

# Beyond Exact Gradients: Convergence of Stochastic Soft-Max Policy Gradient Methods with Entropy Regularization

Yuhao Ding, Junzi Zhang, Hyunin Lee, and Javad Lavaei

**Abstract**—Entropy regularization is an efficient technique for encouraging exploration and preventing a premature convergence of (vanilla) policy gradient methods in reinforcement learning (RL). However, the theoretical understanding of entropy regularized RL algorithms has been limited. In this paper, we revisit the classical entropy regularized policy gradient methods with the soft-max policy parametrization, whose convergence has so far only been established assuming access to exact gradient oracles. To go beyond this scenario, we propose the first set of (nearly) unbiased stochastic policy gradient estimators with trajectory-level entropy regularization, with one being an unbiased visitation measure-based estimator and the other one being a nearly unbiased yet more practical trajectory-based estimator. We prove that although the estimators themselves are unbounded in general due to the additional logarithmic policy rewards introduced by the entropy term, the variances are uniformly bounded. We then propose a two-phase stochastic policy gradient (PG) algorithm that uses a large batch size in the first phase to overcome the challenge of the stochastic approximation due to the non-coercive landscape, and uses a small batch size in the second phase by leveraging the curvature information around the optimal policy. We establish a global optimality convergence result and a sample complexity of  $\tilde{O}(\frac{1}{\epsilon^2})$  for the proposed algorithm. Our result is the first global convergence and sample complexity results for the stochastic entropy-regularized vanilla PG method.

**Index Terms**—Reinforcement learning, policy gradient, stochastic approximation

## I. INTRODUCTION

Entropy regularization is a popular technique to encourage exploration and prevent premature convergence for reinforcement learning (RL) algorithms. It was originally proposed in [1] to improve the performance of REINFORCE, a classical family of vanilla policy gradient (PG) methods widely used in practice. Since then, the entropy regularization technique has been applied to a large set of other RL algorithms, including actor-critic [2, 3], Q-learning [4, 5] and trust-region policy optimization methods [6]. It has been shown that the entropy regularization works satisfactorily with deep learning approximations for achieving an impressive empirical performance boost, provides a substantial improvement in exploration and robustness [3, 5, 7], and connects the policy gradient with

Q-learning under a one-step entropy regularization [4] or a trajectory-level KL regularization<sup>1</sup> [8].

There has been considerable interest in the theoretical understanding of how the entropy regularization exploits the geometry of the optimization landscape. In particular, it has been shown in [9, 10, 11] that entropy regularization makes the regularized objective behave similar to a local quadratic function and thus accelerates the convergence of entropy-regularized PG algorithms. When the exact entropy-regularized PG is available, a linear convergence rate has been established for the entropy-regularized PG algorithms with the natural PG (NPG) or policy mirror descent [10, 11] or without the NPG [9]. However, in practice, the agent does not have access to the exact entropy-regularized PG but only its stochastic estimation from the samples of trajectories. The advantages of entropy regularization have mostly been established for the exact gradient setting. It is not fully understood whether these advantages are only restricted to theoretical analysis in the exact gradient settings and whether any geometric property can be exploited to accelerate convergence to global optimality in the stochastic gradient settings. Recently, it is proven in [11] that the NPG with the entropy regularization has a sample complexity of  $\tilde{O}(\frac{1}{\epsilon^2})$  in the stochastic gradient settings, where the inexactness of the gradient can be reduced to the inexactness of the state-action value functions. However, the literature on the global optimality convergence and the sample complexity of the most fundamental PG, namely REINFORCE and its variants with regularizations, is still limited, despite its simplicity and popularity in practice. The work [9] has recently developed the first set of global convergence results for PG, which focuses on the soft-max policy parametrization by assuming access to exact PG evaluations. However, their result heavily relies on the access to the exact PG evaluations, and it has been shown that the geometric advantages existing in the exact gradient setting may not be preserved in the stochastic setting [12, 13]. It remains an open problem whether a global optimality convergence result and a low sample complexity can be obtained for the PG with entropy regularization in the stochastic gradient setting.

In this paper, we provide an affirmative answer to the above question. In particular, we revisit the classical entropy regularized (vanilla) policy gradient method proposed in the seminal work [1] under the soft-max policy parametrization.

Y. Ding and J. Lavaei are with the Department of Industrial Engineering and Operations Research, University of California, Berkeley, CA 94709 USA (e-mail: yuhao\_ding@berkeley.edu; lavaei@berkeley.edu).

J. Zhang is with Citadel Securities (work done prior to joining Citadel Securities), Chicago, IL 60603 USA (e-mail: saslasroyale@gmail.com).

H. Lee is with the Department of Mechanical Engineering, University of California, Berkeley, CA 94709 USA (e-mail: hyunin@berkeley.edu).

<sup>1</sup>Note that this is related to but different from the widely-used trajectory-level entropy regularization later introduced in [5].

We focus on the modern trajectory-level entropy regularization proposed in [5], which is shown to improve over the original one-step entropy regularization adopted in [1, 2] and [4]. Our contributions are summarized below:

- We begin by proposing two new entropy regularized stochastic PG estimators. The first one is an unbiased visitation measure-based estimator, whereas the second one is a nearly unbiased yet more practical trajectory-based estimator. These (nearly) unbiased stochastic PG estimators are the first likelihood-ratio-based estimators in the literature with a trajectory-level entropy regularization. We show that although the estimators themselves are unbounded in general due to the entropy-induced logarithmic policy rewards, the variances indeed remain uniformly bounded.
- One main challenge on extending the result in [9] to the stochastic PG setting is the non-coercive landscape<sup>2</sup> of the entropy-regularized RL. To overcome this challenge, we propose a two-phase stochastic PG algorithm that uses a large batch size in the first phase and uses a small batch size in the second phase. We establish a global optimality convergence result and a sample complexity of  $\tilde{\mathcal{O}}(\frac{1}{\epsilon^2})$  for the proposed algorithm under the softmax parameterization. Our result is the first to achieve the sample complexity of  $\tilde{\mathcal{O}}(\frac{1}{\epsilon^2})$  for the stochastic entropy-regularized vanilla PG method and matches the sample complexity of the natural PG [11] in terms of dependence on  $\epsilon$ .

#### A. Related work

It has been shown in [7] that the entropy-regularized RL formulation provides a substantial improvement in exploration and robustness. An actor-critic method is proposed in [3] which updates the policy towards the exponential of the new Q-function and projects the improved policy onto the desired set of policies in the policy improvement step. Instead of using the likelihood ratio gradient estimator [14], the gradient estimator for their policy improvement is based on the re-parameterization technique [3, Equation (13)] where the function approximation is inevitable and the theoretical analysis is challenging. In contrast, we focus on the stochastic PG method and classical likelihood ratio estimators.

Stochastic PG estimators with the original one-step entropy regularization have been proposed and adopted in [1, 2, 4]. For trajectory-level entropy regularization, an exact (visitation measure-based) PG formula has been derived in [15] and later re-derived in the soft-max policy parametrization setting in [9], while stochastic PG estimators have not been formally proposed or studied in the literature. [8] provides a stochastic PG estimator for the value function with a related but different regularization term: trajectory-level KL-divergence regularization. However, KL-divergence regularization is far more aggressive and less used in practice in reinforcement learning compared with the entropy regularization [16]).

<sup>2</sup>A continuous function  $f(x)$  that is defined on  $\mathbb{R}^n$  is called coercive if  $\lim_{\|x\| \rightarrow \infty} f(x) = +\infty$ .

The theoretical understanding of policy-based methods has received considerable attention recently [9, 10, 11, 16, 17, 18, 19, 20, 21, 22]. Several techniques have been developed to improve standard PG and achieve a linear convergence rate, such as adding entropy regularization [9, 10, 11, 16], exploiting natural geometries based on Bregman divergences leading to NPG or policy mirror descent [10, 11, 17], and using a geometry-aware normalized PG (GNPG) approach to exploit the non-uniformity of the value function [23]. For the stochastic policy optimization, the existing results have mostly focused on policy mirror ascent methods with the goal of reducing the stochastic analysis to the estimation of the Q-value function [10, 11], as well as incorporating variance reduction techniques to improve the sample complexity of the vanilla PG [24, 25]. The prior literature still lacks globally optimal convergence results and sample complexity for stochastic (vanilla) PG with the entropy regularization.

#### B. Notation

The set of real numbers is shown as  $\mathbb{R}$ .  $u \sim \mathcal{U}$  means that  $u$  is a random vector sampled from the distribution  $\mathcal{U}$ . We use  $|\mathcal{X}|$  to denote the cardinality of a finite set  $\mathcal{X}$ . The notations  $\mathbb{E}_\xi[\cdot]$  and  $\mathbb{E}[\cdot]$  refer to the expectation over the random variable  $\xi$  and over all of the randomness. The notation  $\text{Var}[\cdot]$  refers to the variance.  $\Delta(\mathcal{X})$  denotes the probability simplex over a finite set  $\mathcal{X}$ . For vectors  $x, y \in \mathbb{R}^d$ , let  $\|x\|_1$ ,  $\|x\|_2$  and  $\|x\|_\infty$  denote the  $\ell_1$ -norm,  $\ell_2$ -norm and  $\ell_\infty$ -norm. We use  $\langle x, y \rangle$  to denote the inner product. For a matrix  $A$ , the notation  $A \succeq 0$  means that  $A$  is positive semi-definite. Given a variable  $x$ , the notation  $a = \mathcal{O}(b(x))$  means that  $a \leq C \cdot b(x)$  for some constant  $C > 0$  that is independent of  $x$ . Similarly,  $a = \tilde{\mathcal{O}}(b(x))$  indicates that the previous inequality may also depend on the function  $\log(x)$ , that is,  $a \leq C \cdot b(x) \cdot \log(x)$ , where  $C > 0$  is again independent of  $x$ . We use  $\text{Geom}(x)$  to denote a geometric distribution with the parameter  $x$ . Let the notation  $\mathbb{1}_A$  denote the indicator function of an event  $A \subseteq \Omega$ , i.e.,  $\mathbb{1}_A(\omega) = 1$  if  $\omega \in A$  and  $\mathbb{1}_A(\omega) = 0$  otherwise. For a given stochastic algorithm, let  $\mathcal{F}_t$  denote the  $\sigma$ -field generated by the history of the algorithm up to the iteration  $t$ , just before the randomness at the iteration  $t$  is generated. We define  $\mathbb{E}^t := \mathbb{E}[\cdot | \mathcal{F}_t]$  as the expectation operator conditioned on the  $\sigma$ -field  $\mathcal{F}_t$ .

## II. PRELIMINARIES

**Markov decision processes.** RL is generally modeled as a discounted Markov decision process (MDP) defined by a tuple  $(\mathcal{S}, \mathcal{A}, \mathbb{P}, r, \gamma)$ . Here,  $\mathcal{S}$  and  $\mathcal{A}$  denote the finite state and action spaces;  $\mathbb{P}(s'|s, a)$  is the probability that the agent transits from the state  $s$  to the state  $s'$  under the action  $a \in \mathcal{A}$ ;  $r(s, a)$  is the reward function, i.e., the agent obtains the reward  $r(s_h, a_h)$  after it takes the action  $a_h$  at the state  $s_h$  at time  $h$ ;  $\gamma \in (0, 1)$  is the discount factor. Without loss of generality, we assume that  $r(s, a) \in [0, \bar{r}]$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$ . The policy  $\pi(a|s)$  at the state  $s$  is usually represented by a conditional probability distribution  $\pi_\theta(a|s)$  associated to the parameter  $\theta \in \mathbb{R}^d$ . Let  $\tau^\infty = \{s_0, a_0, s_1, a_1, \dots\}$  denote the data of a sampled trajectory under policy  $\pi_\theta$  with the probability distribution over the trajectory as  $p(\tau^\infty = \cdot | \theta, \rho) :=$

$\rho(s_0) \prod_{h=1}^{\infty} \mathbb{P}(s_{h+1}|s_h, a_h) \pi_{\theta}(a_h|s_h)$ , where  $\rho \in \Delta(\mathcal{S})$  is the probability distribution of the initial state  $s_0$ .

**Value functions and Q-functions.** Given a policy  $\pi$ , one can define the state-action value function  $Q^{\pi} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  as

$$Q^{\pi}(s, a) := \mathbb{E}_{\substack{a_h \sim \pi(\cdot|s_h) \\ s_{h+1} \sim \mathbb{P}(\cdot|s_h, a_h)}} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \middle| s_0 = s, a_0 = a \right].$$

The state-value function  $V^{\pi} : \mathcal{S} \rightarrow \mathbb{R}$  and the advantage function  $A^{\pi} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  can be defined as  $V^{\pi}(s) := \mathbb{E}_{a \sim \pi(\cdot|s)} [Q^{\pi}(s, a)]$ ,  $A^{\pi}(s, a) := Q^{\pi}(s, a) - V^{\pi}(s)$ . The goal is to find an optimal policy in the underlying policy class that maximizes the expected discounted return under the initial state distribution, namely,  $\max_{\theta \in \mathbb{R}^d} V^{\pi_{\theta}}(\rho) := \mathbb{E}_{s_0 \sim \rho} [V^{\pi_{\theta}}(s_0)]$ . For the notational convenience, we will denote  $V^{\pi_{\theta}}(\rho)$  by the shorthand notation  $V^{\theta}(\rho)$ .

**Exploratory initial distribution.** The discounted state visitation distribution  $d_{s_0}^{\pi}$  is defined as  $d_{s_0}^{\pi}(s) := (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}(s_h = s | s_0, \pi)$ , where  $\mathbb{P}(s_h = s | s_0, \pi)$  is the state visitation probability that  $s_h$  is equal to  $s$  under the policy  $\pi$  starting from the state  $s_0$ . The discounted state visitation distribution under the initial distribution  $\rho$  is defined as  $d_{\rho}^{\pi}(s) := \mathbb{E}_{s_0 \sim \rho} [d_{s_0}^{\pi}(s)]$ . Furthermore, the state-action visitation distribution induced by  $\pi$  and the initial state distribution  $\rho$  is defined as  $v_{\rho}^{\pi}(s, a) := d_{\rho}^{\pi}(s) \pi(a|s)$ , which can also be written as  $v_{\rho}^{\pi}(s, a) := (1 - \gamma) \mathbb{E}_{s_0 \sim \rho} \sum_{h=0}^{\infty} \gamma^h \mathbb{P}(s_h = s, a_h = a | s_0, \pi)$ , where  $\mathbb{P}(s_h = s, a_h = a | s_0, \pi)$  is the state-action visitation probability that  $s_h = s$  and  $a_h = a$  under  $\pi$  starting from the state  $s_0$ . To facilitate the presentation of the main results of the paper, we assume that the state distribution  $\rho$  for the performance measure is exploratory [9, 19], i.e.,  $\rho(\cdot)$  adequately covers the entire state distribution:

*Assumption 1:* The state distribution  $\rho$  satisfies  $\rho(s) > 0$  for all  $s \in \mathcal{S}$ .

In practice, when the above assumption is not satisfied, we can optimize under another initial distribution  $\mu$ , i.e., the gradient is taken with respect to the optimization measure  $\mu$ , where  $\mu$  is usually chosen as an exploratory initial distribution that adequately covers the state distribution of some optimal policy. It is shown in [16] that the difficulty of the exploration problem faced by PG algorithms can be captured through the distribution mismatch coefficient defined as  $\left\| \frac{d_{\rho}^{\pi}}{\mu} \right\|_{\infty}$ , where  $\frac{d_{\rho}^{\pi}}{\mu}$  denotes component-wise division.

