 C_p \alpha_{k-d_m} \\
&\quad - D_5 \kappa \rho^{m-1} - (1 + \gamma) L_V C_\psi \epsilon_{app}, \tag{84}
\end{aligned}$$

where the derivation of the last inequality is similar to that of (75). By setting $m = m_K$ in (84), and following the fact that $d_{m_K} \leq M_K$ and $\tau_k \geq 0$, we have

$$\begin{aligned}
I_2^{(1)} &\geq -D_2 M_K C_p \alpha_{k-M_K} - D_3 M_K C_p \alpha_{k-M_K} - D_4 M_K^2 C_p \alpha_{k-M_K} - D_5 \kappa \rho^{m_K-1} - (1 + \gamma) L_V C_\psi \epsilon_{app} \\
&= -((D_2 + D_3) C_p M_K + D_4 C_p M_K^2) \alpha_{k-M_K} - D_5 \kappa \rho^{m_K-1} - (1 + \gamma) L_V C_\psi \epsilon_{app} \\
&\geq -((D_2 + D_3) C_p M_K + D_4 C_p M_K^2) \alpha_{k-M_K} - D_5 \alpha_K - (1 + \gamma) L_V C_\psi \epsilon_{app}, \tag{85}
\end{aligned}$$

where the last inequality is due to the definition of m_K .

Following Lemma 7, for some positive constants D_6, D_7 and D_8 , we bound $I_2^{(2)}$ as

$$\begin{aligned}
I_2^{(2)} &= \mathbb{E} \left\langle \nabla J(\theta_k), \delta(x_{(k)}, \theta_k) \psi_{\theta_k - \tau_k}(s_{(k)}, a_{(k)}) - \nabla J(\theta_k) \right\rangle \\
&\geq -D_6 \mathbb{E} \|\theta_{k-\tau_k} - \theta_{k-d_m}\|_2 - D_7 \mathbb{E} \|\theta_k - \theta_{k-d_m}\|_2 - D_8 \sum_{i=\tau_k}^{d_m} \mathbb{E} \|\theta_{k-i} - \theta_{k-d_m}\|_2 - D_9 \kappa \rho^{m-1} - \epsilon_{sp}.
\end{aligned}$$

Similar to the derivation of (85), we have

$$I_2^{(2)} \geq -(D_6 C_p M_K + D_7 C_p M_K + D_8 C_p M_K^2) \alpha_{k-M_K} - D_9 \alpha_K - \epsilon_{sp}. \tag{86}$$

Collecting the lower bounds of $I_2^{(1)}$ and $I_2^{(2)}$ yields

$$I_2 \geq -D_K \alpha_{k-M_K} - (D_5 + D_9) \alpha_K - (1 + \gamma) L_V C_\psi \epsilon_{app} - \epsilon_{sp} + \mathbb{E} \|\nabla J(\theta_k)\|_2^2, \tag{87}$$

where we define $D_K := C_p (D_4 + D_8) M_K^2 + C_p (D_2 + D_3 + D_6 + D_7) M_K$ for brevity.

Substituting (82) and (87) into (81) yields

$$\begin{aligned}
\mathbb{E}[J(\theta_{k+1})] &\geq \mathbb{E}[J(\theta_k)] - \alpha_k \sqrt{2} (1 + \gamma) C_\psi \sqrt{\mathbb{E} \|\nabla J(\theta_k)\|_2^2} \sqrt{C_\delta^2 K_0^2 \beta_{k-1}^2 + \mathbb{E} \|\omega_k - \omega_k^*\|_2^2} - \alpha_k \epsilon_{sp} \\
&\quad - \alpha_k (D_K \alpha_{k-M_K} + (D_5 + D_9) \alpha_K) - \alpha_k (1 + \gamma) L_V C_\psi \epsilon_{app} + \alpha_k \mathbb{E} \|\nabla J(\theta_k)\|_2^2 - \frac{L_J}{2} C_p^2 \alpha_k^2.
\end{aligned}$$

Rearranging and dividing both sides by α_k yield

$$\begin{aligned}
\mathbb{E} \|\nabla J(\theta_k)\|_2^2 &\leq \frac{1}{\alpha_k} (\mathbb{E}[J(\theta_{k+1})] - \mathbb{E}[J(\theta_k)]) + D_K \alpha_{k-M_K} + (D_5 + D_9) \alpha_K + \frac{L_J}{2} C_p^2 \alpha_k + \epsilon_{sp} \\
&\quad + \sqrt{2} (1 + \gamma) C_\psi \sqrt{\mathbb{E} \|\nabla J(\theta_k)\|_2^2} \sqrt{C_\delta^2 K_0^2 \beta_{k-1}^2 + \mathbb{E} \|\omega_k - \omega_k^*\|_2^2} + (1 + \gamma) L_V C_\psi \epsilon_{app}.
\end{aligned}$$

Taking summation gives

$$\begin{aligned}
\sum_{k=M_K}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 &\leq \underbrace{\sum_{k=M_K}^K \frac{1}{\alpha_k} (\mathbb{E}[J(\theta_{k+1})] - \mathbb{E}[J(\theta_k)])}_{I_3} \\
&\quad + \underbrace{\sum_{k=M_K}^K \left(D_K \alpha_{k-M_K} + \frac{L_J}{2} C_p^2 \alpha_k + (D_5 + D_9) \alpha_K \right)}_{I_4} \\
&\quad + \underbrace{\sqrt{2}(1+\gamma)C_\psi \sum_{k=M_K}^K \sqrt{\mathbb{E} \|\nabla J(\theta_k)\|_2^2} \sqrt{C_\delta^2 K_0^2 \beta_{k-1}^2 + \mathbb{E} \|\omega_k - \omega_k^*\|_2^2}}_{I_5} \\
&\quad + (K - M_K + 1)(1+\gamma)L_V C_\psi \epsilon_{app} + (K - M_K + 1)\epsilon_{sp}. \tag{88}
\end{aligned}$$

in which the upper bounds of I_3 and I_5 have already been given by (68) and (70) respectively.

We bound I_4 as

$$\begin{aligned}
I_4 &= \sum_{k=M_K}^K \left(D_K \alpha_{k-M_K} + \frac{L_J}{2} C_p^2 \alpha_k + (D_5 + D_9) \alpha_K \right) \\
&\leq \sum_{k=M_K}^K \left(D_K \alpha_{k-M_K} + \frac{L_J}{2} C_p^2 \alpha_{k-M_K} + (D_5 + D_9) \alpha_K \right) \\
&= \left(D_K + \frac{L_J}{2} C_p^2 \right) \sum_{k=M_K}^K \alpha_{k-M_K} + (D_5 + D_9)(K - M_K + 1) \alpha_K \\
&= \left(D_K + \frac{L_J}{2} C_p^2 \right) \sum_{k=0}^{K-M_K} \alpha_k + (D_5 + D_9)(K - M_K + 1) \alpha_K \\
&\leq \left(D_K + \frac{L_J}{2} C_p^2 \right) \frac{c_1}{1 - \sigma_1} K^{1-\sigma_1} + c_1 (D_5 + D_9)(K + 1)^{1-\sigma_1} \\
&= \mathcal{O}((K_0^2 \log^2 K) K^{1-\sigma_1}) \tag{89}
\end{aligned}$$

where the last inequality uses $\sum_{k=a}^b k^{-\sigma} \leq \frac{b^{1-\sigma}}{1-\sigma}$, and the last equality is due to the fact that

$$\mathcal{O}(D_K) = \mathcal{O}(M_K^2 + M_K) = \mathcal{O}((K_0 m_K)^2 + K_0 m_K) = \mathcal{O}(K_0^2 \log^2 K).$$

Substituting the upper bounds of I_3 , I_4 and I_5 into (88), and dividing both sides by $K - M_K + 1$ give

$$\begin{aligned}
\frac{1}{K - M_K + 1} \sum_{k=M_K}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 &\leq \frac{\sqrt{2}(1+\gamma)C_\psi}{K - M_K + 1} \sqrt{\sum_{k=M_K}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2} \sqrt{\mathcal{O}(K_0^2 K^{1-2\sigma_2}) + \sum_{k=M_K}^K \mathbb{E} \|\omega_k - \omega_k^*\|_2^2} \\
&\quad + \mathcal{O}\left(\frac{1}{K^{1-\sigma_1}}\right) + \mathcal{O}\left(\frac{K_0^2 \log^2 K}{K^{\sigma_1}}\right) + \mathcal{O}(\epsilon_{app} + \epsilon_{sp}). \tag{90}
\end{aligned}$$