**Soft-max policy parameterization.** In this work, we consider the soft-max parameterization – a widely adopted scheme that naturally ensures that the policy lies in the probability simplex. Specifically, for an unconstrained parameter  $\theta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ ,  $\pi_{\theta}(a|s)$  is chosen to be  $\frac{\exp(\theta_{s,a})}{\sum_{a' \in \mathcal{A}} \exp(\theta_{s,a'})}$ . The soft-max parameterization is generally used for MDPs with finite state and action spaces. It is complete in the sense that every stochastic policy can be represented by this class. For the soft-max parameterization, it can be shown that the gradient and Hessian of the function  $\log \pi_{\theta}(a|s)$  are bounded, i.e., for all  $\theta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ ,  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$ , we have:  $\|\nabla \log \pi_{\theta}(a|s)\|_2 \leq 2$ ,  $\|\nabla^2 \log \pi_{\theta}(a|s)\|_2 \leq 1$ .

**RL with entropy regularization.** Entropy is a commonly used regularization in RL to promote exploration and discourage premature convergence to suboptimal policies [5, 8, 26].

It is far less aggressive in penalizing small probabilities, in comparison to other common regularizations such as log barrier functions [16]. In the entropy-regularized RL (also known as maximum entropy RL), near-deterministic policies are penalized, which is achieved by modifying the value function to

$$V_{\lambda}^{\pi}(\rho) = V^{\pi}(\rho) + \lambda \mathbb{H}(\rho, \pi), \quad (1)$$

where  $\lambda \geq 0$  determines the strength of the penalty and  $\mathbb{H}(\rho, \pi)$  stands for the discounted entropy defined as

$$\mathbb{H}(\rho, \pi) := \mathbb{E}_{\substack{s_0 \sim \rho, a_t \sim \pi(\cdot|s_t) \\ s_{t+1} \sim \mathbb{P}(\cdot|s_t, a_t)}} \left[ \sum_{t=0}^{\infty} -\gamma^t \log \pi(a_t | s_t) \right].$$

Equivalently,  $V_{\lambda}^{\pi}(\rho)$  can be viewed as the weighted value function of  $\pi$  by adjusting the instantaneous reward to be policy-dependent regularized version as  $r^{\lambda}(s, a) := r(s, a) - \lambda \log \pi(a|s)$ , for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . We also define  $V_{\lambda}^{\pi}(s)$  analogously when the initial state is fixed at a given state  $s \in \mathcal{S}$ . The regularized Q-function  $Q_{\lambda}^{\pi}$  of a policy  $\pi$ , also known as the soft Q-function, is related to  $V_{\lambda}^{\pi}$  as (for every  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$ )

$$\begin{aligned} Q_{\lambda}^{\pi}(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V_{\lambda}^{\pi}(s')], \\ V_{\lambda}^{\pi}(s) &= \mathbb{E}_{a \sim \pi(\cdot|s)} [-\lambda \log \pi(a|s) + Q_{\lambda}^{\pi}(s, a)]. \end{aligned}$$

**Bias due to entropy regularization.** Due to the presence of regularization, the optimal solution will be biased with the bias disappearing as  $\lambda \rightarrow 0$ . More precisely, the optimal policy  $\pi_{\lambda}^*$  of the entropy-regularized problem could also be nearly optimal in terms of the unregularized objective function, as long as the regularization parameter  $\lambda$  is chosen to be small. Denote by  $\pi^*$  and  $\pi_{\lambda}^*$  the policies that maximize the objective function and the entropy-regularized objective function with the regularization parameter  $\lambda$ , respectively. Let  $V^*$  and  $V_{\lambda}^*$  represent the resulting optimal objective value function and the optimal regularized objective value function. [11] shows a simple but crucial connection between  $\pi^*$  and  $\pi_{\lambda}^*$  via the following sandwich bound:

$$V^{\pi_{\lambda}^*}(\rho) \leq V^{\pi^*}(\rho) \leq V^{\pi_{\lambda}^*}(\rho) + \frac{\lambda \log |\mathcal{A}|}{1 - \gamma},$$

which holds for all initial distribution  $\rho$ .

### III. STOCHASTIC PG METHODS FOR ENTROPY REGULARIZED RL

#### A. Review: Exact PG methods

The PG method is one of the most popular approaches for a direct policy search in RL [27]. The vanilla PG with exact gradient information and the entropy regularization is summarized in Algorithm 1.

---

#### Algorithm 1 Exact PG method

---

- 1: **Inputs:**  $\{\eta_t\}_{t=1}^T, \theta_1$ .
  - 2: **for**  $t = 1, 2, \dots, T - 1$  **do**
  - 3:    $\theta_{t+1} = \theta_t + \eta_t \nabla V_{\lambda}^{\theta_t}(\rho)$ .
  - 4: **end for**
  - 5: **Outputs:**  $\theta_T$ .
-

The uniform boundedness of the reward function  $r$  implies that the absolute value of the entropy-regularized state-value function and Q-value function are bounded.

*Lemma 1 ([9]):*  $V_\lambda^\theta(s) \leq \frac{\bar{r} + \lambda \log |\mathcal{A}|}{1 - \gamma}$  and  $Q_\lambda^\pi(s, a) \leq \frac{\bar{r} + \lambda \log |\mathcal{A}|}{1 - \gamma}$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$  and  $\theta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ .

Under the soft-max policy parameterization, one can obtain the following expression for the gradient of  $V_\lambda^\pi(s)$  with respect to the policy parameter  $\theta$ :

*Lemma 2 (Proposition 2 in [28]):* The entropy regularized PG with respect to  $\theta$  is

$$\nabla V_\lambda^\theta(\rho) = \frac{1}{1 - \gamma} \mathbb{E}_{s, a \sim \nu_\rho^{\pi_\theta}} \left[ \nabla_\theta \log \pi_\theta(a|s) (Q_\lambda^\theta(s, a) - \lambda \log \pi_\theta(a|s)) \right], \quad (2)$$

where

$$\frac{\partial \log \pi_\theta(a|s)}{\partial \theta_{s', a'}} = \begin{cases} -\pi_\theta(a'|s') \pi_\theta(a|s), & (s', a') \neq (s, a), \\ \pi_\theta(a|s) - \pi_\theta(a|s) \pi_\theta(a|s), & (s', a') = (s, a). \end{cases}$$

Furthermore, the entropy regularized PG is bounded, i.e.,  $\|\nabla V_\lambda^\theta(\rho)\| \leq G$  for all  $\rho \in \Delta(\mathcal{S})$  and  $\theta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ , where  $G := \frac{2(\bar{r} + \lambda \log |\mathcal{A}|)}{(1 - \gamma)^2}$ .

In addition, it is shown that the PG  $\nabla V_\lambda^\theta(\rho)$  is Lipschitz continuous.

*Lemma 3 (Lemmas 7 and 14 in [9]):* The PG  $\nabla V_\lambda^\theta(\rho)$  is Lipschitz continuous with some constant  $L > 0$ , i.e.,  $\|\nabla V_\lambda^{\theta^1}(\rho) - \nabla V_\lambda^{\theta^2}(\rho)\| \leq L \cdot \|\theta^1 - \theta^2\|$ , for all  $\theta^1, \theta^2 \in \mathbb{R}^d$ , where the value of the Lipschitz constant  $L$  is defined as  $L := \frac{8\bar{r} + \lambda(4 + 8 \log |\mathcal{A}|)}{(1 - \gamma)^3}$ .

**Challenges for designing entropy regularized PG estimators.** Existing works either consider one-step entropy regularization [2, 14], KL divergence [8], or the re-parametrization technique [3, 5] (which introduces approximation errors that are difficult to quantify exactly). In general, the regularized reward  $r - \lambda \log \pi_\theta$  is policy-dependent and unbounded even though the original reward  $r$  is uniformly bounded. Hence, the existing estimators for the un-regularized setting must be modified to account for the policy-dependency and unboundedness while maintaining the essential properties of (nearly) unbiasedness and bounded variances. In the subsequent sections, we propose two (nearly) unbiased estimators and show that although the estimators may be unbounded due to unbounded regularized rewards, the variances are indeed bounded. The proofs of the results in this section can be found in Section A of the supplemental materials.

## B. Sampling the unbiased PG

It results from (2) that in order to obtain an unbiased sample of  $\nabla V_\lambda^\theta(\rho)$ , we need to first draw a state-action pair  $(s, a)$  from the distribution  $\nu_\rho^{\pi_\theta}(\cdot, \cdot)$  and then obtain an unbiased estimate of the action-value function  $Q_\lambda^\theta(s, a)$ . For the standard discounted infinite-horizon RL setting with bounded reward functions, [29] proposes an unbiased estimate of the PG using the random horizon with a geometric distribution and the Monte-Carlo rollouts of finite horizons. However, their result cannot be immediately applied to the entropy-regularized RL setting since the entropy-regularized instantaneous reward

$r(s, a) - \lambda \log \pi(a|s)$  could be unbounded when  $\pi(a|s) \rightarrow 0$ . Fortunately, we can still show that an unbiased PG estimator with the bounded variance for the entropy regularized RL can be obtained in a similar fashion as in [29]. In particular, we will use a random horizon that follows a certain geometric distribution in the sampling process. To ensure that  $(s_H, a_H) \sim \nu_\rho^{\pi_\theta}(s, a)$ , we will use the last sample  $(s_H, a_H)$  of a finite sample trajectory  $(s_0, a_0, s_1, a_1, \dots, s_H, a_H)$  to be the sample at which  $Q_\lambda^\theta(\cdot, \cdot)$  is evaluated, where the horizon  $H \sim \text{Geom}(1 - \gamma)$ . Moreover, given  $(s_H, a_H)$ , we will perform Monte-Carlo rollouts for another trajectory with the horizon  $H' \sim \text{Geom}(1 - \gamma^{1/2})$  independent of  $H$ , and estimate the advantage function value  $Q_\lambda^\theta(s, a)$  along the trajectory  $(s'_0, a'_0, \dots, s'_{H'}, a'_{H'})$  with  $s'_0 = s, a'_0 = a$  as follows:

$$\hat{Q}_\lambda^\theta(s, a) = r(s'_0, a'_0) + \sum_{t=1}^{H'} \gamma^{t/2} \cdot (r(s'_t, a'_t) - \lambda \log \pi_\theta(a'_t|s'_t)). \quad (3)$$

The subroutines of sampling one pair  $(s, a)$  from  $\nu_\rho^{\pi_\theta}(\cdot, \cdot)$ , estimating  $\hat{Q}_\lambda^\theta(s, a)$ , and estimating  $\hat{V}_\lambda^\theta(s)$  are summarized as **Sam-SA** and **Est-EntQ** in Algorithms 2 and 3, respectively.

---

**Algorithm 2 Sam-SA:** Sample for  $s, a \sim \nu_\rho^{\pi_\theta}(\cdot, \cdot)$

---

- 1: **Inputs:**  $\rho, \theta, \gamma$ .
  - 2: Draw  $H \sim \text{Geom}(1 - \gamma)$ .
  - 3: Draw  $s_0 \sim \rho$  and  $a_0 \sim \pi_\theta(\cdot|s_0)$ .
  - 4: **for**  $h = 1, 2, \dots, H - 1$  **do**
  - 5:     Simulate the next state  $s_{h+1} \sim \mathbb{P}(\cdot|s_h, a_h)$  and action  $a_{h+1} \sim \pi_{\theta_t}(\cdot|s_{h+1})$ .
  - 6: **end for**
  - 7: **Outputs:**  $s_H, a_H$ .
- 

---

**Algorithm 3 Est-EntQ:** Unbiasedly estimating entropy-regularized Q function

---

- 1: **Inputs:**  $s, a, \gamma, \lambda$  and  $\theta$ .
  - 2: Initialize  $s_0 \leftarrow s, a_0 \leftarrow a, \hat{Q} \leftarrow r(s_0, a_0)$ .
  - 3: Draw  $H' \sim \text{Geom}(1 - \gamma^{1/2})$ .
  - 4: **for**  $h = 0, 1, \dots, H' - 1$  **do**
  - 5:     Simulate the next state  $s_{h+1} \sim \mathbb{P}(\cdot|s_h, a_h)$  and action  $a_{h+1} \sim \pi_\theta(\cdot|s_{h+1})$ .
  - 6:     Collect the instantaneous reward  $r(s_{h+1}, a_{h+1}) - \lambda \log \pi_\theta(a_{h+1}|s_{h+1})$  and add to the value  $\hat{Q}$ :  $\hat{Q} \leftarrow \hat{Q} + \gamma^{(h+1)/2} (r(s_{h+1}, a_{h+1}) - \lambda \log \pi_\theta(a_{h+1}|s_{h+1}))$ .
  - 7: **end for**
  - 8: **Outputs:**  $\hat{Q}$ .
- 

Motivated by the form of PG in (2), we propose the following stochastic estimator:

$$\hat{\nabla} V_\lambda^\theta(\rho) = \frac{1}{1 - \gamma} \nabla_\theta \log \pi_\theta(a_H|s_H) \left( \hat{Q}_\lambda^\theta(s_H, a_H) - \lambda \log \pi_\theta(a_H|s_H) \right), \quad (4)$$

where  $s_H, a_H \leftarrow \text{Sam-SA}(\rho, \theta, \gamma)$  and  $\hat{Q}_\lambda^\theta$  is defined in (3). The following lemma shows that the stochastic PG (4) is an unbiased estimator of  $\nabla V_\lambda^\theta(\rho)$ .

*Lemma 4:* For  $\hat{\nabla}V_\lambda^\theta(\rho)$  defined in (4), we have  $\mathbb{E}[\hat{\nabla}V_\lambda^\theta(\rho)] = \nabla V_\lambda^\theta(\rho)$ .

The next lemma shows that the proposed PG estimator  $\hat{\nabla}V_\lambda^\theta(\rho)$  has a bounded variance even if it is unbounded when  $\pi_\theta$  approaches a deterministic policy.

*Lemma 5:* For  $\hat{\nabla}V_\lambda^\theta(\rho)$  defined in (4), we have  $\text{Var}[\hat{\nabla}V_\lambda^\theta(\rho)] \leq \sigma^2$ , where  $\sigma^2 = \frac{8}{(1-\gamma)^2} \left( \frac{\bar{r}^2 + (\lambda \log |\mathcal{A}|)^2}{(1-\gamma^{1/2})^2} \right)$  and  $\bar{r}$  is the upper bound of the reward.

### C. Sampling the trajectory-based PG

Compared to the unbiased PG with a random horizon in (4), a more practical PG estimator is the trajectory-based PG. To derive the trajectory-based PG for the entropy-regularized RL, we first notice that the gradient  $\nabla V_\lambda^\theta(\rho)$  can also be written as

$$\nabla V_\lambda^\theta(\rho) = \mathbb{E} \left[ \left( \sum_{t=0}^{\infty} \nabla \log \pi_\theta(a_t | s_t) \right) \left( \sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) - \lambda \log \pi_\theta(a_t | s_t)) \right) \right],$$

where the expectation is taken over the trajectory distribution, i.e.,  $\tau^\infty \sim p(\tau^\infty = |\theta)$ .

Since the distribution  $p(\tau^\infty = |\theta)$  is unknown,  $\nabla V_\lambda^\theta(\rho)$  needs to be estimated from samples. The trajectory-based estimators include REINFORCE [14], PGT [30] and GPOMDP [31]. In practice, the truncated versions of these trajectory-based PG estimators are used to approximate the infinite sum in the PG estimator. Let  $\tau^H = \{s_0, a_0, s_1, \dots, s_{H-1}, a_{H-1}, s_H\}$  denote the truncation of the full trajectory  $\tau^\infty =$  of length  $H$ . Then, with the commonly used truncated GPOMDP, the truncated PG estimator for  $\nabla V_\lambda^\theta$  can be written as:

$$\hat{\nabla}V_\lambda^{\theta, H}(\rho) = \sum_{h=0}^{H-1} \left( \sum_{j=0}^h \nabla \log \pi_\theta(a_j | s_j) \right) \gamma^h (r(s_h, a_h) - \lambda \log \pi_\theta(a_h | s_h)). \quad (5)$$

Due to the horizon truncation, the PG estimator (5) may no longer be unbiased, but its bias can be very small with a large horizon  $H$ .