Following the similar steps of those in (72), (90) implies

$$\begin{aligned}
\frac{1}{K - M_K + 1} \sum_{k=M_K}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 &= \mathcal{O}\left(\frac{1}{K^{1-\sigma_1}}\right) + \mathcal{O}\left(\frac{K_0^2 \log^2 K}{K^{\sigma_1}}\right) + \mathcal{O}\left(\frac{K_0^2}{K^{2\sigma_2}}\right) \\
&\quad + \mathcal{O}\left(\frac{1}{K - M_K + 1} \sum_{k=M_K}^K \mathbb{E} \|\omega_k - \omega_{\theta_k}^*\|_2^2\right) + \mathcal{O}(\epsilon_{app} + \epsilon_{sp}).
\end{aligned}$$

Similar to (73), we have

$$\begin{aligned} \frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 &= \mathcal{O} \left(\frac{K_0 \log K}{K} \right) + \mathcal{O} \left(\frac{1}{K - M_K + 1} \sum_{k=M_K}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 \right) \\ &= \mathcal{O} \left(\frac{1}{K - M_K + 1} \sum_{k=M_K}^K \mathbb{E} \|\nabla J(\theta_k)\|_2^2 \right) \end{aligned}$$

which completes the proof. \square

C SUPPORTING LEMMAS

C.1 SUPPORTING LEMMAS FOR THEOREM 3

Lemma 5. *For any $m \geq 1$ and $k \geq (K_0 + 1)m + K_0 + 1$, we have*

$$\begin{aligned} \mathbb{E} \langle \omega_k - \omega_{\theta_k}^*, g(x_{(k)}, \omega_k) - \bar{g}(\theta_k, \omega_k) \rangle &\leq C_4 \mathbb{E} \|\theta_k - \theta_{k-d_m}\|_2 + C_5 \sum_{i=\tau_k}^{d_m} \mathbb{E} \|\theta_{k-i} - \theta_{k-d_m}\|_2 \\ &\quad + C_6 \mathbb{E} \|\omega_k - \omega_{k-d_m}\|_2 + C_7 \kappa \rho^{m-1}, \end{aligned}$$

where constant $d_m \leq (K_0 + 1)m + K_0$, and $C_4 := 2C_\delta L_\omega + 4R_\omega C_\delta |\mathcal{A}| L_\pi (1 + \log_\rho \kappa^{-1} + (1 - \rho)^{-1})$, $C_5 := 4R_\omega C_\delta |\mathcal{A}| L_\pi$ and $C_6 := 4(1 + \gamma)R_\omega + 2C_\delta$, $C_7 := 8R_\omega C_\delta$.

Proof. Consider the collection of random samples $\{x_{(k-K_0-1)}, x_{(k-K_0)}, \dots, x_{(k)}\}$. Suppose $x_{(k)}$ is sampled by worker n , then due to Assumption 1, $\{x_{(k-K_0-1)}, x_{(k-K_0)}, \dots, x_{(k-1)}\}$ will contain at least another sample drawn by worker n . Therefore, $\{x_{(k-(K_0+1)m)}, x_{(k-(K_0+1)m+1)}, \dots, x_{(k-1)}\}$ will contain at least m samples from worker n .

Consider the Markov chain formed by $m + 1$ samples in $\{x_{(k-(K_0+1)m)}, x_{(k-(K_0+1)m+1)}, \dots, x_{(k)}\}$:

$$s_{t-m} \xrightarrow{\theta_{k-d_m}} a_{t-m} \xrightarrow{\mathcal{P}} s_{t-m+1} \xrightarrow{\theta_{k-d_{m-1}}} a_{t-m+1} \cdots s_{t-1} \xrightarrow{\theta_{k-d_1}} a_{t-1} \xrightarrow{\mathcal{P}} s_t \xrightarrow{\theta_{k-d_0}} a_t \xrightarrow{\mathcal{P}} s_{t+1},$$

where $(s_t, a_t, s_{t+1}) = (s_{(k)}, a_{(k)}, s'_{(k)})$, and $\{d_j\}_{j=0}^m$ is some increasing sequence with $d_0 := \tau_k$.

Suppose θ_{k-d_m} was used to generate the k_m th update, then we have $x_{t-m} = x_{(k_m)}$. Following Assumption 1, we have $\tau_{k_m} = k_m - (k - d_m) \leq K_0$. Since $x_{(k_m)}$ is in $\{x_{(k-(K_0+1)m)}, \dots, x_{(k)}\}$, we have $k_m \geq k - (K_0 + 1)m$. Combining these two inequalities, we have

$$d_m \leq (K_0 + 1)m + K_0. \quad (91)$$

Given $(s_{t-m}, a_{t-m}, s_{t-m+1})$ and θ_{k-d_m} , we construct an auxiliary Markov chain as Lemma 3:

$$s_{t-m} \xrightarrow{\theta_{k-d_m}} a_{t-m} \xrightarrow{\mathcal{P}} s_{t-m+1} \xrightarrow{\theta_{k-d_m}} \tilde{a}_{t-m+1} \cdots \tilde{s}_{t-1} \xrightarrow{\theta_{k-d_m}} \tilde{a}_{t-1} \xrightarrow{\mathcal{P}} \tilde{s}_t \xrightarrow{\theta_{k-d_m}} \tilde{a}_t \xrightarrow{\mathcal{P}} \tilde{s}_{t+1}.$$

For brevity, we define

$$\Delta_1(x, \theta, \omega) := \langle \omega - \omega_\theta^*, g(x, \omega) - \bar{g}(\theta, \omega) \rangle.$$

Throughout this proof, we use $\theta, \theta', \omega, \omega', x$ and \tilde{x} as shorthand notations of $\theta_k, \theta_{k-d_m}, \omega_k, \omega_{k-d_m}, x_t$ and $\tilde{x}_t := (\tilde{s}_t, \tilde{a}_t, \tilde{s}_{t+1})$, respectively.

First we decompose $\Delta_1(x, \theta, \omega)$ as

$$\begin{aligned} \Delta_1(x, \theta, \omega) &= \underbrace{\Delta_1(x, \theta, \omega) - \Delta_1(x, \theta', \omega)}_{I_1} + \underbrace{\Delta_1(x, \theta', \omega) - \Delta_1(x, \theta', \omega')}_{I_2} \\ &\quad + \underbrace{\Delta_1(x, \theta', \omega') - \Delta_1(\tilde{x}, \theta', \omega')}_{I_3} + \underbrace{\Delta_1(\tilde{x}, \theta', \omega')}_{I_4}. \end{aligned} \quad (92)$$

We bound I_1 in (92) as

$$\begin{aligned}\Delta_1(x, \theta, \omega) - \Delta_1(x, \theta', \omega) &= \langle \omega - \omega_\theta^*, g(x, \omega) - \bar{g}(\theta, \omega) \rangle - \langle \omega - \omega_{\theta'}^*, g(x, \omega) - \bar{g}(\theta', \omega) \rangle \\ &\leq |\langle \omega - \omega_\theta^*, g(x, \omega) - \bar{g}(\theta, \omega) \rangle - \langle \omega - \omega_{\theta'}^*, g(x, \omega) - \bar{g}(\theta, \omega) \rangle| \\ &\quad + |\langle \omega - \omega_{\theta'}^*, g(x, \omega) - \bar{g}(\theta, \omega) \rangle - \langle \omega - \omega_{\theta'}^*, g(x, \omega) - \bar{g}(\theta', \omega) \rangle|. \end{aligned} \quad (93)$$

For the first term in (93), we have

$$\begin{aligned}|\langle \omega - \omega_\theta^*, g(x, \omega) - \bar{g}(\theta, \omega) \rangle - \langle \omega - \omega_{\theta'}^*, g(x, \omega) - \bar{g}(\theta, \omega) \rangle| &= |\langle \omega_\theta^* - \omega_{\theta'}^*, g(x, \omega) - \bar{g}(\theta, \omega) \rangle| \\ &\leq \|\omega_\theta^* - \omega_{\theta'}^*\|_2 \|g(x, \omega) - \bar{g}(\theta, \omega)\|_2 \\ &\leq 2C_\delta \|\omega_\theta^* - \omega_{\theta'}^*\|_2 \\ &\leq 2C_\delta L_\omega \|\theta - \theta'\|_2, \end{aligned}$$

where the last inequality is due to Proposition 2.

We use $x \sim \theta'$ as shorthand notations to represent that $s \sim \mu_{\theta'}$, $a \sim \pi_{\theta'}$, $s' \sim \mathcal{P}$. For the second term in (93), we have

$$\begin{aligned}&|\langle \omega - \omega_{\theta'}^*, g(x, \omega) - \bar{g}(\theta, \omega) \rangle - \langle \omega - \omega_{\theta'}^*, g(x, \omega) - \bar{g}(\theta', \omega) \rangle| \\ &= |\langle \omega - \omega_{\theta'}^*, \bar{g}(\theta', \omega) - \bar{g}(\theta, \omega) \rangle| \\ &\leq \|\omega - \omega_{\theta'}^*\|_2 \|\bar{g}(\theta', \omega) - \bar{g}(\theta, \omega)\|_2 \\ &\leq 2R_\omega \|\bar{g}(\theta', \omega) - \bar{g}(\theta, \omega)\|_2 \\ &= 2R_\omega \left\| \mathbb{E}_{x \sim \theta'} [g(x, \omega)] - \mathbb{E}_{x \sim \theta} [g(x, \omega)] \right\|_2 \\ &\leq 2R_\omega \sup_x \|g(x, \omega)\|_2 \|\mu_{\theta'} \otimes \pi_{\theta'} \otimes \mathcal{P} - \mu_\theta \otimes \pi_\theta \otimes \mathcal{P}\|_{TV} \\ &\leq 2R_\omega C_\delta \|\mu_{\theta'} \otimes \pi_{\theta'} \otimes \mathcal{P} - \mu_\theta \otimes \pi_\theta \otimes \mathcal{P}\|_{TV} \\ &= 4R_\omega C_\delta d_{TV}(\mu_{\theta'} \otimes \pi_{\theta'} \otimes \mathcal{P}, \mu_\theta \otimes \pi_\theta \otimes \mathcal{P}) \\ &\leq 4R_\omega C_\delta |\mathcal{A}| L_\pi (1 + \log_\rho \kappa^{-1} + (1 - \rho)^{-1}) \|\theta - \theta'\|_2, \end{aligned}$$

where the third inequality follows the definition of TV norm, the second last inequality follows (37), and the last inequality follows [19, Lemma A.1].

Collecting the upper bounds of the two terms in (93) yields

$$I_1 \leq [2C_\delta L_\omega + 4R_\omega C_\delta |\mathcal{A}| L_\pi (1 + \log_\rho \kappa^{-1} + (1 - \rho)^{-1})] \|\theta - \theta'\|_2.$$

Next we bound $\mathbb{E}[I_2]$ in (92) as

$$\begin{aligned}\mathbb{E}[I_2] &= \mathbb{E}[\Delta_1(x, \theta', \omega) - \Delta_1(x, \theta', \omega')] \\ &= \mathbb{E} \langle \omega - \omega_{\theta'}^*, g(x, \omega) - \bar{g}(\theta', \omega) \rangle - \langle \omega' - \omega_{\theta'}^*, g(x, \omega') - \bar{g}(\theta', \omega') \rangle \\ &\leq \mathbb{E} |\langle \omega - \omega_{\theta'}^*, g(x, \omega) - \bar{g}(\theta', \omega) \rangle - \langle \omega - \omega_{\theta'}^*, g(x, \omega') - \bar{g}(\theta', \omega') \rangle| \\ &\quad + \mathbb{E} |\langle \omega - \omega_{\theta'}^*, g(x, \omega') - \bar{g}(\theta', \omega') \rangle - \langle \omega' - \omega_{\theta'}^*, g(x, \omega') - \bar{g}(\theta', \omega') \rangle|. \end{aligned} \quad (94)$$

We bound the first term in (94) as

$$\begin{aligned}&\mathbb{E} |\langle \omega - \omega_{\theta'}^*, g(x, \omega) - \bar{g}(\theta', \omega) \rangle - \langle \omega - \omega_{\theta'}^*, g(x, \omega') - \bar{g}(\theta', \omega') \rangle| \\ &= \mathbb{E} |\langle \omega - \omega_{\theta'}^*, g(x, \omega) - g(x, \omega') + \bar{g}(\theta', \omega') - \bar{g}(\theta', \omega) \rangle| \\ &\leq 2R_\omega (\mathbb{E} \|g(x, \omega) - g(x, \omega')\|_2 + \mathbb{E} \|\bar{g}(\theta', \omega') - \bar{g}(\theta', \omega)\|_2) \\ &\leq 2R_\omega \left(\mathbb{E} \|g(x, \omega) - g(x, \omega')\|_2 + \mathbb{E} \left\| \mathbb{E}_{x \sim \theta'} [g(x, \omega')] - \mathbb{E}_{x \sim \theta'} [g(x, \omega)] \right\|_2 \right) \\ &= 2R_\omega \left(\mathbb{E} \|(\gamma\phi(s') - \phi(s))^\top (\omega - \omega')\|_2 + \mathbb{E} \left\| \mathbb{E}_{x \sim \theta'} [(\gamma\phi(s') - \phi(s))^\top] (\omega' - \omega) \right\|_2 \right) \\ &\leq 2R_\omega ((1 + \gamma) \mathbb{E} \|\omega - \omega'\|_2 + (1 + \gamma) \mathbb{E} \|\omega - \omega'\|_2) \\ &= 4R_\omega (1 + \gamma) \mathbb{E} \|\omega - \omega'\|_2. \end{aligned}$$

We bound the second term in (94) as

$$\begin{aligned} & \mathbb{E} |\langle \omega - \omega_{\theta'}^*, g(x, \omega') - \bar{g}(\theta', \omega') \rangle - \langle \omega' - \omega_{\theta'}^*, g(x, \omega') - \bar{g}(\theta', \omega') \rangle| \\ &= \mathbb{E} |\langle \omega - \omega', g(x, \omega') - \bar{g}(\theta', \omega') \rangle| \\ &\leq 2C_\delta \mathbb{E} \|\omega - \omega'\|_2. \end{aligned}$$

Collecting the upper bounds of the two terms in (94) yields

$$\mathbb{E}[I_2] \leq (4(1 + \gamma)R_\omega + 2C_\delta) \mathbb{E} \|\omega_k - \omega_{k-d_m}\|_2. \quad (95)$$

We first bound I_3 as

$$\begin{aligned} \mathbb{E}[I_3 | \theta', \omega', s_{t-m+1}] &= \mathbb{E} [\Delta_1(x, \theta', \omega') - \Delta_1(\tilde{x}, \theta', \omega') | \theta', \omega', s_{t-m+1}] \\ &\leq |\mathbb{E} [\Delta_1(x, \theta', \omega') | \theta', \omega', s_{t-m+1}] - \mathbb{E} [\Delta_1(\tilde{x}, \theta', \omega') | \theta', \omega', s_{t-m+1}]| \\ &\leq \sup_x |\Delta_1(x, \theta', \omega')| \|\mathbb{P}(x \in \cdot | \theta', \omega', s_{t-m+1}) - \mathbb{P}(\tilde{x} \in \cdot | \theta', \omega', s_{t-m+1})\|_{TV} \\ &\leq 8R_\omega C_\delta d_{TV} (\mathbb{P}(x \in \cdot | \theta', s_{t-m+1}), \mathbb{P}(\tilde{x} \in \cdot | \theta', s_{t-m+1})), \end{aligned} \quad (96)$$

where the second last inequality follows the definition of TV norm, and the last inequality follows

$$|\Delta_1(x, \theta', \omega')| \leq \|\omega' - \omega_{\theta'}^*\|_2 \|g(x, \omega') - \bar{g}(\theta', \omega')\|_2 \leq 4R_\omega C_\delta.$$

By following (27) in Lemma 3, we have

$$d_{TV} (\mathbb{P}(x \in \cdot | \theta', s_{t-m+1}), \mathbb{P}(\tilde{x} \in \cdot | \theta', s_{t-m+1})) \leq \frac{1}{2} |\mathcal{A}| L_\pi \sum_{i=\tau_k}^{d_m} \mathbb{E} [\|\theta_{k-i} - \theta_{k-d_m}\|_2 | \theta', s_{t-m+1}].$$

Substituting the last inequality into (96), then taking total expectation on both sides yield

$$\mathbb{E}[I_3] \leq 4R_\omega C_\delta |\mathcal{A}| L_\pi \sum_{i=\tau_k}^{d_m} \mathbb{E} \|\theta_{k-i} - \theta_{k-d_m}\|_2. \quad (97)$$