*Lemma 6:* For  $\hat{\nabla}V_\lambda^{\theta, H}(\rho)$  defined in (5), we have

$$\left\| \mathbb{E}[\hat{\nabla}V_\lambda^{\theta, H}(\rho)] - \nabla V_\lambda^\theta(\rho) \right\|_2 \leq \frac{2(\bar{r} + \lambda \log |\mathcal{A}|)\gamma^H}{(1-\gamma)} \left( H + \frac{1}{1-\gamma} \right).$$

From Lemma 6, we can observe that the bias is proportional to  $\gamma^H$  and thus can be controlled to be arbitrarily small with a constant horizon up to some logarithmic term. We then show that the truncated PG estimator  $\hat{\nabla}V_\lambda^{\theta, H}$  has a bounded variance even if it may be unbounded when  $\pi_\theta$  approaches a deterministic policy.

*Lemma 7:* For  $\hat{\nabla}V_\lambda^{\theta, H}(\rho)$  defined in (5), we have

$$\text{Var}(\hat{\nabla}V_\lambda^{\theta, H}(\rho)) \leq \frac{12\bar{r}^2 + 24\lambda^2(\log |\mathcal{A}|)^2}{(1-\gamma)^4}.$$

### D. Batched PG algorithms

In practice, we can sample and compute a batch of independently and identically distributed PG estimators  $\{\hat{\nabla}V_\lambda^{\theta, i}(\rho)\}_{i=1}^B$  where  $B$  is the batch size, in order to reduce the estimation variance. To maximize the entropy-regularized objective

function (1), we can then update the policy parameter  $\theta$  by iteratively running gradient-ascent-based algorithms, i.e.,  $\theta_{t+1} = \theta_t + \frac{\eta_t}{B} \sum_{i=1}^B \hat{\nabla}V_\lambda^{\theta, i}(\rho)$ , where  $\eta_t > 0$  is the step size. The details of the unbiased PG algorithm with a random horizon for the entropy-regularized RL are provided in Algorithm 4.

---

#### Algorithm 4 Ent-RPG: Random-horizon PG for Entropy-regularized RL

---

```

1: Inputs:  $\rho, \lambda, \theta_1, B, T, \{\eta_t\}_{t=1}^T$ .
2: for  $t = 1, 2, \dots, T$  do
3:   for  $i = 1, 2, \dots, B$  do
4:      $s_{H_t}^i, a_{H_t}^i \leftarrow \text{SamSA}(\rho, \theta_t, \gamma)$ .
5:      $\hat{Q}_\lambda^{\theta_t, i} \leftarrow \text{Est-EntQ}(s_{H_t}^i, a_{H_t}^i, \theta_t, \gamma, \lambda)$ .
6:   end for
7:    $\theta_{t+1} \leftarrow \theta_t + \frac{\eta_t}{(1-\gamma)B} \sum_{i=1}^B \left[ \nabla_\theta \log \pi_{\theta_t}(a_{H_t}^i | s_{H_t}^i) \right. \\ \left. \left( \hat{Q}_\lambda^{\theta_t, i} - \lambda \log \pi_{\theta_t}(s_{H_t}^i | a_{H_t}^i) \right) \right]$ 
8: end for
9: Outputs:  $\theta_T$ .

```

---

*Remark 1:* For the simplicity of the presentation, we focus on deriving the stochastic PG estimator for the soft-max policy parameterization. However, our results in this section (and also the stationary point convergence result in Section IV-C below) can be easily extended to the general parameterization  $\pi_\theta$  as long as  $\|\nabla \log \pi_\theta(a|s)\|_2$  and  $\|\nabla^2 \log \pi_\theta(a|s)\|_2$  are bounded for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .

Due to space restrictions and in order to facilitate the presentation of the main ideas, we will mainly focus on the analysis of the unbiased PG estimator in (4) for the rest of the paper. Similar results hold for the trajectory-based PG estimator in (5) since its bias is exponentially small with respect to the horizon (see Lemma 6). The proofs of the results of this section can be found in Appendix A. We leave the formal discussion of these results as future work.

## IV. NON-COERCIVE LANDSCAPE

In this section, we first review some key results for the entropy-regularized RL with the exact PG and highlight the difficulty of generalizing these results to the stochastic PG setting, due to the non-coercive landscape.

### A. Review: Linear convergence with exact PG

A key result from [9] shows that, under the soft-max parameterization, the entropy-regularized value function  $V_\lambda^\theta(\rho)$  in (1) satisfies a non-uniform Łojasiewicz inequality as follows:

*Lemma 8 (Lemma 15 in [9]):* It holds that

$$\left\| \nabla V_\lambda^\theta(\rho) \right\|_2^2 \geq C(\theta) (V_\lambda^{\theta^*}(\rho) - V_\lambda^\theta(\rho)),$$

where

$$C(\theta) = \frac{2\lambda}{|\mathcal{S}|} \min_s \rho(s) \min_{s,a} \pi_\theta(a|s)^2 \left\| \frac{d_\rho^{\pi_\lambda^*}}{\rho} \right\|_\infty^{-1}.$$

Furthermore, it is shown in [9] that the action probabilities under the soft-max parameterization are uniformly bounded away from zero if the exact PG is available.

*Lemma 9 (Lemma 16 in [9]):* Using the exact PG (Algorithm 1) with  $\eta_t = \eta \leq \frac{2}{L}$  for the entropy regularized objective, it holds that  $\inf_{t \geq 1} \min_{s,a} \pi_{\theta_t}(a|s) > 0$ .

*Remark 2:* Note that by Algorithm 1,  $\inf_{t \geq 1} \min_{s,a} \pi_{\theta_t}(a|s)$  is only dependent on the initialization  $\theta_1$  and step-size  $\eta$  (apart from problem dependent constants). Hence hereafter we denote  $c_{\theta_1, \eta} = \inf_{t \geq 1} \min_{s,a} \pi_{\theta_t}(a|s)$ .

With Lemmas 3, 8 and 9, it is shown in Theorem 6 of [9] that the convergence rate for the entropy regularized PG is  $O(e^{-Ct})$ , where the value of  $C$  depends on  $\inf_{t \geq 1} \min_{s,a} \pi_{\theta_t}(a|s) > 0$  and  $\{\theta_t\}_{t=1}^{\infty}$  is generated by Algorithm 1. With a bad initialization  $\theta_1$ ,  $\min_{s,a} \pi_{\theta_1}(a|s)$  could be very small and result in a slow convergence rate. When studying the stochastic PG, this issue of bad initialization will create more severe challenges on the convergence, which we will discuss in the following sections.

One main challenge is the boundedness of iterations under the stochastic PG. The iterates of stochastic gradient methods may indeed escape to infinity in general, rendering the entire scheme of stochastic approximation useless [32, 33]. In particular, when using the stochastic truncated PG for the entropy regularized RL, the key result of Lemma 9 may no longer hold true. This in turn results in the loss of gradient domination condition in guaranteeing the global convergence.

### B. Landscape of a simple bandit example

To have a better understanding of the landscape of the entropy-regularized value function, we visualize its landscape in this section. For the simplicity of the visualization, we use a simple bandit example (corresponding to  $\gamma = 0$ ) with 2 actions, 2 parameters  $(\theta_1, \theta_2)$ , the reward vector  $r = [2, 1]$  and the regularization parameter  $\lambda = 1$ . Then, the entropy-regularized value function can be written as  $\pi_{\theta}^{\top} (r - \log \pi_{\theta})$ .

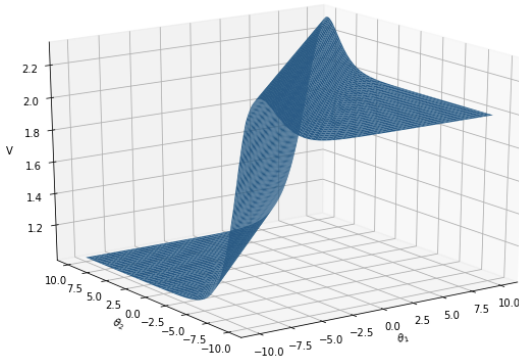


Fig. 1. Landscape of  $\pi_{\theta}^{\top} (r - \log \pi_{\theta})$ .

As shown in Figure 1, the entropy-regularized value function is not coercive. When  $\theta_1$  goes to positive (negative) infinity and  $\theta_2$  goes to negative (positive) infinity, the landscape will become highly flat. It can also be seen that there is a line space for  $(\theta_1, \theta_2)$  at which the entropy-regularized value function is maximum.

When the stochastic PG is used, the search direction may be dominated by the gradient estimation noise at the region where the landscape is highly flat. This may further lead to the

failure of the globally optimal convergence for the stochastic PG algorithm if the initial point is at the flat region.

### C. Convergence to the first-order stationary point

Before presenting our main result, we first show that the stochastic PG proposed in Algorithm 4 asymptotically converges to a region where the PG vanishes almost surely if a specific adaptive step-size sequence is used.

*Lemma 10:* Suppose that the sequence  $\{\theta_t\}_{t=1}^{\infty}$  is generated by Algorithm 4 for the entropy regularized objective with the step-sizes satisfying  $\sum_{t=1}^{\infty} \eta_t = \infty$ ,  $\sum_{t=1}^{\infty} \eta_t^2 < \infty$  and  $\eta_t \leq \frac{2}{L}$  for all  $t = 1, 2, \dots$ . It holds that  $\lim_{t \rightarrow \infty} \|\nabla V_{\lambda}^{\theta_t}(\rho)\|_2 = 0$  with probability 1.

This result follows from classic results for the Robbins-Monro algorithm [34, 35, 36] when an unbiased PG estimator with the bounded variance, as in Algorithm 4, is used in the update rule. No requirement on the batch size  $B$  is needed in Lemma 10. We now provide the proof of Lemma 10 below.

*Proof.*

To prove Lemma 10, it suffices to check the conditions in Proposition 3 of [34] for the objective function  $V_{\lambda}^{\theta}(\rho)$  and the update rule  $\theta_{t+1} = \theta_t + \eta_t(u_t + w_t)$ , where  $u_t = \nabla V_{\lambda}^{\theta_t}(\rho)$  and  $w_t = \hat{\nabla} V_{\lambda}^{\theta_t}(\rho) - \nabla V_{\lambda}^{\theta_t}(\rho)$ .

- 1) From Lemma 3, we know that Condition 1 in Proposition 3 of [34] is satisfied with  $L = \frac{8\bar{r} + \lambda(4 + 8 \log |\mathcal{A}|)}{(1-\gamma)^3}$ .
- 2) Condition 1 in Proposition 3 of [34] is satisfied by the definition of  $\theta_t$  and  $\nabla V_{\lambda}^{\theta}(\rho)$ .
- 3) Condition 1 in Proposition 3 of [34] is satisfied with  $c_1 = 1$  and  $c_2 = 1$ .
- 4) From Lemma 4 and 5, we know that Condition 4 in Proposition 3 of [34] is satisfied with  $A = \frac{8}{(1-\gamma)^2} \left( \frac{\bar{r}^2 + (\lambda \log |\mathcal{A}|)^2}{(1-\gamma^{1/2})^2} \right)$ .
- 5) Condition 1 in Proposition 3 of [34] is satisfied by the definition of  $\eta_t$ .

In addition, it results from Lemma 1 that the entropy-regularized value function  $V_{\lambda}^{\theta}(\rho)$  is bounded. Thus, by Proposition 3 of [34], we must have  $\lim_{t \rightarrow \infty} \nabla V_{\lambda}^{\theta_t}(\rho) = 0$  with probability 1. This completes the proof.  $\square$

However, since the entropy-regularized value function  $V_{\lambda}^{\theta}(\rho)$  is not coercive in  $\theta$  and it may be the case that the gradient  $\nabla V_{\lambda}^{\theta_t}(\rho)$  diminishing to 0 corresponds to  $\theta_t$  going to infinity instead of converging to a stationary point. In addition, the existing results [32, 35, 36] on the almost surely stationary point convergence rely on the assumption that the trajectories of the process are bounded, i.e.,  $\sup_{t \geq 0} \|\theta_t\| < \infty$ , almost surely. This assumption is proven to hold when the function is coercive [37]. However, when the function is not coercive, as in our problem, it is very challenging to characterize the trade-off between the gradient information and the estimation error without additional assumptions.

## V. MAIN RESULT

To overcome the non-coercive landscape challenge, we propose a two-phase stochastic PG algorithm (Algorithm 5).

In the first phase, we will use a large batch size to control the estimation error to guarantee that the stochastic PG is informative even in the regime where the landscape is almost flat. After a certain number of iterations, which is a constant with respect to the optimality gap  $\epsilon$ , the iteration will reach a region where the landscape has enough curvature information. Then, in the second phase, a small batch size is enough to guarantee a fast convergence to the optimal policy.

Before presenting the main result, we first introduce some helpful definitions. Let  $D(\theta_t) = V_\lambda^{\theta_t^*}(\rho) - V_\lambda^{\theta^*}(\rho)$  denote the sub-optimality gap. Since the optimal policy of (1) is unique [11], there must exist a continuum of optimal solutions

$$\Theta^* := \{\theta^* \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} : \frac{\exp(\theta_{s,a}^*)}{\sum_{a'} \exp(\theta_{s,a'}^*)} = \pi_\lambda^*(a | s), \forall s \in \mathcal{S}, a \in \mathcal{A}\}.$$

In addition, we use  $\pi_{\theta^*}$  and  $\pi_\lambda^*$  interchangeably to denote the optimal policy of the entropy-regularized RL. Let  $\{\bar{\theta}_t\}_{t=1}^T$  denote the iterates of the algorithm with the exact PG (Algorithm 1) with  $\eta_t = \eta \leq \frac{1}{2L}$  starting from the initial point  $\theta_1$ . For the soft-max parameterization, we have  $\theta_{s,a} = \log \pi_{\theta}(a | s) + C_s$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , where  $\{C_s\}_{s=1}^{|\mathcal{S}|}$  are some constants. Then, we have

$$\min_{\theta^* \in \Theta^*} \|\bar{\theta}_t - \theta^*\|_2 = \|\log \pi_{\bar{\theta}_t} - \log \pi_\lambda^*\|_2, \quad \text{for all } t = \{1, 2, \dots\}.$$

Furthermore, by Lemma 9, we can define  $\bar{\Delta} := \|\log c_{\bar{\theta}_1, \eta} - \log \pi_\lambda^*\|_2$ , where  $c_{\bar{\theta}_1, \eta} > 0$  is defined in Remark 2. Note that  $\bar{\Delta}$  is only dependent on  $\theta_1$  and  $\eta$  (apart from problem dependent constants), and  $\|\log \pi_{\bar{\theta}_t} - \log \pi_\lambda^*\|_2 \leq \bar{\Delta}$  for any  $\theta_1$  and  $\eta \leq \frac{1}{2L}$ . In addition, with a fair degree of hindsight and for some  $\delta > 0$ , we define the stopping time for the iterates  $\{\theta_t\}_{t=1}^T$  as

$$\tau := \min \left\{ t \mid \min_{\theta^* \in \Theta^*} \|\theta_t - \theta^*\|_2 > \left(1 + \frac{1}{\delta}\right) \bar{\Delta} \right\}, \quad (6)$$

which is the index of the first iterate that exits the bounded region

$$\mathcal{G}_\delta^0 := \left\{ \theta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} : \min_{\theta^* \in \Theta^*} \|\theta - \theta^*\|_2 \leq \left(1 + \frac{1}{\delta}\right) \bar{\Delta} \right\}.$$

Finally, we define  $d(\theta_t) = \min_{\theta^* \in \Theta^*} \|\theta_t - \theta^*\|_2$ . We are now ready to present the main result.