Next we bound I_4 . Define $\bar{x} := (\bar{s}, \bar{a}, \bar{s}')$ where $\bar{s} \sim \mu_{\theta'}$, $\bar{a} \sim \pi_{\theta'}$ and $\bar{s}' \sim \mathcal{P}$. It is immediate that

$$\begin{aligned} \mathbb{E} [\Delta_1(\bar{x}, \theta', \omega') | \theta', \omega', s_{t-m+1}] &= \langle \omega' - \omega_{\theta'}^*, \mathbb{E} [g(\bar{x}, \omega') | \theta', \omega', s_{t-m+1}] - \bar{g}(\theta', \omega') \rangle \\ &= \langle \omega' - \omega_{\theta'}^*, \bar{g}(\theta', \omega') - \bar{g}(\theta', \omega') \rangle = 0. \end{aligned} \quad (98)$$

Then we have

$$\begin{aligned} \mathbb{E}[I_4 | \theta', \omega', s_{t-m+1}] &= \mathbb{E} [\Delta_1(\tilde{x}, \theta', \omega') - \Delta_1(\bar{x}, \theta', \omega') | \theta', \omega', s_{t-m+1}] \\ &\leq |\mathbb{E} [\Delta_1(\tilde{x}, \theta', \omega') | \theta', \omega', s_{t-m+1}] - \mathbb{E} [\Delta_1(\bar{x}, \theta', \omega') | \theta', \omega', s_{t-m+1}]| \\ &\leq \sup_x |\Delta_1(x, \theta', \omega')| \|\mathbb{P}(\tilde{x} \in \cdot | \theta', s_{t-m+1}) - \mathbb{P}(\bar{x} \in \cdot | \theta', s_{t-m+1})\|_{TV} \\ &\leq 8R_\omega C_\delta d_{TV} (\mathbb{P}(\tilde{x} \in \cdot | \theta', s_{t-m+1}), \mathbb{P}(\bar{x} \in \cdot | \theta', s_{t-m+1})) \\ &= 8R_\omega C_\delta d_{TV} (\mathbb{P}(\tilde{x} \in \cdot | \theta', s_{t-m+1}), \mu_{\theta'} \otimes \pi_{\theta'} \otimes \mathcal{P}), \end{aligned} \quad (99)$$

where the second inequality follows the definition of TV norm, and the third inequality follows (98).

The auxiliary Markov chain with policy $\pi_{\theta'}$ starts from initial state s_{t-m+1} , and \tilde{s}_t is the $(m-1)$ th state on the chain. Following Lemma 1, we have:

$$\begin{aligned} & d_{TV} (\mathbb{P}(\tilde{x} \in \cdot | \theta', s_{t-m+1}), \mu_{\theta'} \otimes \pi_{\theta'} \otimes \mathcal{P}) \\ &= d_{TV} (\mathbb{P}((\tilde{s}_t, \tilde{a}_t, \tilde{s}_{t+1}) \in \cdot | \theta', s_{t-m+1}), \mu_{\theta'} \otimes \pi_{\theta'} \otimes \mathcal{P}) \leq \kappa \rho^{m-1}. \end{aligned}$$

Substituting the last inequality into (99) and taking total expectation on both sides yield

$$\mathbb{E}[I_4] \leq 8R_\omega C_\delta \kappa \rho^{m-1}.$$

Taking total expectation on (92) and collecting bounds of I_1, I_2, I_3, I_4 yield

$$\begin{aligned} \mathbb{E} [\Delta_1(x, \theta, \omega)] &\leq C_4 \mathbb{E} \|\theta_k - \theta_{k-d_m}\|_2 + C_5 \sum_{i=\tau_k}^{d_m} \mathbb{E} \|\theta_{k-i} - \theta_{k-d_m}\|_2 \\ &\quad + C_6 \mathbb{E} \|\omega_k - \omega_{k-d_m}\|_2 + C_7 \kappa \rho^{m-1}, \end{aligned}$$

where $C_4 := 2C_\delta L_\omega + 4R_\omega C_\delta |\mathcal{A}| L_\pi (1 + \log \kappa^{-1} + (1 - \rho)^{-1})$, $C_5 := 4R_\omega C_\delta |\mathcal{A}| L_\pi$, $C_6 := 4(1 + \gamma)R_\omega + 2C_\delta$ and $C_7 := 8R_\omega C_\delta$. \square

C.2 SUPPORTING LEMMAS FOR THEOREM 4

Lemma 6. For any $m \geq 1$ and $k \geq (K_0 + 1)m + K_0 + 1$, we have

$$\begin{aligned} \mathbb{E} \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_k^*) - \delta(x_{(k)}, \theta_k) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle &\geq -D_2 \mathbb{E} \|\theta_{k-\tau_k} - \theta_{k-d_m}\|_2 \\ &\quad - D_3 \mathbb{E} \|\theta_k - \theta_{k-d_m}\|_2 - D_4 \sum_{i=\tau_k}^{d_m} \mathbb{E} \|\theta_{k-i} - \theta_{k-d_m}\|_2 - D_5 \kappa \rho^{m-1} - L_V C_\psi (1 + \gamma) \epsilon_{app}, \end{aligned}$$

where $D_2 := 2L_V L_\psi C_\delta$, $D_3 := (2C_\delta C_\psi L_J + L_V C_\psi (L_\omega + L_V)(1 + \gamma))$, $D_4 := 2L_V C_\psi C_\delta |\mathcal{A}| L_\pi$ and $D_5 := 4L_V C_\psi C_\delta$.

Proof. For the worker that contributes to the k th update, we construct its Markov chain:

$$s_{t-m} \xrightarrow{\theta_{k-d_m}} a_{t-m} \xrightarrow{\mathcal{P}} s_{t-m+1} \xrightarrow{\theta_{k-d_{m-1}}} a_{t-m+1} \cdots s_{t-1} \xrightarrow{\theta_{k-d_1}} a_{t-1} \xrightarrow{\mathcal{P}} s_t \xrightarrow{\theta_{k-d_0}} a_t \xrightarrow{\mathcal{P}} s_{t+1},$$

where $(s_t, a_t, s_{t+1}) = (s_{(k)}, a_{(k)}, s'_{(k)})$, and $\{d_j\}_{j=0}^m$ is some increasing sequence with $d_0 := \tau_k$. By (91) in Lemma 5, we have $d_m \leq (K_0 + 1)m + K_0$.

Given $(s_{t-m}, a_{t-m}, s_{t-m+1})$ and θ_{k-d_m} , we construct an auxiliary Markov chain:

$$s_{t-m} \xrightarrow{\theta_{k-d_m}} a_{t-m} \xrightarrow{\mathcal{P}} s_{t-m+1} \xrightarrow{\theta_{k-d_m}} \tilde{a}_{t-m+1} \cdots \tilde{s}_{t-1} \xrightarrow{\theta_{k-d_m}} \tilde{a}_{t-1} \xrightarrow{\mathcal{P}} \tilde{s}_t \xrightarrow{\theta_{k-d_m}} \tilde{a}_t \xrightarrow{\mathcal{P}} \tilde{s}_{t+1}.$$

First we have

$$\begin{aligned} &\left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_k^*) - \delta(x_{(k)}, \theta_k) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\ &= \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_k^*) - \delta(x_{(k)}, \theta_k) \right) \left(\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)}) \right) \right\rangle \\ &\quad + \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_k^*) - \delta(x_{(k)}, \theta_k) \right) \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)}) \right\rangle. \end{aligned} \quad (100)$$

We first bound the first term in (100) as

$$\begin{aligned} &\left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_k^*) - \delta(x_{(k)}, \theta_k) \right) \left(\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)}) \right) \right\rangle \\ &\geq -\|J(\theta_k)\|_2 \|\hat{\delta}(x_{(k)}, \omega_k^*) - \delta(x_{(k)}, \theta_k)\| \|\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)})\|_2 \\ &\geq -\|J(\theta_k)\|_2 \left(\|\hat{\delta}(x_{(k)}, \omega_k^*)\| + \|\delta(x_{(k)}, \theta_k)\| \right) \|\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)})\|_2 \\ &\geq -L_V \left(\|\hat{\delta}(x_{(k)}, \omega_k^*)\| + \|\delta(x_{(k)}, \theta_k)\| \right) \|\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)})\|_2 \\ &\geq -2L_V C_\delta \|\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)})\|_2 \\ &\geq -2L_V L_\psi C_\delta \|\theta_{k-\tau_k} - \theta_{k-d_m}\|_2, \end{aligned} \quad (101)$$

where the last inequality follows Assumption 3 and second last inequality follows

$$\begin{aligned} |\hat{\delta}(x, \omega_\theta^*)| &\leq |r(x)| + \gamma \|\phi(s')\|_2 \|\omega_\theta^*\|_2 + \|\phi(s)\|_2 \|\omega_\theta^*\|_2 \leq r_{\max} + (1 + \gamma) R_\omega \leq C_\delta, \\ |\delta(x, \theta)| &\leq |r(x)| + \gamma |V_{\pi_\theta}(s')| + |V_{\pi_\theta}(s)| \leq r_{\max} + (1 + \gamma) R_{\max} \leq C_\delta. \end{aligned}$$

Substituting (101) into (100) gives

$$\begin{aligned} &\left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_k^*) - \delta(x_{(k)}, \theta_k) \right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\ &\geq -2L_V L_\psi C_\delta \|\theta_{k-\tau_k} - \theta_{k-d_m}\|_2 + \left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_k^*) - \delta(x_{(k)}, \theta_k) \right) \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)}) \right\rangle. \end{aligned} \quad (102)$$

Then we start to bound the second term in (102). For brevity, we define

$$\Delta_2(x, \theta) := \left\langle \nabla J(\theta), \left(\hat{\delta}(x, \omega_\theta^*) - \delta(x, \theta) \right) \psi_{\theta_{k-d_m}}(s, a) \right\rangle.$$

In the following proof, we use $\theta, \theta', \omega_\theta^*, \omega_{\theta'}^*, x$ and \tilde{x} as shorthand notations for $\theta_k, \theta_{k-d_m}, \omega_k^*, \omega_{k-d_m}^*, x_t$ and \tilde{x}_t respectively. We also define $\bar{x} := (\bar{s}, \bar{a}, \bar{s}')$, where $\bar{s} \sim \mu_{\theta'}$, $\bar{a} \sim \pi_{\theta'}$ and $\bar{s}' \sim \mathcal{P}$.

We decompose the second term in (102) as

$$\Delta_2(x, \theta) = \underbrace{\Delta_2(x, \theta) - \Delta_2(x, \theta')}_{I_1} + \underbrace{\Delta_2(x, \theta') - \Delta_2(\tilde{x}, \theta')}_{I_2} + \underbrace{\Delta_2(\tilde{x}, \theta') - \Delta_2(\bar{x}, \theta')}_{I_3} + \underbrace{\Delta_2(\bar{x}, \theta')}_{I_4}.$$

We bound the term I_1 as

$$\begin{aligned} I_1 &= \left\langle \nabla J(\theta), \left(\hat{\delta}(x, \omega_\theta^*) - \delta(x, \theta) \right) \psi_{\theta'}(s, a) \right\rangle - \left\langle \nabla J(\theta'), \left(\hat{\delta}(x, \omega_{\theta'}^*) - \delta(x, \theta') \right) \psi_{\theta'}(s, a) \right\rangle \\ &= \left\langle \nabla J(\theta), \left(\hat{\delta}(x, \omega_\theta^*) - \delta(x, \theta) \right) \psi_{\theta'}(s, a) \right\rangle - \left\langle \nabla J(\theta'), \left(\hat{\delta}(x, \omega_\theta^*) - \delta(x, \theta) \right) \psi_{\theta'}(s, a) \right\rangle \\ &\quad + \left\langle \nabla J(\theta'), \left(\hat{\delta}(x, \omega_\theta^*) - \delta(x, \theta) \right) \psi_{\theta'}(s, a) \right\rangle - \left\langle \nabla J(\theta'), \left(\hat{\delta}(x, \omega_{\theta'}^*) - \delta(x, \theta') \right) \psi_{\theta'}(s, a) \right\rangle. \end{aligned}$$

For the first term in I_1 , we have

$$\begin{aligned} &\left\langle \nabla J(\theta), \left(\hat{\delta}(x, \omega_\theta^*) - \delta(x, \theta) \right) \psi_{\theta'}(s, a) \right\rangle - \left\langle \nabla J(\theta'), \left(\hat{\delta}(x, \omega_\theta^*) - \delta(x, \theta) \right) \psi_{\theta'}(s, a) \right\rangle \\ &= \left\langle \nabla J(\theta) - \nabla J(\theta'), \left(\hat{\delta}(x, \omega_\theta^*) - \delta(x, \theta) \right) \psi_{\theta'}(s, a) \right\rangle \\ &\geq -\|\nabla J(\theta) - \nabla J(\theta')\|_2 \|\hat{\delta}(x, \omega_\theta^*) - \delta(x, \theta)\|_2 \|\psi_{\theta'}(s, a)\|_2 \\ &\geq -2C_\delta C_\psi \|\nabla J(\theta) - \nabla J(\theta')\|_2 \\ &\geq -2C_\delta C_\psi L_J \|\theta - \theta'\|_2, \end{aligned}$$

where the last inequality is due to the L_J -Lipschitz of policy gradient shown in Proposition 1.

For the second term in I_1 , we have

$$\begin{aligned} &\left\langle \nabla J(\theta'), \left(\hat{\delta}(x, \omega_\theta^*) - \delta(x, \theta) \right) \psi_{\theta'}(s, a) \right\rangle - \left\langle \nabla J(\theta'), \left(\hat{\delta}(x, \omega_{\theta'}^*) - \delta(x, \theta') \right) \psi_{\theta'}(s, a) \right\rangle \\ &= \left\langle \nabla J(\theta'), \left(\hat{\delta}(x, \omega_\theta^*) - \hat{\delta}(x, \omega_{\theta'}^*) + \delta(x, \theta') - \delta(x, \theta) \right) \psi_{\theta'}(s, a) \right\rangle \\ &\geq -L_V C_\psi \left| \hat{\delta}(x, \omega_\theta^*) - \hat{\delta}(x, \omega_{\theta'}^*) + \delta(x, \theta') - \delta(x, \theta) \right| \\ &\geq -L_V C_\psi \left| \gamma \phi(s')^\top (\omega_\theta^* - \omega_{\theta'}^*) + \phi(s)^\top (\omega_{\theta'}^* - \omega_\theta^*) + \gamma V_{\pi_{\theta'}}(s') - \gamma V_{\pi_\theta}(s') + V_{\pi_\theta}(s) - V_{\pi_{\theta'}}(s) \right| \\ &\geq -L_V C_\psi (\gamma \|\omega_\theta^* - \omega_{\theta'}^*\|_2 + \|\omega_{\theta'}^* - \omega_\theta^*\|_2 + \gamma |V_{\pi_{\theta'}}(s') - V_{\pi_\theta}(s')| + |V_{\pi_\theta}(s) - V_{\pi_{\theta'}}(s)|) \\ &\geq -L_V C_\psi (\gamma L_\omega \|\theta - \theta'\|_2 + L_\omega \|\theta - \theta'\|_2 + \gamma L_V \|\theta - \theta'\|_2 + L_V \|\theta - \theta'\|_2) \\ &= -L_V C_\psi (L_\omega + L_V)(1 + \gamma) \|\theta - \theta'\|_2, \end{aligned}$$

where the last inequality is due to the L_ω -Lipschitz continuity of ω_θ^* shown in Proposition 2 and L_V -Lipschitz continuity of $V_{\pi_\theta}(s)$ shown in Lemma 4. Collecting the upper bounds of I_1 yields

$$I_1 \geq -(2C_\delta C_\psi L_J + L_V C_\psi (L_\omega + L_V)(1 + \gamma)) \|\theta - \theta'\|_2.$$

First we bound I_2 as

$$\begin{aligned} \mathbb{E}[I_2 | \theta', s_{t-m+1}] &= \mathbb{E}[\Delta_2(x, \theta') - \Delta_2(\tilde{x}, \theta') | \theta', s_{t-m+1}] \\ &\geq -|\mathbb{E}[\Delta_2(x, \theta') | \theta', s_{t-m+1}] - \mathbb{E}[\Delta_2(\tilde{x}, \theta') | \theta', s_{t-m+1}]| \\ &\geq -\sup_x |\Delta_2(x, \theta')| \|\mathbb{P}(x \in \cdot | \theta', s_{t-m+1}) - \mathbb{P}(\tilde{x} \in \cdot | \theta', s_{t-m+1})\|_{TV} \\ &\geq -4L_V C_\psi C_\delta d_{TV} (\mathbb{P}(x \in \cdot | \theta', s_{t-m+1}), \mathbb{P}(\tilde{x} \in \cdot | \theta', s_{t-m+1})) \\ &\geq -2L_V C_\psi C_\delta |\mathcal{A}| L_\pi \sum_{i=\tau_k}^{d_m} \mathbb{E}[\|\theta_{k-i} - \theta_{k-d_m}\|_2 | \theta', s_{t-m+1}], \end{aligned} \tag{103}$$

where the second inequality is due to the definition of TV norm, the last inequality follows (27) in Lemma 3, and the second last inequality follows the fact that

$$|\Delta_2(x, \theta')| \leq \|\nabla J(\theta')\|_2 \|\hat{\delta}(x, \omega_{\theta'}^*) - \delta(x, \theta')\| \|\psi_{\theta'}(s, a)\|_2 \leq 2L_V C_\delta C_\psi. \tag{104}$$