*Theorem 1:* Consider an arbitrary tolerance level  $\delta > 0$  and a small enough tolerance level  $\epsilon > 0$ . For every initial point  $\theta_1$ , if  $\theta_{T+1}$  is generated by Algorithm 5 with

$$\begin{aligned} T_1 &\geq \left( \frac{6D(\theta_1)}{\delta\epsilon_0} \right)^{\frac{8L}{C_\delta^0 \ln 2}}, \quad T_2 \geq \frac{t_0 \epsilon_0}{6\delta\epsilon} - t_0, \quad T = T_1 + T_2, \\ B_1 &\geq \max \left\{ \frac{30\sigma^2}{C_\delta^0 \epsilon_0 \delta}, \frac{6\sigma T_1 \log T_1}{\bar{\Delta} L} \right\}, \quad B_2 \geq \frac{\sigma^2 \ln(T_2 + t_0)}{6C_\alpha \delta \epsilon}, \\ \eta_t &= \eta \leq \min \left\{ \frac{\log T_1}{T_1 L}, \frac{8}{C_\delta^0} \right\} \text{ for } t \leq T_1, \\ \eta_t &= \frac{1}{t - T_1 + t_0} \text{ for } t > T_1 \end{aligned}$$

where

$$\epsilon_0 = \min \left\{ \left( \frac{\lambda \min_s \rho(s)}{6 \ln 2} \right)^2 \left( \alpha \exp\left(\frac{-\bar{r}}{(1-\gamma)\lambda}\right) \right)^4, 1 \right\}, \quad (7)$$

$$t_0 \geq \sqrt{\frac{3\sigma^2}{2\delta\epsilon_0}}, \quad (8)$$

$$C_\alpha := \frac{2\lambda}{|\mathcal{S}|} \min_s \rho(s) (1-\alpha)^2 \min_{s,a} \pi_{\theta^*}(a|s)^2 \left\| \frac{d_\rho^{\pi_\lambda^*}}{\rho} \right\|_\infty^{-1} > 0, \quad (9)$$

$$C_\delta^0 = \frac{2\lambda}{|\mathcal{S}|} \left\| \frac{d_\rho^{\pi_\lambda^*}}{\rho} \right\|_\infty^{-1} \min_s \rho(s) \min_{\theta \in \mathcal{G}_\delta^0} \min_{s,a} \pi_\theta(a|s)^2, \quad (10)$$

and  $\sigma$  is defined in Lemma 5, then we have  $\mathbb{P}(D(\theta_{T+1}) \leq \epsilon) \geq 1 - \delta$ . In total, it requires  $\tilde{\mathcal{O}}(\epsilon^{-2})$  samples to obtain an  $\epsilon$ -optimal policy with high probability.

---

**Algorithm 5** Two-phase stochastic PG for entropy regularized RL

---

- 1: **Inputs:**  $\rho, \lambda, \theta_1, B_1, B_2, T_1, T, \{\eta_t\}_{t=1}^T$ .
  - 2: **for**  $t = 1, 2, \dots, T$  **do**
  - 3:   **if**  $t \leq T_1$  **then**
  - 4:      $B = B_1$
  - 5:   **else**
  - 6:      $B = B_2$
  - 7:   **end if**
  - 8:   Run lines 3-7 in Algorithm 4
  - 9: **end for**
  - 10: **Outputs:**  $\theta_T$ ;
- 

## A. Discussion

In Theorem 1, we have derived strong last-iterate complexity bounds (in contrast to the predominant running-min and ergodic complexity bounds in the reinforcement learning literature), with the desirable  $\tilde{\mathcal{O}}(1/\epsilon^2)$  dependency on the targeting tolerance  $\epsilon$ . That being said, the polynomial dependency on  $1/\delta$  and exponential dependencies on other problem- and algorithm-dependent constants also indicate that our bounds may not be tight in general.

The convergence analysis of the stochastic softmax PG with the entropy regularization is challenging<sup>3</sup> due to the weaker regularization effect of the entropy regularization (compared to the log-barrier regularization adopted in previous works on global optimality convergence of policy gradient methods [16, 38]), as well as the ‘‘softmax gravity well’’ induced by the softmax parameterization which has also been observed in the exact gradient setting [9, 39]. In particular, it only entails uniform gradient domination properties for policies that are bounded below uniformly (*cf.* Lemma 8). We thus need to control the trajectory to ensure that  $\pi_{\theta_t}$  remains in the region where it is uniformly bounded from below for all  $t$ . However, even with large batches, it is generally

<sup>3</sup>Note that similar difficulties in generalization from exact policy gradients to stochastic policy gradients have been observed in [13], which states that ‘‘unlike the exact gradient setting, geometric information cannot be easily exploited in the stochastic case for accelerating policy optimization without detrimental consequences or impractical assumptions’’.

difficult to control stochastic trajectories, which eventually leads to the polynomial dependency on  $1/\delta$  and the exponential dependencies on some constants. If large batches are not used, then the trajectories would be even harder to control and no guarantees may be attained unless additional structural assumptions are enforced on the underlying MDP. In addition, because of the different batch sizes and analysis techniques used in two phases, the conditions for the step-size  $\eta_t$  are also specific for the corresponding phase.

In the next three sections, we provide the proof of Theorem 1. We begin by showing that the iterates will converge to a neighborhood of the optimal solution with high probability in Section VII, and then utilize the curvature information around the optimal policy to guarantee that the action probabilities will still remain uniformly bounded with high probability in Section VIII. We then combine the two steps to prove Theorem 1 in Section IX. For a roadmap of the main ideas behind the proof and the utilization of different lemmas in the paper, please refer to Figure 2.

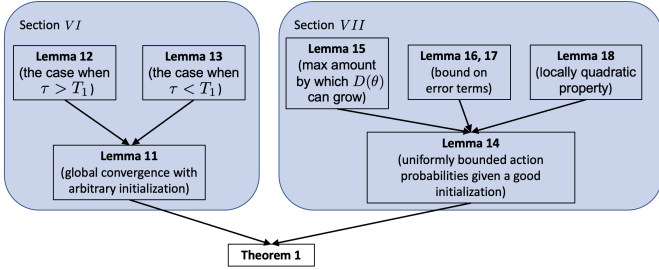


Fig. 2. Roadmap of the main ideas behind the proof of Theorem 1 and their connection to various lemmas in the paper.

## VI. EXPERIMENT

We compare our proposed PG estimator with the estimator given in [8] for PG regularized by KL-divergence between the current policy and the reference policy.

The two state-value estimators are evaluated in a cartpole environment of the Mujoco package. Experiments are performed to compare the performance of the estimators for two different batch sizes  $B \in \{8, 16\}$  when  $\lambda = 0.1$  (see Fig. 3). For each batch size, experiments are repeated with five different seeds. For each subfigure, the solid lines are the means of the experiments among five different seeds, and the shaded area is a confidence interval within one standard deviation.

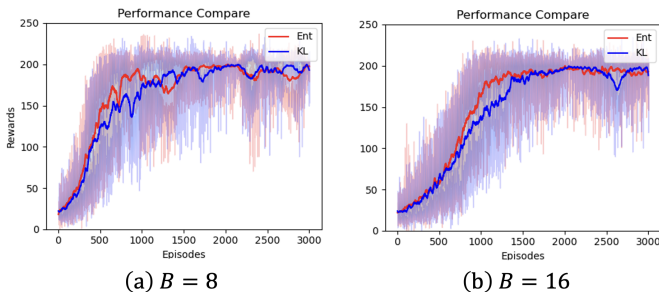


Fig. 3. Rewards comparison among two different value estimators.

The red line is our proposed method, and the blue line is the KL-divergence-based estimator given in [8]. Fig. 3 shows that our proposed estimator converges to the rewards faster than the estimator of [8], which supports the results of this paper that the proposed state value estimator can better evaluate the policy than the KL-divergence-based estimator.

## VII. GLOBAL CONVERGENCE WITH ARBITRARY INITIALIZATION

In this section, we provide the first step towards the proof of Theorem 1. In particular, we will prove that after the first phase of Algorithm 5, the iterates will converge to a neighborhood of the optimal solution with high probability due to the use of a large batch size.

With a large batch size, we can show that if the iterations with the exact PG are bounded, then the iterations with the unbiased stochastic PG will remain bounded with high probability. This will further imply that the unbiased stochastic PG will converge to the neighborhood of the globally optimal policy with high probability. This is a non-trivial result involving the stopping/hitting time analysis, as presented below.

*Lemma 11:* Consider arbitrary tolerance levels  $\delta > 0$  and  $\epsilon_0 > 0$ . For every initial point  $\theta_1$ , if  $\theta_T$  is generated by Algorithm 4 with  $\eta_t = \eta \leq \min\left\{\frac{\log T_1}{T_1 L}, \frac{8}{C_\delta^0}\right\}$ ,  $T_1 = \left(\frac{6D(\theta_1)}{\delta \epsilon_0}\right)^{\frac{8L}{C_\delta^0 \ln 2}}$ , and  $B_1 = \max\left\{\frac{30\sigma^2}{C_\delta^0 \epsilon_0 \delta}, \frac{6\sigma}{\Delta L} \cdot T_1 \cdot \log T_1\right\}$ , then we have  $\mathbb{P}(D(\theta_{T_1}) \leq \epsilon_0) \geq 1 - \delta/2$ .

### A. Helpful lemmas

To prove Lemma 11, we consider the case when  $\tau > T_1$  and the case when  $\tau \leq T_1$  separately, where  $\tau$  is defined in (6). When  $\tau > T_1$ , we can use Lemma 8 to show that  $D(\theta_t)$  is linearly convergent up to some aggregated estimation error.

*Lemma 12:* If  $\eta_t = \eta \leq \min\left\{\frac{1}{2L}, \frac{8}{C_\delta^0}\right\}$ , then  $\mathbb{E}[D(\theta_{T_1}) \mathbf{1}_{\tau > T_1}] \leq \left(1 - \frac{\eta C_\delta^0}{8}\right)^{T_1 - 1} D(\theta_1) + \frac{5\sigma^2}{8C_\delta^0 B_1}$ .

*Proof.* Let  $e_t = \nabla V_\lambda^{\theta_t}(\rho) - u_t$ , where  $u_t = \frac{1}{B_1} \sum_{i=1}^{B_1} \hat{\nabla} V_\lambda^{\theta_t, i}(\rho)$  and  $\hat{\nabla} V_\lambda^{\theta_t, i}(\rho)$  is an unbiased estimator of  $\nabla V_\lambda^{\theta_t}(\rho)$ . Since  $\nabla V_\lambda^\theta(\rho)$  is  $L$ -smooth due to Lemma 3, it follows from Lemma 19 in the supplementary material:

$$\begin{aligned} & \mathbb{E}^t [D(\theta_{t+1}) - D(\theta_t)] \mathbf{1}_{\tau > t} \\ &= \mathbb{E}^t [V_\lambda^{\theta_t}(\rho) - V_\lambda^{\theta_{t+1}}(\rho)] \mathbf{1}_{\tau > t} \\ &\leq \mathbb{E}^t \left[ -\frac{\eta}{8} \|u_t\|_2^2 + \frac{3\eta}{4} \|e_t\|_2^2 \right] \mathbf{1}_{\tau > t} \\ &\leq \mathbb{E}^t \left[ -\frac{\eta}{8} \|u_t - \nabla V_\lambda^{\theta_t}(\rho) + \nabla V_\lambda^{\theta_t}(\rho)\|_2^2 + \frac{3\eta}{4} \|e_t\|_2^2 \right] \mathbf{1}_{\tau > t} \\ &= \mathbb{E}^t \left[ -\frac{\eta}{8} \|\nabla V_\lambda^{\theta_t}(\rho)\|_2^2 + \frac{5\eta}{8} \|e_t\|_2^2 \right] \mathbf{1}_{\tau > t} \\ &\leq \mathbb{E}^t \left[ -\frac{\eta C(\theta_t)}{8} D(\theta_t) + \frac{5\eta}{8} \|e_t\|_2^2 \right] \mathbf{1}_{\tau > t}, \end{aligned}$$

for every  $\eta \leq \frac{1}{2L}$ , where the second inequality uses the fact that  $u_t$  is an unbiased estimator of  $\nabla V_\lambda^{\theta_t}(\rho)$  and the last inequality is due to Lemma 8. We now consider two cases:



- Case 1: Assume that  $\tau > t$ , which implies that  $\theta_t \in \mathcal{G}_\delta^0$  and  $C(\theta_t) \geq C_\delta^0$ . Then, we have  $\mathbb{E}[D(\theta_{t+1})|\mathcal{F}_t] \leq \left(1 - \frac{\eta C_\delta^0}{8}\right) D(\theta_t) + \frac{5\eta}{8} \mathbb{E}[\|e_t\|_2^2|\mathcal{F}_t]$ .
- Case 2: Assume that  $\tau \leq t$  which leads to  $\mathbb{E}[D(\theta_{t+1})|\mathcal{F}_t]\mathbf{1}_{\tau>t} = 0$ .

Now combining the above two cases yields the inequality

$$\begin{aligned} & \mathbb{E}[D(\theta_{t+1})|\mathcal{F}_t]\mathbf{1}_{\tau>t} \\ & \leq \left\{ \left(1 - \frac{\eta C_\delta^0}{8}\right) D(\theta_t) + \frac{5\eta}{8} \mathbb{E}[\|e_t\|_2^2|\mathcal{F}_t] \right\} \mathbf{1}_{\tau>t} \\ & \leq \left(1 - \frac{\eta C_\delta^0}{8}\right) D(\theta_t)\mathbf{1}_{\tau>t} + \frac{5\eta}{8} \mathbb{E}[\|e_t\|_2^2|\mathcal{F}_t]. \end{aligned}$$

In addition, conditioning on  $\mathcal{F}_t$  yields that

$$\begin{aligned} \mathbb{E}[D(\theta_{t+1})\mathbf{1}_{\tau>t+1}|\mathcal{F}_t] & \leq \mathbb{E}[D(\theta_{t+1})\mathbf{1}_{\tau>t}|\mathcal{F}_t] \\ & = \mathbb{E}[D(\theta_{t+1})|\mathcal{F}_t]\mathbf{1}_{\tau>t}, \end{aligned}$$

where the last equality uses the fact that  $\tau$  is a stopping time and the random variable  $\mathbf{1}_{\tau>t}$  is determined completely by the sigma-field  $\mathcal{F}_t$ . Taking the expectations over the sigma-field  $\mathcal{F}_t$  and then arguing inductively gives rise to

$$\begin{aligned} & \mathbb{E}[D(\theta_{t+1})\mathbf{1}_{\tau>t+1}] \\ & \leq \prod_{i=0}^t \left(1 - \frac{\eta C_\delta^0}{8}\right) D(\theta_1) + \sum_{i=0}^t \left(1 - \frac{\eta C_\delta^0}{8}\right)^i \frac{5\eta}{8} \mathbb{E}[\|e_i\|_2^2] \\ & \leq \left(1 - \frac{\eta C_\delta^0}{8}\right)^t D(\theta_1) + \frac{5\sigma^2}{C_\delta^0 B_1}. \end{aligned}$$

By setting  $t+1 = T_1$ , we obtain that  $\mathbb{E}[D(\theta_{T_1})\mathbf{1}_{\tau>T_1}] \leq \left(1 - \frac{\eta C_\delta^0}{8}\right)^{T_1-1} D(\theta_1) + \frac{5\sigma^2}{C_\delta^0 B_1}$ . This completes the proof.  $\square$

We now establish that  $\{\theta_t\}_{t=1}^T$  will be bounded with high probability if the large batch size is used.

*Lemma 13:* It holds that  $\mathbb{P}(\tau \leq T_1) \leq \frac{\delta \eta T_1 \cdot (1+\eta L)^{T_1-1} \cdot \sigma}{\Delta B_1}$ .

*Proof.* By the triangle inequality and the fact that the iterations of the algorithm with the exact PG are bounded by  $\bar{\Delta}$ , we have

$$d(\theta_t) \leq \|\theta_t - \bar{\theta}_t\|_2 + \min_{\theta^* \in \Theta^*} \|\theta^* - \bar{\theta}_t\|_2 \leq \|\theta_t - \bar{\theta}_t\|_2 + \bar{\Delta}.$$

Using the update rule of the algorithm with the exact PG  $\nabla V_\lambda^{\bar{\theta}_i}(\rho)$  and the stochastic PG  $u_i = \frac{1}{B_1} \sum_{j=1}^{B_1} \hat{\nabla} V_\lambda^{\theta_i, j}(\rho)$ , one can write

$$\begin{aligned} d(\theta_i) & = \left\| \left( \theta_1 + \sum_{i=1}^{t-1} \eta_i u_i \right) - \left( \theta_1 + \sum_{i=1}^{t-1} \eta_i \nabla V_\lambda^{\bar{\theta}_i}(\rho) \right) \right\|_2 + \bar{\Delta} \\ & \leq \sum_{i=1}^{t-1} \eta_i \left\| u_i - \nabla V_\lambda^{\bar{\theta}_i}(\rho) \right\|_2 + \bar{\Delta} \\ & = \sum_{i=1}^{t-1} \eta_i \left\| u_i - \nabla V_\lambda^{\theta_i}(\rho) + \nabla V_\lambda^{\theta_i}(\rho) - \nabla V_\lambda^{\bar{\theta}_i}(\rho) \right\|_2 + \bar{\Delta} \\ & \leq \sum_{i=1}^{t-1} \eta_i \|e_i\|_2 + \sum_{i=1}^{t-1} \eta_i L \|\theta_i - \bar{\theta}_i\|_2 + \bar{\Delta}. \end{aligned}$$

By expanding  $\|\theta_i - \bar{\theta}_i\|_2$  recursively, it can be concluded that

$$\begin{aligned} & d(\theta_t) \\ & \leq \sum_{i=1}^{t-1} \eta_i \|e_i\|_2 + \eta_{t-1} L \|\theta_{t-1} - \bar{\theta}_{t-1}\|_2 + \sum_{i=1}^{t-2} \eta_i L \|\theta_i - \bar{\theta}_i\|_2 + \bar{\Delta} \\ & \leq \sum_{i=1}^{t-1} \eta_i \|e_i\|_2 + \eta_{t-1} L \sum_{i=1}^{t-2} \eta_i \|e_i\|_2 + \eta_{t-1} L^2 \sum_{i=1}^{t-2} \eta_i \|\theta_i - \bar{\theta}_i\|_2 \\ & \quad + \sum_{i=1}^{t-2} \eta_i L \|\theta_i - \bar{\theta}_i\|_2 + \bar{\Delta} \\ & = \sum_{i=1}^{t-1} \eta_i \|e_i\|_2 + \eta_{t-1} L \sum_{i=1}^{t-2} \eta_i \|e_i\|_2 \\ & \quad + \sum_{i=1}^{t-2} (\eta_i L + \eta_{t-1} \eta_i L^2) \|\theta_i - \bar{\theta}_i\|_2 + \bar{\Delta} \\ & \leq \sum_{i=1}^{t-1} \eta_i \|e_i\|_2 + \eta_{t-1} L \sum_{i=1}^{t-2} \eta_i \|e_i\|_2 + \bar{\Delta} \\ & \quad + (\eta_{t-2} L + \eta_{t-1} \eta_{t-2} L^2) \sum_{i=1}^{t-3} \eta_i \|e_i\|_2 \\ & \quad + \sum_{i=1}^{t-3} ((\eta_{t-2} L + \eta_{t-1} \eta_{t-2} L^2) \eta_i L \\ & \quad \quad \quad + (\eta_i L + \eta_{t-1} \eta_i L^2)) \|\theta_i - \bar{\theta}_i\|_2 \\ & \leq \bar{\Delta} + \sum_{i=1}^{t-1} \eta_i \|e_i\|_2 + \eta_{t-1} L \sum_{i=1}^{t-2} \eta_i \|e_i\|_2 \\ & \quad + (\eta_{t-2} L + \eta_{t-1} \eta_{t-2} L^2) \sum_{i=1}^{t-3} \eta_i \|e_i\|_2 \\ & \quad + \sum_{i=1}^{t-4} ((\eta_{t-2} L + \eta_{t-1} \eta_{t-2} L^2) \eta_{t-3} L \\ & \quad \quad \quad + (\eta_{t-3} L + \eta_{t-1} \eta_{t-3} L^2)) \|e_i\|_2 \\ & \quad + \sum_{i=1}^{t-4} ((\eta_{t-2} \eta_{t-3} L^2 + \eta_{t-1} \eta_{t-2} \eta_{t-3} L^3) \eta_i L \\ & \quad \quad \quad + (\eta_{t-3} L + \eta_{t-1} \eta_{t-3} L^2) \eta_i) \|\theta_i - \bar{\theta}_i\|_2 \\ & = \sum_{i=1}^{t-1} \eta_i \prod_{j=i+1}^{t-1} (1 + \eta_j L) \|e_i\|_2 + \bar{\Delta}. \end{aligned}$$

Then, by the definition of  $\tau$  in (6) and Markov inequality, we obtain

$$\begin{aligned} \mathbb{P}(\tau \leq T_1) & = \mathbb{P}\left( \max_{t \in \{1, \dots, T_1\}} d(\theta_t) \geq \left(1 + \frac{1}{\delta}\right) \bar{\Delta} \right) \\ & \leq \mathbb{P}\left( \sum_{i=1}^{T_1-1} \eta_i \prod_{j=i+1}^{T_1-1} (1 + \eta_j L) \|e_i\|_2 + \bar{\Delta} \geq \left(1 + \frac{1}{\delta}\right) \bar{\Delta} \right) \\ & \leq \frac{\sum_{i=1}^{T_1-1} \eta_i \prod_{j=i+1}^{T_1-1} (1 + \eta_j L) \mathbb{E}[\|e_i\|_2]}{\frac{1}{\delta} \bar{\Delta}} \\ & \leq \frac{\delta \eta (1 + \eta L)^{T_1-1} \sum_{i=1}^{T_1-1} \mathbb{E}[\|e_i\|_2]}{\bar{\Delta}}, \end{aligned}$$

where we use the fact that  $\eta_t = \eta$  for all  $t \in \{1, 2, \dots\}$ . Furthermore, since  $\mathbb{E}[\|e_i\|_2] \leq \sqrt{\mathbb{E}[\|e_i\|_2^2]} \leq \frac{\sigma}{B_1}$ , we have  $\mathbb{P}(\tau \leq T_1) \leq \frac{\delta \eta T_1 \cdot (1+\eta L)^{T_1-1} \cdot \sigma}{\Delta B_1}$ . This completes the proof.  $\square$

### B. Proof of Lemma 11

By combining Lemmas 12 and 13, we obtain that

$$\begin{aligned}
& \mathbb{P}(D(\theta_{T_1}) \geq \epsilon_0) \\
& \leq \mathbb{P}(\tau > T_1, D(\theta_{T_1}) \geq \epsilon_0) + \mathbb{P}(\tau \leq T_1, D(\theta_{T_1}) \geq \epsilon_0) \\
& \leq \frac{\mathbb{E}[\mathbf{1}_{\tau > T_1} D(\theta_{T_1})]}{\epsilon_0} + \mathbb{P}(\tau \leq T_1) \\
& \leq \left(1 - \frac{\eta C_\delta^0}{8}\right)^{T_1-1} \frac{D(\theta_1)}{\epsilon_0} + \frac{5\sigma^2}{C_\delta^0 B_1 \epsilon_0} \\
& \quad + \frac{\delta \cdot \eta \cdot T_1 \cdot (1 + \eta L)^{T_1-1} \cdot \sigma}{\bar{\Delta} B_1} \\
& \leq \left(1 - \frac{\eta C_\delta^0}{8}\right)^{\frac{8}{\eta C_\delta^0} \frac{\eta C_\delta^0 T_1}{8}} \frac{D(\theta_1)}{\epsilon_0} \\
& \quad + \frac{5\sigma^2}{C_\delta^0 B_1 \epsilon_0} + \frac{\delta \cdot \eta \cdot T_1 \cdot (1 + \eta L)^{T_1-1} \cdot \sigma}{\bar{\Delta} B_1} \\
& \leq \frac{1}{2} \frac{\eta C_\delta^0 T_1}{8} \frac{D(\theta_1)}{\epsilon_0} + \frac{5\sigma^2}{C_\delta^0 B_1 \epsilon_0} + \frac{\delta \cdot \eta \cdot T_1 \cdot (1 + \eta L)^{T_1-1} \cdot \sigma}{\bar{\Delta} B_1},
\end{aligned}$$

where the second inequality holds due to the Markov inequality, and the last inequality holds because of  $(1 - \frac{1}{m})^m \leq \frac{1}{2}$  for all  $m \geq 1$  and  $\frac{8}{\eta C_\delta^0} \geq 1$ . For any  $x \in \mathbb{R}$  that satisfies  $x > 0$ , the inequality  $(\log x)/x - 1/2 < 0$  also holds. Therefore, for any  $T_1 > 0$ , the inequality  $\frac{\log T_1}{T_1 L} - \frac{1}{2L} < 0$  always holds. By taking  $\eta \leq \min\left\{\frac{\log T_1}{T_1 L}, \frac{8}{C_\delta^0}\right\}$ , we obtain

$$\begin{aligned}
& \mathbb{P}(D(\theta_{T_1}) \geq \epsilon_0) \\
& \leq \frac{1}{2} \frac{C_\delta^0 \log T_1}{8L} \frac{D(\theta_1)}{\epsilon_0} + \frac{5\sigma^2}{C_\delta^0 B_1 \epsilon_0} + \frac{\delta \cdot \log T_1 \cdot (1 + \frac{\log T_1}{T_1})^{T_1-1} \cdot \sigma}{\bar{\Delta} B_1 L} \\
& \leq \frac{1}{2} \frac{C_\delta^0 \log T_1}{8L} \frac{D(\theta_1)}{\epsilon_0} + \frac{5\sigma^2}{C_\delta^0 B_1 \epsilon_0} \\
& \quad + \frac{\delta \cdot \log T_1 \cdot (1 + \frac{\log T_1}{T_1})^{\frac{T_1}{\log T_1} \cdot \log T_1} \cdot \sigma}{\bar{\Delta} B_1 L} \\
& \leq \frac{1}{2} \frac{C_\delta^0 \log T_1}{8L} \frac{D(\theta_1)}{\epsilon_0} + \frac{5\sigma^2}{C_\delta^0 B_1 \epsilon_0} + \frac{\delta \cdot \log T_1 \cdot T_1 \cdot \sigma}{\bar{\Delta} B_1 L} \\
& \leq \frac{1}{T_1^{\frac{\ln 2 C_\delta^0}{8L}}} \frac{D(\theta_1)}{\epsilon_0} + \frac{5\sigma^2}{C_\delta^0 B_1 \epsilon_0} + \frac{\delta \cdot \log T_1 \cdot T_1 \cdot \sigma}{\bar{\Delta} B_1 L},
\end{aligned}$$

where we have used  $(1+x)^{1/x} \leq e$  in the third inequality and  $a^{\ln b} = b^{\ln a}$  in the last inequality. To guarantee  $\mathbb{P}(D(\theta_{T_1}) \geq \epsilon_0) \leq \delta/2$ , it suffices to have

$$T_1 = \left(\frac{6D(\theta_1)}{\delta \epsilon_0}\right)^{\frac{8L}{C_\delta^0 \ln 2}}, B_1 = \max\left\{\frac{30\sigma^2}{C_\delta^0 \epsilon_0 \delta}, \frac{6\sigma}{\bar{\Delta} L} \cdot T_1 \cdot \log T_1\right\}.$$

This completes the proof.

### VIII. UNIFORMLY BOUNDED ACTION PROBABILITIES GIVEN A GOOD INITIALIZATION

In this section, we will show how to utilize the curvature information around the optimal policy to guarantee that the action probabilities will still remain uniformly bounded with

high probability, which serves as the second step towards the proof of Theorem 1.

*Lemma 14:* Given a tolerance level  $\delta > 0$ , let  $\pi_\lambda^*$  be the optimal policy of  $V_\lambda^\theta(\rho)$ . Assume further that the random variable  $\{\theta_t\}_{t=1}^{T_2}$  is generated from Algorithm 4 with a step-size sequence of the form  $\eta_t = 1/(t + t_0)$  and a batch-size sequence  $B \geq \frac{1}{\eta_t}$  for all  $t = 1, 2, \dots, T_2$ . If  $t_0 \geq \sqrt{\frac{3\sigma^2}{2\delta\epsilon_0}}$ , and  $\pi_{\theta_1}$  is initialized in a neighborhood  $\mathcal{U}_1$  such that

$$\mathcal{U}_1 = \{\pi \in \Delta(\mathcal{A})^{|\mathcal{S}|} : D(\pi) \leq \epsilon_0\}, \quad (11)$$

where  $\epsilon_0 = \min\left\{\left(\frac{\lambda \min_s \rho(s)}{6 \ln 2}\right)^2 \left(\alpha \exp\left(\frac{-\bar{r}}{(1-\gamma)\lambda}\right)\right)^4, 1\right\}$  and the constant  $\alpha \in (0, 1)$ , then the event <sup>4</sup>

$$\Omega_{\alpha,1}^{T_2} = \left\{\min_{s,a} \pi_{\theta_t}(a|s) \geq (1-\alpha) \min_{s,a} \pi_\lambda^*(a|s), \forall t = 1, 2, \dots, T_2\right\} \quad (12)$$

occurs with probability at least  $1 - \delta/6$ .

#### A. Helpful lemmas

To prove Lemma 14, we first characterize the maximum amount by which  $D(\theta_t)$  can grow at each step.

*Lemma 15:* Suppose that  $\{\theta_t\}$  is generated by Algorithm 4 with  $0 < \eta_t \leq \frac{(1-\gamma)^3}{16\bar{r} + \lambda(8+16\log|\mathcal{A}|)}$  for all  $t \geq 1$ . We have

$$D(\theta_{t+1}) \leq \left(1 - \frac{\eta_t C(\theta_t)}{4}\right) D(\theta_t) - \frac{\eta_t}{2} \xi_t + \frac{\eta_t}{4} \|e_t\|_2^2, \quad (13)$$

where  $\xi_t = \langle e_t, \nabla V_\lambda^{\theta_t}(\rho) \rangle$  and  $e_t = \hat{\nabla} V_\lambda^{\theta_t}(\rho) - \nabla V_\lambda^{\theta_t}(\rho)$ .

*Proof.* Since  $\nabla V_\lambda^\theta(\rho)$  is  $L$ -smooth in light of Lemma 3, it follows from Lemma 19 that

$$\begin{aligned}
& D(\theta_{t+1}) - D(\theta_t) \\
& \leq -\frac{\eta_t}{4} \|\hat{\nabla} V_\lambda^\theta(\rho)\|_2^2 + \frac{\eta_t}{2} \|e_t\|_2^2 \\
& \leq -\frac{\eta_t}{4} \|\hat{\nabla} V_\lambda^{\theta_t}(\rho) - \nabla V_\lambda^{\theta_t}(\rho) + \nabla V_\lambda^{\theta_t}(\rho)\|_2^2 + \frac{\eta_t}{2} \|e_t\|_2^2 \\
& = -\frac{\eta_t}{4} \|\hat{\nabla} V_\lambda^{\theta_t}(\rho) - \nabla V_\lambda^{\theta_t}(\rho)\|_2^2 - \frac{\eta_t}{4} \|\hat{\nabla} V_\lambda^{\theta_t}(\rho)\|_2^2 \\
& \quad - \frac{\eta_t}{2} \langle e_t, \nabla V_\lambda^{\theta_t}(\rho) \rangle + \frac{\eta_t}{2} \|e_t\|_2^2 \\
& = -\frac{\eta_t}{4} \|\nabla V_\lambda^{\theta_t}(\rho)\|_2^2 - \frac{\eta_t}{2} \langle e_t, \nabla V_\lambda^{\theta_t}(\rho) \rangle + \frac{\eta_t}{4} \|e_t\|_2^2 \\
& \leq -\frac{\eta_t C(\theta_t)}{4} D(\theta_t) - \frac{\eta_t}{2} \langle e_t, \nabla V_\lambda^{\theta_t}(\rho) \rangle + \frac{\eta_t}{4} \|e_t\|_2^2,
\end{aligned}$$

for every  $\eta_t \leq \frac{1}{2L}$ , where the last inequality is due to Lemma 8.  $\square$

The quantity by which  $D(\theta_t)$  can grow at each step can be large for any given  $t$  but we will show that, with high probability, the aggregation of these errors remains controllably small under the stated conditions on the step-sizes and batch size.