Taking total expectation on both sides of (103) yields

$$\mathbb{E}[I_2] \geq -2L_V C_\psi C_\delta |\mathcal{A}| L_\pi \sum_{i=\tau_k}^{d_m} \mathbb{E}[\|\theta_{k-i} - \theta_{k-d_m}\|_2].$$

Next we bound I_3 as

$$\begin{aligned}
\mathbb{E}[I_3|\theta', s_{t-m+1}] &= \mathbb{E}[\Delta_2(\tilde{x}, \theta') - \Delta_2(\bar{x}, \theta')|\theta', s_{t-m+1}] \\
&\geq -|\mathbb{E}[\Delta_2(\tilde{x}, \theta')|\theta', s_{t-m+1}] - \mathbb{E}[\Delta_2(\bar{x}, \theta')|\theta', s_{t-m+1}]| \\
&\geq -\sup_x |\Delta_2(x, \theta')| \|\mathbb{P}(\tilde{x} \in \cdot|\theta', s_{t-m+1}) - \mathbb{P}(\bar{x} \in \cdot|\theta', s_{t-m+1})\|_{TV} \\
&\geq -4L_V C_\psi C_\delta d_{TV}(\mathbb{P}(\tilde{x} \in \cdot|\theta', s_{t-m+1}), \mu_{\theta'} \otimes \pi_{\theta'} \otimes \mathcal{P}), \tag{105}
\end{aligned}$$

where the second inequality is due to the definition of TV norm, and the last inequality follows (104).

The auxiliary Markov chain with policy $\pi_{\theta'}$ starts from initial state s_{t-m+1} , and \tilde{s}_t is the $(m-1)$ th state on the chain. Following Lemma 1, we have:

$$d_{TV}(\mathbb{P}(\tilde{x} \in \cdot|\theta', s_{t-m+1}), \mu_{\theta'} \otimes \pi_{\theta'} \otimes \mathcal{P}) = d_{TV}(\mathbb{P}((\tilde{s}_t, \tilde{a}_t, \tilde{s}_{t+1}) \in \cdot|\theta', s_{t-m+1}), \mu_{\theta'} \otimes \pi_{\theta'} \otimes \mathcal{P}) \leq \kappa \rho^{m-1}.$$

Substituting the last inequality into (105) and taking total expectation on both sides yield

$$\mathbb{E}[I_3] \geq -4L_V C_\psi C_\delta \kappa \rho^{m-1}$$

We bound I_4 as

$$\begin{aligned}
\mathbb{E}[I_4|\theta'] &= \mathbb{E}\left[\left\langle \nabla J(\theta'), \left(\hat{\delta}(\bar{x}, \omega_{\theta'}^*) - \delta(\bar{x}, \theta')\right) \psi_{\theta'}(s, a) \right\rangle \middle| \theta' \right] \\
&\geq -L_V C_\psi \mathbb{E}\left[\left|\hat{\delta}(\bar{x}, \omega_{\theta'}^*) - \delta(\bar{x}, \theta')\right| \middle| \theta' \right] \\
&= -L_V C_\psi \mathbb{E}\left[\left|\gamma(\phi(\bar{s}')^\top \omega_{\theta'}^* - V_{\pi_{\theta'}}(\bar{s}')) + V_{\pi_{\theta'}}(\bar{s}) - \phi(\bar{s})^\top \omega_{\theta'}^*\right| \middle| \theta' \right] \\
&\geq -L_V C_\psi (\gamma \mathbb{E}\left[|\phi(\bar{s}')^\top \omega_{\theta'}^* - V_{\pi_{\theta'}}(\bar{s}')| \middle| \theta' \right] + \mathbb{E}\left[|V_{\pi_{\theta'}}(\bar{s}) - \phi(\bar{s})^\top \omega_{\theta'}^*| \middle| \theta' \right]) \\
&\geq -L_V C_\psi \left(\gamma \sqrt{\mathbb{E}\left[|\phi(\bar{s}')^\top \omega_{\theta'}^* - V_{\pi_{\theta'}}(\bar{s}')|^2 \middle| \theta' \right]} + \sqrt{\mathbb{E}\left[|V_{\pi_{\theta'}}(\bar{s}) - \phi(\bar{s})^\top \omega_{\theta'}^*|^2 \middle| \theta' \right]} \right) \\
&= -L_V C_\psi \left(\gamma \sqrt{\mathbb{E}_{\bar{s}' \sim \mu_{\theta'}}\left[|\phi(\bar{s}')^\top \omega_{\theta'}^* - V_{\pi_{\theta'}}(\bar{s}')|^2\right]} + \sqrt{\mathbb{E}_{\bar{s} \sim \mu_{\theta'}}\left[|V_{\pi_{\theta'}}(\bar{s}) - \phi(\bar{s})^\top \omega_{\theta'}^*|^2\right]} \right) \\
&\geq -L_V C_\psi (1 + \gamma) \epsilon_{app}
\end{aligned}$$

where the second last inequality follows Jensen's inequality.

Taking total expectation on both sides of (102), and collecting lower bounds of I_1, I_2, I_3 and I_4 yield

$$\begin{aligned}
&\mathbb{E}\left\langle \nabla J(\theta_k), \left(\hat{\delta}(x_{(k)}, \omega_k^*) - \delta(x_{(k)}, \theta_k)\right) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) \right\rangle \\
&\geq -D_2 \mathbb{E}\|\theta_{k-\tau_k} - \theta_{k-d_m}\|_2 - D_3 \mathbb{E}\|\theta_k - \theta_{k-d_m}\|_2 - D_4 \sum_{i=\tau_k}^{d_m} \mathbb{E}\|\theta_{k-i} - \theta_{k-d_m}\|_2 \\
&\quad - D_5 \kappa \rho^{m-1} - L_V C_\psi (1 + \gamma) \epsilon_{app},
\end{aligned}$$

where $D_2 := 2L_V L_\psi C_\delta$, $D_3 := (2C_\delta C_\psi L_J + L_V C_\psi (L_\omega + L_V)(1 + \gamma))$, $D_4 := 2L_V C_\psi C_\delta |\mathcal{A}| L_\pi$ and $D_5 := 4L_V C_\psi C_\delta$. \square

Lemma 7. For any $m \geq 1$ and $k \geq (K_0 + 1)m + K_0 + 1$, we have

$$\begin{aligned}
&\mathbb{E}\left\langle \nabla J(\theta_k), \delta(x_{(k)}, \theta_k) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \nabla J(\theta_k) \right\rangle \\
&\geq -D_6 \mathbb{E}\|\theta_{k-\tau_k} - \theta_{k-d_m}\|_2 - D_7 \mathbb{E}\|\theta_k - \theta_{k-d_m}\|_2 - D_8 \sum_{i=\tau_k}^{d_m} \mathbb{E}\|\theta_{k-i} - \theta_{k-d_m}\|_2 - D_9 \kappa \rho^{m-1} - \epsilon_{sp},
\end{aligned}$$

where $d_m \leq (K_0 + 1)m + K_0$, $D_6 := L_V C_\delta L_\psi$, $D_7 := C_p L_J + (1 + \gamma) L_V^2 C_\psi + 2L_V L_J$, $D_8 := L_V C_p |\mathcal{A}| L_\pi$ and $D_9 := 2L_V C_p$.

Proof. For the worker that contributes to the k th update, we construct its Markov chain:

$$s_{t-m} \xrightarrow{\theta_{k-d_m}} a_{t-m} \xrightarrow{\mathcal{P}} s_{t-m+1} \xrightarrow{\theta_{k-d_{m-1}}} a_{t-m+1} \cdots s_{t-1} \xrightarrow{\theta_{k-d_1}} a_{t-1} \xrightarrow{\mathcal{P}} s_t \xrightarrow{\theta_{k-d_0}} a_t \xrightarrow{\mathcal{P}} s_{t+1},$$

where $(s_t, a_t, s_{t+1}) = (s_{(k)}, a_{(k)}, s'_{(k)})$, and $\{d_j\}_{j=0}^m$ is some increasing sequence with $d_0 := \tau_k$. By (91) in Lemma 5, we have $d_m \leq (K_0 + 1)m + K_0$.