Similar as the techniques used in [23, 37, 41, 42, 43], we now encode the error terms in (13) as  $M_n = \sum_{t=1}^n \eta_t \xi_t$  and  $S_n = \sum_{t=1}^n \frac{\eta_t}{4} \|e_t\|_2^2$ .

<sup>4</sup>The  $\sigma$ -field of this event is the Cartesian product of the natural Borel  $\sigma$ -field on the underlying MDPs [40, Section 2.1.6].

For given  $\pi_{\theta^n}$ , the parameter  $\xi_n$  is fully determined by a trajectory obtained at iteration  $n$ , i.e.,  $\tau_n = \{s_0^n a_0^n, \dots, s_H^n\}$ . For a stationary environment, the trajectory  $\tau_n$  is fully determined by the policy at time  $n$ , namely  $\pi^n$ , which is parameterized by  $\theta_n$  where  $\theta_n$  is fully determined by the update rule by algorithm 4 (**Ent-RPG**). Note that the input of the update rule is fully determined by the information up to the  $n-1$ . Then,  $\theta_n$  is measurable with respect to  $\mathcal{F}_{n-1}$ , which subsequently concludes that  $\{\xi_1, \xi_2, \dots, \xi_n\}$  are also measurable regarding with  $\mathcal{F}_{n-1}$ . Therefore  $\mathbb{E}^{n-1}[\xi_n] = 0$  since

$$\begin{aligned} \mathbb{E}^{n-1}[\xi_n] &= \mathbb{E} \left[ \left\langle \hat{\nabla} V_\lambda^{\theta^n}(\rho) - \nabla V_\lambda^{\theta^n}(\rho), \nabla V_\lambda^{\theta^n}(\rho) \right\rangle \middle| \mathcal{F}_{n-1} \right] \\ &= \mathbb{E} \left[ \left\langle \nabla V_\lambda^{\theta^n}(\rho), \nabla V_\lambda^{\theta^n}(\rho) \right\rangle \middle| \mathcal{F}_{n-1} \right] \\ &\quad - \mathbb{E} \left[ \left\| \nabla V_\lambda^{\theta^n}(\rho) \right\|^2 \middle| \mathcal{F}_{n-1} \right] \\ &= 0 \end{aligned}$$

holds. Then, we have  $\mathbb{E}^{n-1}[M_n] = M_{n-1}$ . Therefore,  $M_n$  is a zero-mean martingale; likewise,  $\mathbb{E}^{n-1}[S_n] \geq S_{n-1}$ , and therefore,  $S_n$  is a submartingale. The difficulty of controlling the errors in  $M_n$  and  $S_n$  lies in the fact that the estimation error  $e_n$  may be unbounded. Because of this, we need to take a less direct, step-by-step approach to bound the total error increments conditioned on the event that  $D(\theta_n)$  remains close to  $D(\theta^*)$ . We begin by introducing the ‘‘cumulative mean square error’’  $R_n = M_n^2 + S_n$ . By construction, we have

$$\begin{aligned} R_n &= (M_{n-1} + \eta_n \xi_n)^2 + S_{n-1} + \frac{1}{4} \eta_n \|e_n\|^2 \\ &= R_{n-1} + 2M_{n-1} \eta_n \xi_n + \eta_n^2 \xi_n^2 + \frac{1}{4} \eta_n \|e_n\|^2. \end{aligned}$$

Hence,  $\mathbb{E}^{n-1}[R_n] = R_{n-1} + 2M_{n-1} \eta_n \mathbb{E}^{n-1}[\xi_n] + \eta_n^2 \mathbb{E}^{n-1}[\xi_n^2] + \frac{1}{4} \eta_n \mathbb{E}^{n-1}[\|e_n\|^2] \geq R_{n-1}$ , i.e.,  $R_n$  is a submartingale. With a fair degree of hindsight, we define  $\mathcal{U}$  as:

$$\mathcal{U} = \left\{ \pi \in \Delta(\mathcal{A})^{|\mathcal{S}|} : D(\pi) \leq 2\epsilon_0 + \sqrt{\epsilon_0} \right\}. \quad (14)$$

To condition it further, we also define the events

$$\begin{aligned} \Omega_n &\equiv \Omega_n(\epsilon_0) = \{ \pi_{\theta_t} \in \mathcal{U} \text{ for all } t = 1, 2, \dots, n \} \\ E_n &\equiv E_n(\epsilon_0) = \{ R_t \leq \epsilon_0 \text{ for all } t = 1, 2, \dots, n \} \end{aligned}$$

By definition, we also have  $\Omega_0 = E_0 = \Omega$  (because the set-building index set for  $k$  is empty in this case, and every statement is true for the elements of the empty set). These events will play a crucial role in the sequel as indicators of whether  $\pi_{\theta_t}$  has escaped the vicinity of  $\pi_\lambda^*$ .

For brevity, we write  $\mathcal{F}_n = \sigma(\theta_1, \dots, \theta_n)$  for the natural filtration of  $\theta_n$ . Now, we are ready to state the next lemma.

*Lemma 16:* Let  $\pi_\lambda^*$  be the optimal policy. Then, for all  $n \in \{1, 2, \dots\}$ , the following statements hold:

- 1)  $\Omega_{n+1} \subseteq \Omega_n$  and  $E_{n+1} \subseteq E_n$ .
- 2)  $E_{n-1} \subseteq \Omega_n$ .
- 3) Consider the ‘‘large noise’’ event

$$\begin{aligned} \tilde{E}_n &\equiv E_{n-1} \setminus E_n = E_{n-1} \cap \{ R_n > \epsilon_0 \} \\ &= \{ R_t \leq \epsilon_0 \text{ for all } t = 1, 2, \dots, n-1 \text{ and } R_n > \epsilon_0 \} \end{aligned}$$

and let  $\tilde{R}_n = R_n \mathbb{1}_{E_{n-1}}$  denote the cumulative error subject to the noise being ‘‘small’’ until time  $n$ . Then,

$$\mathbb{E}[\tilde{R}_n] \leq \mathbb{E}[\tilde{R}_{n-1}] + G^2 \sigma^2 \eta_n^2 + \frac{\eta_n \sigma^2}{4B} - \epsilon_0 \mathbb{P}(\tilde{E}_{n-1}). \quad (15)$$

By convention, we write  $\tilde{E}_0 = \emptyset$  and  $\tilde{R}_0 = 0$ .

*Proof.* Statement 1 is obviously true. For Statement 2, we proceed inductively:

- For the base case  $n = 1$ , we have  $\Omega_1 = \{ \pi_{\theta_1} \in \mathcal{U} \} \supseteq \{ \pi_{\theta_1} \in \mathcal{U}_1 \} = \Omega$  because  $\pi_{\theta_1}$  is initialized in  $\mathcal{U}_1 \subseteq \mathcal{U}$ . Since  $E_0 = \Omega$ , our claim follows.
- For the inductive step, assume that  $E_{n-1} \subseteq \Omega_n$  for some  $n \geq 1$ . To show that  $E_n \subseteq \Omega_{n+1}$ , we fix a realization in  $E_n$  such that  $R_t \leq \epsilon$  for all  $t = 1, 2, \dots, n$ . Since  $E_n \subseteq E_{n-1}$ , the inductive hypothesis posits that  $\Omega_n$  also occurs, i.e.,  $\pi_{\theta_t} \in \mathcal{U}$  for all  $t = 1, 2, \dots, n$ ; hence, it suffices to show that  $\pi_{\theta_{n+1}} \in \mathcal{U}$ . To that end, given that  $\pi_{\theta_t} \in \mathcal{U}$  for all  $t = 1, 2, \dots, n$ , the distance estimate (13) readily gives  $D(\theta_{t+1}) \leq D(\theta_t) + \eta_t \xi_t + \frac{\eta_t}{4} \|e_t\|_2^2$  for all  $t = 1, 2, \dots, n$ . Therefore, after telescoping, we obtain

$$\begin{aligned} D(\theta_{n+1}) &\leq D(\theta_1) + M_n + S_n \leq D(\theta_1) + \sqrt{R_n} + R_n \\ &\leq \epsilon + \sqrt{\epsilon} + \epsilon \\ &= 2\epsilon + \sqrt{\epsilon} \end{aligned}$$

by the inductive hypothesis. This completes the induction.

For Statement 3, we decompose  $\tilde{R}_n$  as

$$\begin{aligned} \tilde{R}_n &= R_n \mathbb{1}_{E_{n-1}} \\ &= R_{n-1} \mathbb{1}_{E_{n-1}} + (R_n - R_{n-1}) \mathbb{1}_{E_{n-1}} \\ &= R_{n-1} \mathbb{1}_{E_{n-2}} - R_{n-1} \mathbb{1}_{\tilde{E}_{n-1}} + (R_n - R_{n-1}) \mathbb{1}_{E_{n-1}} \\ &= \tilde{R}_{n-1} + (R_n - R_{n-1}) \mathbb{1}_{E_{n-1}} - R_{n-1} \mathbb{1}_{\tilde{E}_{n-1}} \end{aligned}$$

where we have used the fact that  $E_{n-1} = E_{n-2} \setminus \tilde{E}_{n-1}$  so  $\mathbb{1}_{E_{n-1}} = \mathbb{1}_{E_{n-2}} - \mathbb{1}_{\tilde{E}_{n-1}}$  (recall that  $E_{n-1} \subseteq E_{n-2}$ ). Then, by the definition of  $R_n$ , we have

$$R_n - R_{n-1} = 2M_{n-1} \eta_n \xi_n + \eta_n^2 \xi_n^2 + \frac{1}{4} \eta_n \|e_n\|^2$$

and therefore

$$\begin{aligned} \mathbb{E}[(R_n - R_{n-1}) \mathbb{1}_{E_{n-1}}] &= \\ 2\eta_n \mathbb{E}[M_{n-1} \xi_n \mathbb{1}_{E_{n-1}}] + \eta_n^2 \mathbb{E}[\xi_n^2 \mathbb{1}_{E_{n-1}}] + \frac{1}{4} \eta_n \mathbb{E}[\|e_n\|^2 \mathbb{1}_{E_{n-1}}]. \end{aligned} \quad (16)$$

However, since  $E_{n-1}$  and  $M_{n-1}$  are both  $\mathcal{F}_n$ -measurable, we have the following estimates:

- For the term in (16), by the unbiasedness of the gradient estimator shown in Lemma 4, we have:  $\mathbb{E}[M_{n-1} \xi_n \mathbb{1}_{E_{n-1}}] = \mathbb{E}[M_{n-1} \mathbb{1}_{E_{n-1}} \mathbb{E}[\xi_n | \mathcal{F}_n]] = 0$ .
- The second term in (16) is where the conditioning on  $E_{n-1}$  plays the most important role. It holds that:

$$\begin{aligned} \mathbb{E}[\xi_n^2 \mathbb{1}_{E_{n-1}}] &= \mathbb{E} \left[ \mathbb{1}_{E_{n-1}} \mathbb{E} \left[ \left\langle e_n, \nabla V_\lambda^{\theta^n}(\rho) \right\rangle^2 \middle| \mathcal{F}_n \right] \right] \\ &\leq \mathbb{E} \left[ \mathbb{1}_{E_{n-1}} \left\| \nabla V_\lambda^{\theta^n}(\rho) \right\|^2 \mathbb{E} \left[ \|e_n\|^2 \middle| \mathcal{F}_n \right] \right] \\ &\leq \mathbb{E} \left[ \mathbb{1}_{\Omega_n} \left\| \nabla V_\lambda^{\theta^n}(\rho) \right\|^2 \mathbb{E} \left[ \|e_n\|^2 \middle| \mathcal{F}_n \right] \right] \\ &\leq G^2 \sigma^2 \end{aligned}$$

where the first inequality is due to the Cauchy-Schwarz inequality, the second inequality follows from  $E_{n-1} \subseteq \Omega_n$  and the last inequality results from Lemmas 2 and 5.

- Finally, for the third term in (16), we have:

$$\frac{\eta_n}{4} \mathbb{E} \left[ \|e_n\|_2^2 \mathbb{1}_{E_{n-1}} \right] \leq \frac{\eta_n \sigma^2}{4B}. \quad (17)$$

Thus, putting together all of the above, we obtain  $\mathbb{E}[(R_n - R_{n-1}) \mathbb{1}_{E_{n-1}}] \leq G^2 \sigma^2 \eta_n^2 + \frac{\eta_n \sigma^2}{4B}$ . Since  $R_{n-1} > \varepsilon$  if  $\tilde{E}_{n-1}$  occurs, we obtain  $\mathbb{E}[R_{n-1} \mathbb{1}_{\tilde{E}_{n-1}}] \geq \varepsilon \mathbb{E}[\mathbb{1}_{\tilde{E}_{n-1}}] = \varepsilon \mathbb{P}(\tilde{E}_{n-1})$ . This completes the proof of Statement 3.  $\square$

With the above results, we can show that the cumulative mean square error  $R_n$  is small with high probability at all times.

*Lemma 17:* Consider an arbitrary tolerance level  $\delta > 0$ . If Algorithm 4 is run with a step-size schedule of the form  $\eta_t = 1/(t + t_0)$  where  $t_0 \geq \sqrt{\frac{3\sigma^2}{2\delta\epsilon_0}}$  and a batch size schedule  $B_t \geq \frac{1}{\eta_t}$ , we have  $\mathbb{P}(E_n) \geq 1 - \delta/6$ , for all  $n = 1, 2, \dots$

*Proof.* We begin by bounding the probability of the ‘‘large noise’’ event  $\tilde{E}_n = E_{n-1} \setminus E_n$  as follows:

$$\begin{aligned} \mathbb{P}(\tilde{E}_n) &= \mathbb{P}(E_{n-1} \setminus E_n) = \mathbb{P}(E_{n-1} \cap \{R_n > \varepsilon\}) \\ &= \mathbb{E}[\mathbb{1}_{E_{n-1}} \times \mathbb{1}_{\{R_n > \varepsilon\}}] \\ &\leq \mathbb{E}[\mathbb{1}_{E_{n-1}} \times (R_n/\varepsilon)] = \mathbb{E}[\tilde{R}_n]/\varepsilon, \end{aligned}$$

which is derived by using the fact that  $R_n \geq 0$  (so  $\mathbb{1}_{\{R_n > \varepsilon\}} \leq R_n/\varepsilon$ ). Now, by summing up (15), we conclude that  $\mathbb{E}[\tilde{R}_n] \leq \mathbb{E}[\tilde{R}_0] + \frac{\sigma^2}{4B} \sum_{t=1}^n \eta_t - \varepsilon \sum_{t=1}^n \mathbb{P}(\tilde{E}_{t-1})$ . Hence, combining the above results, we obtain the estimate

$$\sum_{t=1}^n \mathbb{P}(\tilde{E}_k) \leq \frac{\sigma^2}{4B\epsilon_0} \sum_{t=1}^n \eta_t \leq \frac{\sigma^2}{4\epsilon_0} \sum_{t=1}^n \eta_t^2 \leq \frac{\sigma^2 \Gamma}{4\epsilon_0},$$

where  $\Gamma = \sum_{t=1}^{\infty} \eta_t^2 = \sum_{t=1}^{\infty} (t + t_0)^{-2}$ , and we have used the relations that  $\tilde{R}_0 = 0$  and  $\tilde{E}_0 = \emptyset$  (by convention). By choosing  $t_0 \geq \sqrt{\frac{3\sigma^2}{2\delta\epsilon_0}}$ , we ensure that  $\frac{\sigma^2 \Gamma}{4\epsilon_0} < \delta/6$ ; moreover, since the events  $\tilde{E}_t$  are disjoint for all  $t = 1, 2, \dots$ , we obtain  $\mathbb{P}(\bigcup_{t=1}^n \tilde{E}_t) = \sum_{t=1}^n \mathbb{P}(\tilde{E}_t) \leq \delta/6$ . Hence,  $\mathbb{P}(E_n) = \mathbb{P}(\bigcap_{t=1}^n \tilde{E}_t^c) \geq 1 - \delta/6$  as claimed.  $\square$

Furthermore, we can show that the entropy-regularized value function  $V_\lambda^\theta(\rho)$  is locally quadratic around the optimal policy  $\pi_{\theta^*}$ .

*Lemma 18:* For every policy  $\pi_\theta$ , we have

$$D(\theta) \geq \frac{\lambda \min_s \rho(s)}{2 \ln 2} |\pi_\theta(a|s) - \pi_{\theta^*}(a|s)|^2, \quad \forall s \in \mathcal{S}, a \in \mathcal{A}.$$

*Proof.* It follows from the soft sub-optimality difference lemma

(Lemma 26 in [9]) that

$$\begin{aligned} &V_\lambda^{\theta^*}(\rho) - V_\lambda^\theta(\rho) \\ &= \frac{1}{1-\gamma} \sum_s [d_\rho^{\pi_\theta}(s) \cdot \lambda \cdot D_{\text{KL}}(\pi_\theta(\cdot|s) \|\pi_{\theta^*}(\cdot|s))] \\ &\geq \frac{1}{1-\gamma} \sum_s \left[ d_\rho^{\pi_\theta}(s) \cdot \lambda \cdot \frac{1}{2 \ln 2} \|\pi_\theta(\cdot|s) - \pi_{\theta^*}(\cdot|s)\|_1^2 \right] \\ &\geq \frac{\lambda}{2 \ln 2} \sum_s [\rho(s) \cdot \|\pi_\theta(\cdot|s) - \pi_{\theta^*}(\cdot|s)\|_1^2] \\ &\geq \frac{\lambda}{2 \ln 2} \sum_s [\rho(s) \cdot \|\pi_\theta(\cdot|s) - \pi_{\theta^*}(\cdot|s)\|_2^2] \\ &\geq \frac{\lambda}{2 \ln 2} [\rho(s) \|\pi_\theta(\cdot|s) - \pi_{\theta^*}(\cdot|s)\|_2^2] \quad \forall s \in \mathcal{S} \\ &\geq \frac{\lambda \min_s \rho(s)}{2 \ln 2} |\pi_\theta(a|s) - \pi_{\theta^*}(a|s)|^2, \quad \forall s \in \mathcal{S}, a \in \mathcal{A}, \end{aligned}$$

where the first inequality is due to Theorem 11.6 in [44] stating that

$$D_{\text{KL}}[P(\cdot) | Q(\cdot)] \geq \frac{1}{2 \ln 2} \|P(\cdot) - Q(\cdot)\|_1^2$$

for every two discrete distributions  $P(\cdot)$  and  $Q(\cdot)$ . Moreover, the second inequality is due to  $d_\rho^{\pi_\theta}(s) \geq (1-\gamma)\rho(s)$  and the third inequality is due to the equivalence between  $\ell_1$ -norm and  $\ell_2$ -norm. This completes the proof.  $\square$

## B. Proof of Lemma 14

Since the sequence  $\Omega_n$  is decreasing and  $\Omega_n \supseteq E_{n-1}$  (by the second part of Lemma 16), Lemma 17 yields that  $\mathbb{P}(\Omega_{T_2}) \geq \inf_n \mathbb{P}(\Omega_n) \geq \inf_n \mathbb{P}(E_{n-1}) \geq 1 - \delta/6$  provided that  $t_0$  is chosen large enough.

Now, it remains to show that  $\Omega_{T_2} \subseteq \Omega_{\alpha,1}^{T_2}$ . We fix a realization in  $\Omega_{T_2}$  such that  $D(\theta_t) \leq 2\epsilon_0 + \sqrt{\epsilon_0}$  for all  $t = 1, 2, \dots, T_2$ . By Lemma 18, we have

$$\begin{aligned} &|\pi_{\theta_t}(a|s) - \pi_{\theta^*}(a|s)| \\ &\leq \sqrt{\frac{2D(\theta_t) \ln 2}{\lambda \min_s \rho(s)}} \leq \sqrt{\frac{2(2\epsilon_0 + \sqrt{\epsilon_0}) \ln 2}{\lambda \min_s \rho(s)}} \\ &\leq \sqrt{\frac{6\sqrt{\epsilon_0} \ln 2}{\lambda \min_s \rho(s)}} \leq \alpha \exp\left(\frac{-\bar{r}}{(1-\gamma)\lambda}\right) \leq \alpha \min_{s,a} \pi_{\theta^*}(a|s), \end{aligned}$$

where the second inequality is due to the condition that the event  $\Omega_{T_2}$  occurs, the third inequality is due to  $\epsilon_0 \leq \sqrt{\epsilon_0}$  when  $\epsilon_0 \leq 1$ , the fourth inequality is due to the definition of  $\epsilon_0$ , and the last inequality is due to Theorem 1 in [45] where it holds that  $\log \pi_\lambda^*(a|s) = \frac{1}{\lambda} (Q^{\pi_\lambda^*}(s,a) - V^{\pi_\lambda^*}(s)) \geq \frac{-\bar{r}}{(1-\gamma)\lambda}$ ,  $\forall (s,a) \in \mathcal{S} \times \mathcal{A}$ .

Now, it can be easily verified that  $\pi_{\theta_t}(a|s) \geq \pi_{\theta^*}(a|s) - \alpha \min_{s,a} \pi_{\theta^*}(a|s)$ . For every  $t \in \{1, 2, \dots, T_2\}$ , let  $\bar{s}, \bar{a} = \operatorname{argmin}_{s,a} \pi_{\theta_t}(a|s)$ . One can write

$$\begin{aligned} \min_{s,a} \pi_{\theta_t}(a|s) &= \pi_{\theta_t}(\bar{a}|\bar{s}) \geq \pi_{\theta^*}(\bar{a}|\bar{s}) - \alpha \min_{s,a} \pi_{\theta^*}(a|s) \\ &\geq (1-\alpha) \min_{s,a} \pi_{\theta^*}(\bar{a}|\bar{s}), \end{aligned}$$

where the last inequality is due to  $\pi(a|s) \geq \min_{s,a} \pi(a|s)$  for every  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$ . Thus, we obtain  $\mathbb{P}(\Omega_{\alpha,1}^{T_2}) \geq \mathbb{P}(\Omega_{T_2}) \geq 1 - \delta/6$ . This completes the proof.

## IX. PROOF OF THEOREM 1

From Lemma 11, we conclude that, with a large batch size, the iterations will converge to a neighborhood of the optimal solution with high probability. From Lemma 14, we know that, with a good initialization, the policies will remain in the interior of the probability simplex with high probability. By combining the above two results, we are now ready to prove the sample complexity of the stochastic PG for entropy-regularized RL.

From Lemma 11, we can conclude that  $\mathbb{P}(D(\theta_{T_1}) \leq \epsilon_0) \geq 1 - \delta$  after the first phase. We then establish the algorithm's sample complexity when the initial policy of the second phase satisfies the good initialization condition  $\mathbb{P}(D(\theta_{T_1}) \leq \epsilon_0) \geq 1 - \delta$ . It follows from Lemma 15 that

$$\begin{aligned} & D(\theta_{t+1}) \mathbb{1}_{\Omega_{\alpha, T_1}^t} \\ \leq & \left(1 - \frac{\eta_t C(\theta_t)}{4}\right) D(\theta_t) \mathbb{1}_{\Omega_{\alpha, T_1}^t} - \frac{\eta_t}{2} \xi_t \mathbb{1}_{\Omega_{\alpha, T_1}^t} + \frac{\eta_t}{4} \|e_t\|_2^2 \mathbb{1}_{\Omega_{\alpha, T_1}^t}, \end{aligned}$$

for all  $t \geq T_1$ , where  $\xi_t = \langle e_t, \nabla V_\lambda^{\theta_t}(\rho) \rangle$  and  $\Omega_{\alpha, T_1}^t$  is defined in (12). When the event  $\Omega_{\alpha, T_1}^t$  occurs, we have  $C(\theta_t) \geq C_\alpha$ , where  $C_\alpha$  is defined in (9). By taking the expectation, we have

$$\begin{aligned} & \mathbb{E} \left[ -\frac{\eta_t}{2} \xi_t \mathbb{1}_{\Omega_{\alpha, T_1}^t} + \frac{\eta_t}{4} \|e_t\|_2^2 \mathbb{1}_{\Omega_{\alpha, T_1}^t} \right] \\ = & \mathbb{E} \left[ \mathbb{1}_{\Omega_{\alpha, T_1}^t} \mathbb{E} \left[ -\frac{\eta_t}{2} \xi_t + \frac{\eta_t}{4} \|e_t\|_2^2 \middle| \mathcal{F}_t \right] \right] \\ = & \mathbb{E} \left[ \mathbb{1}_{\Omega_{\alpha, T_1}^t} \mathbb{E} \left[ \frac{\eta_t}{4} \|e_t\|_2^2 \middle| \mathcal{F}_t \right] \right] \leq \frac{\eta_t \sigma^2}{4B}, \end{aligned}$$

where the first equality is because  $\Omega_{\alpha, T_1}^t$  is deterministic conditioning on  $\mathcal{F}_t$ , the second equality is due to the unbiasedness of  $\xi_t$  conditioning on  $\mathcal{F}_t$ , and the first inequality is due to (17). Therefore,  $\mathbb{E}[D(\theta_{t+1}) \mathbb{1}_{\Omega_{\alpha, T_1}^t}] \leq \left(1 - \frac{\eta_t C_\alpha}{4}\right) \mathbb{E}[D(\theta_t) \mathbb{1}_{\Omega_{\alpha, T_1}^t}] + \frac{\eta_t \sigma^2}{4B}$ . Arguing inductively yields that

$$\begin{aligned} & \mathbb{E}[D(\theta_{T+1}) \mathbb{1}_{\Omega_{\alpha, T_1}^T}] \\ \leq & \prod_{i=1}^{T_2} \left(1 - \frac{\eta_{T_1+i} C_\alpha}{4}\right) D(\theta_{T_1}) + \sum_{i=1}^{T_2} \left(1 - \frac{\eta_{T_1+i} C_\alpha}{4}\right)^i \frac{\eta_{T_1+i} \sigma^2}{4B} \\ \leq & \prod_{i=1}^{T_2} \left(1 - \frac{\eta_{T_1+i} C_\alpha}{4}\right) D(\theta_{T_1}) + \sum_{i=1}^{T_2} \frac{\eta_{T_1+i} \sigma^2}{4B}. \end{aligned}$$

By taking  $\eta_{T_1+i} = \frac{4}{C_\alpha(i+t_0)}$ , we obtain that

$$\begin{aligned} \mathbb{E}[D(\theta_{T+1}) \mathbb{1}_{\Omega_{\alpha, T_1}^T}] & \leq \prod_{i=1}^{T_2} \left(\frac{i+t_0-1}{i+t_0}\right) D(\theta_{T_1}) + \frac{\sigma^2}{C_\alpha B} \sum_{i=1}^{T_2} \frac{1}{i+t_0} \\ & \leq \frac{t_0}{T_2+t_0} D(\theta_{T_1}) + \frac{\sigma^2 \ln(T_2+t_0)}{BC_\alpha}. \end{aligned}$$

By the law of total probability and the Markov inequality,

we obtain that

$$\begin{aligned} & \mathbb{P}(D(\theta_{T+1}) \geq \epsilon) \\ = & \mathbb{P}(D(\theta_{T+1}) \geq \epsilon, \Omega_{\alpha, T_1}^T) + \mathbb{P}(D(\theta_{T+1}) \geq \epsilon, (\Omega_{\alpha, T_1}^T)^c) \\ = & \mathbb{P}(D(\theta_{T+1}) \geq \epsilon \mid \Omega_{\alpha, T_1}^T) \mathbb{P}(\Omega_{\alpha, T_1}^T) \\ & + \mathbb{P}(D(\theta_{T+1}) \geq \epsilon \mid (\Omega_{\alpha, T_1}^T)^c) \mathbb{P}((\Omega_{\alpha, T_1}^T)^c) \\ \leq & \frac{\mathbb{E}[D(\theta_{T+1}) \mid \Omega_{\alpha, T_1}^T]}{\epsilon} \mathbb{P}(\Omega_{\alpha, T_1}^T) \\ & + \mathbb{P}(D(\theta_{T+1}) \geq \epsilon \mathbb{1}_{\Omega_{\alpha, T_1}^T}^c) \mathbb{P}((\Omega_{\alpha, T_1}^T)^c) \\ \leq & \frac{\mathbb{E}[D(\theta_{T+1}) \mathbb{1}_{\Omega_{\alpha, T_1}^T}]}{\epsilon} + \delta/6 \\ \leq & \frac{t_0}{(T_2+t_0)\epsilon} D(\theta_{T_1}) + \frac{\sigma^2 \ln(T_2+t_0)}{BC_\alpha \epsilon} + \delta/6, \end{aligned}$$

where the second inequality follows from Lemma 14. To guarantee  $\mathbb{P}(D(\theta_{T+1}) \geq \epsilon) \leq \frac{\delta}{2}$ , it suffices to have  $T_2 = \frac{t_0 D(\theta_{T_1})}{6\delta\epsilon} - t_0$ ,  $B = \frac{\sigma^2 \ln(T_2+t_0)}{6C_\alpha \delta\epsilon}$ . This completes the proof.

## X. CONCLUSION

In this work, we studied the global convergence and the sample complexity of stochastic PG methods for the entropy-regularized RL with the soft-max parameterization. We proposed two new (nearly) unbiased PG estimators for the entropy-regularized RL and proved that they have a bounded variance even though they could be unbounded. In addition, we developed a two-phase stochastic PG algorithm to overcome the non-coercive landscape challenge. This work provided the first global convergence result for stochastic PG methods for the entropy-regularized RL and obtained the sample complexity of  $\tilde{O}(\frac{1}{\epsilon^2})$ , where  $\epsilon$  is the optimality threshold. This work paves the way for a deeper understanding of other stochastic PG methods with entropy-related regularization, including those with trajectory-level KL regularization and policy reparameterization. An important future direction is to study the dependence of the sample complexity of the entropy-regularized RL with respect to the dimension of the state space and improve the bound.

## ACKNOWLEDGMENT

This work was funded by grants from AFOSR, ARO, ONR, NSF and C3.ai Digital Transformation Institute.