Given $(s_{t-m}, a_{t-m}, s_{t-m+1})$ and θ_{k-d_m} , we construct an auxiliary Markov chain:

$$s_{t-m} \xrightarrow{\theta_{k-d_m}} a_{t-m} \xrightarrow{\mathcal{P}} s_{t-m+1} \xrightarrow{\theta_{k-d_m}} \tilde{a}_{t-m+1} \cdots \tilde{s}_{t-1} \xrightarrow{\theta_{k-d_m}} \tilde{a}_{t-1} \xrightarrow{\mathcal{P}} \tilde{s}_t \xrightarrow{\theta_{k-d_m}} \tilde{a}_t \xrightarrow{\mathcal{P}} \tilde{s}_{t+1}.$$

First we have

$$\begin{aligned} & \left\langle \nabla J(\theta_k), \delta(x_{(k)}, \theta_k) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \nabla J(\theta_k) \right\rangle \\ &= \left\langle \nabla J(\theta_k), \delta(x_{(k)}, \theta_k) \left(\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)}) \right) \right\rangle \\ & \quad + \left\langle \nabla J(\theta_k), \delta(x_{(k)}, \theta_k) \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)}) - \nabla J(\theta_k) \right\rangle. \end{aligned} \quad (106)$$

We bound the first term in (106) as

$$\begin{aligned} & \left\langle \nabla J(\theta_k), \delta(x_{(k)}, \theta_k) \left(\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)}) \right) \right\rangle \\ & \geq -\|\nabla J(\theta_k)\|_2 \|\delta(x_{(k)}, \theta_k)\|_2 \|\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)})\|_2 \\ & \geq -L_V \|\delta(x_{(k)}, \theta_k)\|_2 \|\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)})\|_2 \\ & \geq -L_V C_\delta \|\psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)})\|_2 \\ & \geq -L_V C_\delta L_\psi \|\theta_{k-\tau_k} - \theta_{k-d_m}\|_2, \end{aligned} \quad (107)$$

where the last inequality follows Assumption 3, and the second last inequality follows the fact that

$$|\delta(x, \theta)| \leq |r(x)| + \gamma |V_{\pi_\theta}(s')| + |V_{\pi_\theta}(s)| \leq r_{\max} + (1 + \gamma)R_{\max} \leq C_\delta.$$

Substituting (107) into (106) gives

$$\begin{aligned} & \left\langle \nabla J(\theta_k), \delta(x_{(k)}, \theta_k) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \nabla J(\theta_k) \right\rangle \\ & \geq -L_V C_\delta L_\psi \|\theta_{k-\tau_k} - \theta_{k-d_m}\|_2 + \left\langle \nabla J(\theta_k), \delta(x_{(k)}, \theta_k) \psi_{\theta_{k-d_m}}(s_{(k)}, a_{(k)}) - \nabla J(\theta_k) \right\rangle. \end{aligned} \quad (108)$$

Then we start to bound the second term in (108). For brevity, we define

$$\Delta_3(x, \theta) := \left\langle \nabla J(\theta), \delta(x, \theta) \psi_{\theta_{k-d_m}}(s, a) - \nabla J(\theta) \right\rangle.$$

Throughout the following proof, we use θ, θ', x and \tilde{x} as shorthand notations of $\theta_k, \theta_{k-d_m}, x_t$ and \tilde{x}_t respectively.

We decompose $\Delta_3(x, \theta)$ as

$$\Delta_3(x, \theta) = \underbrace{\Delta_3(x, \theta) - \Delta_3(x, \theta')}_{I_1} + \underbrace{\Delta_3(x, \theta') - \Delta_3(\tilde{x}, \theta')}_{I_2} + \underbrace{\Delta_3(\tilde{x}, \theta')}_{I_3}.$$

We first bound I_1 as

$$\begin{aligned} |I_1| &= |\Delta_3(x, \theta) - \Delta_3(x, \theta')| \\ &= \left| \langle \nabla J(\theta), \delta(x, \theta) \psi_{\theta'}(s, a) \rangle - \|\nabla J(\theta)\|_2^2 - \langle \nabla J(\theta'), \delta(x, \theta') \psi_{\theta'}(s, a) \rangle + \|\nabla J(\theta')\|_2^2 \right| \\ &\leq \left| \langle \nabla J(\theta), \delta(x, \theta) \psi_{\theta'}(s, a) \rangle - \langle \nabla J(\theta'), \delta(x, \theta') \psi_{\theta'}(s, a) \rangle \right| + \left| \|\nabla J(\theta')\|_2^2 - \|\nabla J(\theta)\|_2^2 \right| \\ &\leq \left| \langle \nabla J(\theta), \delta(x, \theta) \psi_{\theta'}(s, a) \rangle - \langle \nabla J(\theta'), \delta(x, \theta') \psi_{\theta'}(s, a) \rangle \right| + \|\nabla J(\theta') + \nabla J(\theta)\|_2 \|\nabla J(\theta') - \nabla J(\theta)\|_2 \\ &\leq \left| \langle \nabla J(\theta), \delta(x, \theta) \psi_{\theta'}(s, a) \rangle - \langle \nabla J(\theta'), \delta(x, \theta') \psi_{\theta'}(s, a) \rangle \right| + 2L_V L_J \|\theta - \theta'\|_2, \end{aligned} \quad (109)$$

where the last equality is due to L_V -Lipschitz of value function and L_J -Lipschitz of policy gradient. We bound the first term in (109) as

$$\begin{aligned}
& |\langle \nabla J(\theta), \delta(x, \theta) \psi_{\theta'}(s, a) \rangle - \langle \nabla J(\theta'), \delta(x, \theta') \psi_{\theta'}(s, a) \rangle| \\
& \leq |\langle \nabla J(\theta), \delta(x, \theta) \psi_{\theta'}(s, a) \rangle - \langle \nabla J(\theta), \delta(x, \theta') \psi_{\theta'}(s, a) \rangle| \\
& \quad + |\langle \nabla J(\theta), \delta(x, \theta') \psi_{\theta'}(s, a) \rangle - \langle \nabla J(\theta'), \delta(x, \theta') \psi_{\theta'}(s, a) \rangle| \\
& = |\langle \nabla J(\theta), (\delta(x, \theta) - \delta(x, \theta')) \psi_{\theta'}(s, a) \rangle| + |\langle \nabla J(\theta) - \nabla J(\theta'), \delta(x, \theta') \psi_{\theta'}(s, a) \rangle| \\
& \leq L_V C_\psi |\delta(x, \theta) - \delta(x, \theta')| + C_p \|\nabla J(\theta) - \nabla J(\theta')\|_2 \\
& = L_V C_\psi |\gamma(V_{\pi_\theta}(s') - V_{\pi_{\theta'}}(s')) + V_{\pi_{\theta'}}(s) - V_{\pi_\theta}(s)| + C_p \|\nabla J(\theta) - \nabla J(\theta')\|_2 \\
& \leq L_V C_\psi (\gamma |V_{\pi_\theta}(s') - V_{\pi_{\theta'}}(s')| + |V_{\pi_{\theta'}}(s) - V_{\pi_\theta}(s)|) + C_p \|\nabla J(\theta) - \nabla J(\theta')\|_2 \\
& \leq L_V C_\psi (\gamma L_V \|\theta - \theta'\|_2 + L_V \|\theta' - \theta\|) + C_p L_J \|\theta - \theta'\|_2 \\
& = (C_p L_J + (1 + \gamma) L_V^2 C_\psi) \|\theta - \theta'\|_2.
\end{aligned}$$

Substituting the above inequality into (109) gives the lower bound of I_1 :

$$I_1 \geq -(C_p L_J + (1 + \gamma) L_V^2 C_\psi + 2 L_V L_J) \|\theta - \theta'\|_2.$$

First we bound I_2 as

$$\begin{aligned}
\mathbb{E}[I_2 | \theta', s_{t-m+1}] &= \mathbb{E}[\Delta_3(x, \theta') - \Delta_3(\tilde{x}, \theta') | \theta', s_{t-m+1}] \\
&\geq -|\mathbb{E}[\Delta_3(x, \theta') | \theta', s_{t-m+1}] - \mathbb{E}[\Delta_3(\tilde{x}, \theta') | \theta', s_{t-m+1}]| \\
&\geq -\sup_x |\Delta_3(x, \theta')| \|\mathbb{P}(x \in \cdot | \theta', s_{t-m+1}) - \mathbb{P}(\tilde{x} \in \cdot | \theta', s_{t-m+1})\|_{TV} \\
&\geq -2 L_V (C_p + L_V) d_{TV}(\mathbb{P}(x \in \cdot | \theta', s_{t-m+1}), \mathbb{P}(\tilde{x} \in \cdot | \theta', s_{t-m+1})) \\
&\geq -L_V (C_p + L_V) |\mathcal{A}| L_\pi \sum_{i=\tau_k}^{d_m} \mathbb{E}[\|\theta_{k-i} - \theta_{k-d_m}\|_2 | \theta', s_{t-m+1}], \quad (110)
\end{aligned}$$

where the second inequality is due to the definition of TV norm, the last inequality is due to (27) in Lemma 3, and the second last inequality follows the fact that

$$|\Delta_3(x, \theta')| \leq \|\nabla J(\theta)\|_2 (\|\delta(x, \theta) \psi_{\theta_{k-d_m}}(s, a)\|_2 + \|\nabla J(\theta)\|_2) \leq L_V (C_p + L_V). \quad (111)$$

Taking total expectation on both sides of (110) yields

$$\mathbb{E}[I_2] \geq -L_V (C_p + L_V) |\mathcal{A}| L_\pi \sum_{i=\tau_k}^{d_m} \mathbb{E}[\|\theta_{k-i} - \theta_{k-d_m}\|_2].$$

Define $\bar{x} := (\bar{s}, \bar{a}, \bar{s}')$, where $\bar{s} \sim \mu_{\theta'}$, $\bar{a} \sim \pi_{\theta'}$ and $\bar{s}' \sim \mathcal{P}$. Then we can further decompose I_3 as

$$\mathbb{E}[I_3 | \theta', s_{t-m+1}] = \mathbb{E}[\Delta_3(\tilde{x}, \theta') - \Delta_3(\bar{x}, \theta') | \theta', s_{t-m+1}] + \mathbb{E}[\Delta_3(\bar{x}, \theta') | \theta', s_{t-m+1}]. \quad (112)$$

The first term in (112) can be bounded as

$$\begin{aligned}
& \mathbb{E}[\Delta_3(\tilde{x}, \theta') - \Delta_3(\bar{x}, \theta') | \theta', s_{t-m+1}] \\
& \geq -|\mathbb{E}[\Delta_3(\tilde{x}, \theta') | \theta', s_{t-m+1}] - \mathbb{E}[\Delta_3(\bar{x}, \theta') | \theta', s_{t-m+1}]| \\
& \geq -\sup_x |\Delta_3(x, \theta')| \|\mathbb{P}(\tilde{x} \in \cdot | \theta', s_{t-m+1}) - \mathbb{P}(\bar{x} \in \cdot | \theta', s_{t-m+1})\|_{TV} \\
& \geq -2 L_V (C_p + L_V) d_{TV}(\mathbb{P}(\tilde{x} \in \cdot | \theta', s_{t-m+1}), \mathbb{P}(\bar{x} \in \cdot | \theta', s_{t-m+1})) \\
& = -2 L_V (C_p + L_V) d_{TV}(\mathbb{P}(\tilde{x} \in \cdot | \theta', s_{t-m+1}), \mu_{\theta'} \otimes \pi_{\theta'} \otimes \mathcal{P}) \quad (113)
\end{aligned}$$

where the second inequality follows the definition of TV-norm, and the third one follows (111).

The auxiliary Markov chain with policy $\pi_{\theta'}$ starts from initial state s_{t-m+1} , and \tilde{s}_t is the $(m-1)$ th state on the chain. Following Lemma 1, we have:

$$\begin{aligned}
d_{TV}(\mathbb{P}(\tilde{x} \in \cdot | \theta', s_{t-m+1}), \mu_{\theta'} \otimes \pi_{\theta'} \otimes \mathcal{P}) &= d_{TV}(\mathbb{P}((\tilde{s}_t, \tilde{a}_t, \tilde{s}_{t+1}) \in \cdot | \theta', s_{t-m+1}), \mu_{\theta'} \otimes \pi_{\theta'} \otimes \mathcal{P}) \\
&\leq \kappa \rho^{m-1}.
\end{aligned}$$

Substituting the last inequality into (113) yields

$$\mathbb{E} [\Delta_3(\tilde{x}, \theta') - \Delta_3(\bar{x}, \theta') | \theta', s_{t-m+1}] \geq -2L_V(C_p + L_V)\kappa\rho^{m-1}. \quad (114)$$

The second term in (112) can be bounded as

$$\begin{aligned} \mathbb{E}[\Delta_3(\bar{x}, \theta') | \theta', s_{t-m+1}] &= \mathbb{E}[\langle \nabla J(\theta'), \delta(\bar{x}, \theta') \psi_{\theta'}(\bar{s}, \bar{a}) - \nabla J(\theta') \rangle | \theta', s_{t-m+1}] \\ &\geq -L_V \mathbb{E} \|\delta(\bar{x}, \theta') \psi_{\theta'}(\bar{s}, \bar{a}) | \theta', s_{t-m+1}] - \nabla J(\theta')\|_2 \\ &\geq -L_V R_{\max} C_\psi \|\mu_{\theta'} - d_{\theta'}\|_{TV} \\ &\geq -4L_V R_{\max} C_\psi \left(\log_\rho \kappa^{-1} + \frac{1}{1-\rho} \right) (1-\gamma) := -\epsilon_{sp} \end{aligned} \quad (115)$$

where the last inequality follows Lemma 2.

Substituting the lower bounds in (114) and (115) into (112) yields

$$\mathbb{E}[I_3 | \theta', s_{t-m+1}] \geq -2L_V(C_p + L_V)\kappa\rho^{m-1} - \epsilon_{sp}.$$

Taking total expectation on $\Delta_3(x, \theta)$ and collecting lower bounds of I_1, I_2, I_3 yield

$$\begin{aligned} \mathbb{E}[\Delta_3(x, \theta)] &\geq -(C_p L_J + (1+\gamma)L_V^2 C_\psi + 2L_V L_J) \mathbb{E} \|\theta_k - \theta_{k-d_m}\|_2 - \epsilon_{sp} \\ &\quad - L_V(C_p + L_V) |\mathcal{A}| L_\pi \sum_{i=\tau_k}^{d_m} \mathbb{E} \|\theta_{k-i} - \theta_{k-d_m}\|_2 - 2L_V(C_p + L_V)\kappa\rho^{m-1} \end{aligned}$$

Taking total expectation on (108) and substituting the above inequality into it yield

$$\begin{aligned} &\mathbb{E} \left\langle \nabla J(\theta_k), \delta(x_{(k)}, \theta_k) \psi_{\theta_{k-\tau_k}}(s_{(k)}, a_{(k)}) - \nabla J(\theta_k) \right\rangle \\ &\geq -D_6 \mathbb{E} \|\theta_{k-\tau_k} - \theta_{k-d_m}\|_2 - D_7 \mathbb{E} \|\theta_k - \theta_{k-d_m}\|_2 - D_8 \sum_{i=\tau_k}^{d_m} \mathbb{E} \|\theta_{k-i} - \theta_{k-d_m}\|_2 - D_9 \kappa \rho^{m-1} - \epsilon_{sp}, \end{aligned}$$

where $D_6 := L_V C_\delta L_\psi$, $D_7 := C_p L_J + (1+\gamma)L_V^2 C_\psi + 2L_V L_J$, $D_8 := L_V(C_p + L_V) |\mathcal{A}| L_\pi$, $D_9 := 2L_V(C_p + L_V)$. \square

D EXPERIMENT DETAILS

Hardware device. The tests on synthetic environment and CartPole was performed in a 16-core CPU computer. The test on Atari game was run in a 4 GPU computer.

Parameterization. For the synthetic environment, we used linear value function approximation and tabular softmax policy [38]. For CartPole, we used a 3-layer MLP with 128 neurons and sigmoid activation function in each layer. The first two layers are shared for both actor and critic network. For the Atari seaquest game, we used a convolution-LSTM network. For network details, see [44].

Hyper-parameters. For the synthetic environment tests, we run Algorithm 1 with actor step size $\alpha_k = \frac{0.05}{(1+k)^{0.6}}$ and critic step size $\beta_k = \frac{0.05}{(1+k)^{0.4}}$. In tests of CartPole, we run Algorithm 1 with a minibatch of 20 samples. We update the actor network with a step size of $\alpha_k = \frac{0.01}{(1+k)^{0.6}}$ and critic network with a step size of $\beta_k = \frac{0.01}{(1+k)^{0.4}}$. See Table 1 for hyper-parameters to generate the Atari game results in Figure 4.

Hyper-parameters	Value
Number of workers	16
Optimizer	Adam
Step size	0.00015
Batch size	20
Discount factor	0.99
Entropy coefficient	0.01
Frame size	80×80
Frame skip rate	4
Grayscale	Yes
Training reward clipping	[-1,1]

Table 1: Hyper-parameters of A3C-TD(0) in the Atari seaquest game.