## REFERENCES

- [1] R. J. Williams and J. Peng, "Function optimization using connectionist reinforcement learning algorithms," *Connection Science*, vol. 3, no. 3, pp. 241–268, 1991.
- [2] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *International conference on machine learning*. PMLR, 2016, pp. 1928–1937.
- [3] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *International*

- conference on machine learning*. PMLR, 2018, pp. 1861–1870.
- [4] B. O’Donoghue, R. Munos, K. Kavukcuoglu, and V. Mnih, “Combining policy gradient and Q-learning,” *arXiv preprint arXiv:1611.01626*, 2016.
- [5] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine, “Reinforcement learning with deep energy-based policies,” in *International Conference on Machine Learning*. PMLR, 2017, pp. 1352–1361.
- [6] H. Zang, X. Li, L. Zhang, P. Zhao, and M. Wang, “Teac: Intergrating trust region and max entropy actor critic for continuous control,” <https://openreview.net/references/pdf?id=bzTQQZQ6ix>, 2020.
- [7] B. D. Ziebart, *Modeling purposeful adaptive behavior with the principle of maximum causal entropy*. Carnegie Mellon University, 2010.
- [8] J. Schulman, P. Abbeel, and X. Chen, “Equivalence between policy gradients and soft Q-learning,” *CoRR*, vol. abs/1704.06440, 2017.
- [9] J. Mei, C. Xiao, C. Szepesvari, and D. Schuurmans, “On the global convergence rates of softmax policy gradient methods,” in *International Conference on Machine Learning*. PMLR, 2020, pp. 6820–6829.
- [10] G. Lan, “Policy mirror descent for reinforcement learning: Linear convergence, new sampling complexity, and generalized problem classes,” *Mathematical programming*, pp. 1–48, 2022.
- [11] S. Cen, C. Cheng, Y. Chen, Y. Wei, and Y. Chi, “Fast global convergence of natural policy gradient methods with entropy regularization,” *Operations Research*, 2021.
- [12] W. Chung, V. Thomas, M. C. Machado, and N. Le Roux, “Beyond variance reduction: Understanding the true impact of baselines on policy optimization,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 1999–2009.
- [13] J. Mei, B. Dai, C. Xiao, C. Szepesvari, and D. Schuurmans, “Understanding the effect of stochasticity in policy optimization,” *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [14] R. J. Williams, “Simple statistical gradient-following algorithms for connectionist reinforcement learning,” *Machine learning*, vol. 8, no. 3-4, pp. 229–256, 1992.
- [15] Z. Ahmed, N. Le Roux, M. Norouzi, and D. Schuurmans, “Understanding the impact of entropy on policy optimization,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 151–160.
- [16] A. Agarwal, S. M. Kakade, J. D. Lee, and G. Mahajan, “On the theory of policy gradient methods: Optimality, approximation, and distribution shift,” *J. Mach. Learn. Res.*, vol. 22, no. 98, pp. 1–76, 2021.
- [17] L. Xiao, “On the convergence rates of policy gradient methods,” *arXiv preprint arXiv:2201.07443*, 2022.
- [18] L. Shani, Y. Efroni, and S. Mannor, “Adaptive trust region policy optimization: Global convergence and faster rates for regularized MDPs,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, 2020, pp. 5668–5675.
- [19] J. Bhandari and D. Russo, “Global optimality guarantees for policy gradient methods,” *arXiv preprint arXiv:1906.01786*, 2019.
- [20] J. Zhang, C. Ni, C. Szepesvari, M. Wang *et al.*, “On the convergence and sample efficiency of variance-reduced policy gradient method,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 2228–2240, 2021.
- [21] J. Zhang, A. Koppel, A. S. Bedi, C. Szepesvari, and M. Wang, “Variational policy gradient method for reinforcement learning with general utilities,” in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 4572–4583.
- [22] G. Li, Y. Wei, Y. Chi, Y. Gu, and Y. Chen, “Softmax policy gradient methods can take exponential time to converge,” in *Conference on Learning Theory*. PMLR, 2021, pp. 3107–3110.
- [23] J. Mei, Y. Gao, B. Dai, C. Szepesvari, and D. Schuurmans, “Leveraging non-uniformity in first-order non-convex optimization,” *International Conference on Machine Learning*, 2021.
- [24] Y. Liu, K. Zhang, T. Basar, and W. Yin, “An improved analysis of (variance-reduced) policy gradient and natural policy gradient methods,” *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [25] Y. Ding, J. Zhang, and J. Lavaei, “On the global optimum convergence of momentum-based policy gradient,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2022, pp. 1910–1934.
- [26] B. Eysenbach and S. Levine, “If MaxEnt RL is the answer, what is the question?” *arXiv preprint arXiv:1910.01913*, 2019.
- [27] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [28] S. Cayci, N. He, and R. Srikant, “Linear convergence of entropy-regularized natural policy gradient with linear function approximation,” in *International conference on machine learning*. PMLR, 2021.
- [29] K. Zhang, A. Koppel, H. Zhu, and T. Basar, “Global convergence of policy gradient methods to (almost) locally optimal policies,” *SIAM Journal on Control and Optimization*, vol. 58, no. 6, pp. 3586–3612, 2020.
- [30] R. S. Sutton, D. A. McAllester, S. P. Singh, Y. Mansour *et al.*, “Policy gradient methods for reinforcement learning with function approximation,” in *NIPs*, vol. 99. Citeseer, 1999, pp. 1057–1063.
- [31] J. Baxter and P. L. Bartlett, “Infinite-horizon policy-gradient estimation,” *Journal of Artificial Intelligence Research*, vol. 15, pp. 319–350, 2001.
- [32] M. Benaïm, “Dynamics of stochastic approximation algorithms,” in *Seminaire de probabilites XXXIII*. Springer, 1999, pp. 1–68.
- [33] V. S. Borkar, *Stochastic approximation: a dynamical systems viewpoint*. Springer, 2009, vol. 48.
- [34] D. P. Bertsekas and J. N. Tsitsiklis, “Gradient convergence in gradient methods with errors,” *SIAM Journal on Optimization*, vol. 10, no. 3, pp. 627–642, 2000.
- [35] M. Benaïm and M. W. Hirsch, “Asymptotic pseudotra-

jectories and chain recurrent flows, with applications,” *Journal of Dynamics and Differential Equations*, vol. 8, no. 1, pp. 141–176, 1996.

- [36] H. J. Kushner and D. S. Clark, *Stochastic approximation methods for constrained and unconstrained systems*. Springer Science & Business Media, 2012, vol. 26.
- [37] P. Mertikopoulos, N. Hallak, A. Kavis, and V. Cevher, “On the almost sure convergence of stochastic gradient descent in non-convex problems,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 1117–1128, 2020.
- [38] J. Zhang, J. Kim, B. O’Donoghue, and S. Boyd, “Sample efficient reinforcement learning with REINFORCE,” *35th AAAI Conference on Artificial Intelligence*, 2021.
- [39] J. Mei, C. Xiao, B. Dai, L. Li, C. Szepesvári, and D. Schuurmans, “Escaping the gravitational pull of softmax,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 21 130–21 140, 2020.
- [40] M. L. Puterman, *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- [41] Y.-G. Hsieh, F. Iutzeler, J. Malick, and P. Mertikopoulos, “On the convergence of single-call stochastic extragradient methods,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [42] —, “Explore aggressively, update conservatively: Stochastic extragradient methods with variable stepsize scaling,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 16 223–16 234, 2020.
- [43] P. Mertikopoulos and Z. Zhou, “Learning in games with continuous action sets and unknown payoff functions,” *Mathematical Programming*, vol. 173, no. 1, pp. 465–507, 2019.
- [44] T. M. Cover, *Elements of information theory*. John Wiley & Sons, 1999.
- [45] O. Nachum, M. Norouzi, K. Xu, and D. Schuurmans, “Bridging the gap between value and policy based reinforcement learning,” *Advances in neural information processing systems*, vol. 30, 2017.
- [46] R. Yuan, R. M. Gower, and A. Lazaric, “A general sample complexity analysis of vanilla policy gradient,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2022, pp. 3332–3380.

## APPENDIX A

### PROPERTIES OF STOCHASTIC POLICY GRADIENT

#### A. Proof of Lemma 2

*Proof.* The gradient 2 follows from Proposition 2 in [28]. In particular, they consider the softmax parameterization with linear function approximation:

$$\pi_\theta(a | s) = \frac{\exp(\theta^\top \phi_{s,a})}{\sum_{a' \in \mathcal{A}} \exp(\theta^\top \phi_{s,a'})}$$

Our tabular softmax parameterization setting is a special case by taking  $\phi_{s,a} = e(s,a)$  where  $e(s,a) \in \mathbb{R}^{|\mathcal{S}| \cdot |\mathcal{A}|}$  is a vector where the  $(s,a)$ -th entry is equal to 1 while all other entries equal to 0.

The second part of the statement of Lemma 2 holds due to the following simple algebra:

$$\begin{aligned} \left\| \frac{\partial V_\lambda^\theta(\rho)}{\partial \theta} \right\| &\leq \frac{1}{1-\gamma} \max_{a,s} \|\nabla \log \pi_\theta(a | s)\| \\ &\quad \times \max_s \left\| \sum_a \pi_\theta(a | s) [Q_\lambda^\theta(s,a) - \lambda \log \pi_\theta(a | s)] \right\| \\ &= \frac{1}{1-\gamma} \max_{a,s} \|\nabla \log \pi_\theta(a | s)\| \times \max_s \|V_\lambda^\theta(s)\| \\ &\leq \frac{2(\bar{r} + \lambda \log |\mathcal{A}|)}{(1-\gamma)^2}. \end{aligned}$$

□

#### B. Proof of Lemma 4

*Proof.* We first show the unbiasedness of the Q-estimate, i.e.,  $\mathbb{E}[\hat{Q}_\lambda^\theta(s,a) | \theta, s, a] = Q_\lambda^\theta(s,a)$  for all  $(s,a) \in \mathcal{S} \times \mathcal{A}$  and  $\theta \in \mathbb{R}^d$ . In particular, from the definition of  $Q_\lambda^\theta(s,a)$ , we have

$$\begin{aligned} &\mathbb{E}[\hat{Q}_\lambda^\theta(s,a) | \theta, s, a] \\ &= \mathbb{E} \left[ r(s_0, a_0) + \sum_{h=1}^{H'} \gamma^{h/2} \cdot (r(s_h, a_h) - \lambda \log \pi_\theta(a_h | s_h)) \right. \\ &\quad \left. | \theta, s_0 = s, a_0 = a \right] \\ &= \mathbb{E} \left[ r(s_0, a_0) + \sum_{h=1}^{\infty} \mathbb{1}_{H' \geq h} \gamma^{h/2} \cdot (r(s_h, a_h) \right. \\ &\quad \left. - \lambda \log \pi_\theta(a_h | s_h)) | \theta, s_0 = s, a_0 = a \right], \end{aligned}$$

where we have replaced  $H'$  by  $\infty$  since we use the indicator function  $\mathbb{1}$  such that the summand for  $h \geq H'$  is null. In addition, by the law of total expectation, we have

$$\begin{aligned} &\mathbb{E}[\hat{Q}_\lambda^\theta(s,a) | \theta, s, a] \\ &= \mathbb{E}_{H'} \left[ \mathbb{E}_\tau \left[ r(s_0, a_0) + \sum_{h=1}^{\infty} \mathbb{1}_{H' \geq h} \gamma^{h/2} \cdot (r(s_h, a_h) \right. \right. \\ &\quad \left. \left. - \lambda \log \pi_\theta(a_h | s_h)) | \theta, s_0 = s, a_0 = a, H' \right] \right], \end{aligned} \tag{18}$$

where the trajectory  $\tau$  equal to  $\{s_0, a_0, s_1, a_1, \dots\}$ . The inner expectation over  $\tau$  can be written as

$$\begin{aligned} &\mathbb{E}_\tau \left[ r(s_0, a_0) + \sum_{h=1}^{\infty} \mathbb{1}_{H' \geq h} \gamma^{h/2} \cdot (r(s_h, a_h) \right. \\ &\quad \left. - \lambda \log \pi_\theta(a_h | s_h)) | \theta, s_0 = s, a_0 = a, H' \right] \\ &= \sum_\tau \left[ r(s_0, a_0) + \sum_{h=1}^{\infty} \mathbb{1}_{H' \geq h} \gamma^{h/2} \cdot (r(s_h, a_h) \right. \\ &\quad \left. - \lambda \log \pi_\theta(a_h | s_h)) \right] \cdot \mathbb{P}(\tau) | \theta, s_0 = s, a_0 = a, H' \\ &= r(s_0, a_0) + \sum_\tau \sum_{h=1}^{\infty} \left[ \mathbb{1}_{H' \geq h} \gamma^{h/2} \cdot (r(s_h, a_h) \right. \\ &\quad \left. - \lambda \log \pi_\theta(a_h | s_h)) \right] \cdot \mathbb{P}(\tau) | \theta, s_0 = s, a_0 = a, H'. \end{aligned} \tag{19}$$

By the definition of the probability over the sample trajectory  $\mathbb{P}(\tau)$ , for every  $h \in \{0, 1, 2, \dots\}$ , it holds that

$$\begin{aligned} & |(r(s_h, a_h) - \lambda \log \pi_\theta(a_h | s_h)) \cdot \mathbb{P}(\tau)| \\ &= |(r(s_h, a_h) - \lambda \log \pi_\theta(a_h | s_h)) \cdot \pi_\theta(a_h | s_h) \cdot \mathbb{P}(s_1 | s_0, a_0) \\ &\quad \cdot \pi_\theta(a_1 | s_1) \dots \cdot \mathbb{P}(s_h | s_{h-1}, a_{h-1}) \cdot \mathbb{P}(s_{h+1} | s_h, a_h) \dots \\ &\quad \cdot \mathbb{P}(s_{H'} | s_{H'-1}, a_{H'-1}) \cdot \pi_\theta(a_{H'} | s_{H'})| \\ &\leq \bar{r} + \frac{\lambda}{e}. \end{aligned}$$

where the last inequality follows from  $\mathbb{P}(s' | s, a) \leq 1$ ,  $\pi_\theta(a | s) \leq 1$  for all  $s, s' \in \mathcal{S}$  and  $a \in \mathcal{A}$  together with  $|x \log x| \leq \frac{1}{e}$  for  $x \in [0, 1]$ . Thus, for each trajectory  $\tau$  and  $N > 0$ , we have

$$\begin{aligned} & \sum_{h=1}^N [\mathbb{1}_{H' \geq h \geq 0} \gamma^{h/2} \cdot (r(s_h, a_h) - \lambda \log \pi_\theta(a_h | s_h))] \cdot \mathbb{P}(\tau) \\ &\leq \frac{1}{1 - \gamma^{1/2}} \left( \bar{r} + \frac{\lambda}{e} \right). \end{aligned} \quad (20)$$

Since left-hand side of (20) is non-decreasing and the limit as  $N \rightarrow \infty$  exists, by the Monotone Convergence Theorem, we can interchange the limit with the summation over the trajectory  $\tau$  in (19) as follows:

$$\begin{aligned} & \mathbb{E}_\tau \left[ r(s_0, a_0) + \sum_{h=1}^{\infty} \mathbb{1}_{H' \geq h \geq 0} \gamma^{h/2} \cdot (r(s_h, a_h) \right. \\ &\quad \left. - \lambda \log \pi_\theta(a_h | s_h)) \mid \theta, s_0 = s, a_0 = a, H' \right] \\ &= r(s_0, a_0) + \sum_{h=1}^{\infty} \sum_{\tau} [\mathbb{1}_{H' \geq h \geq 0} \gamma^{h/2} \cdot (r(s_h, a_h) \\ &\quad - \lambda \log \pi_\theta(a_h | s_h))] \cdot \mathbb{P}(\tau) \mid \theta, s_0 = s, a_0 = a, H' \\ &= r(s_0, a_0) + \sum_{h=1}^{\infty} \mathbb{E}_\tau [\mathbb{1}_{H' \geq h \geq 0} \gamma^{h/2} (r(s_h, a_h) \\ &\quad - \lambda \log \pi_\theta(a_h | s_h)) \mid \theta, s_0 = s, a_0 = a, H']. \end{aligned}$$

In addition, for every  $N > 0$ , we have

$$\begin{aligned} & \sum_{h=1}^N \mathbb{E}_\tau [\mathbb{1}_{H' \geq h \geq 0} \gamma^{h/2} (r(s_h, a_h) - \lambda \log \pi_\theta(a_h | s_h)) \\ &\quad \mid \theta, s_0 = s, a_0 = a, H'] + r(s_0, a_0) \\ &\leq \gamma^{1/2} \sum_{h=0}^{\infty} \mathbb{E}_\tau [\gamma^{(h+1)/2} (r(s_h, a_h) - \lambda \log \pi_\theta(a_h | s_h)) \\ &\quad \mid \theta, s_0 = s, a_0 = a, H'] + r(s_0, a_0) \\ &\leq r(s_0, a_0) + \gamma^{1/2} \mathbb{E}_{s_1} [V_{\lambda, \gamma/2}^\theta(s_1) \mid s_0, a_0] \\ &\leq \bar{r} + \frac{\gamma/2(\bar{r} + \lambda \log |\mathcal{A}|)}{1 - \gamma/2} \\ &\leq \frac{\bar{r} + \lambda \log |\mathcal{A}|}{1 - \gamma/2}, \end{aligned} \quad (21)$$

where the third inequality is due to the boundedness of the entropy-regularized value function in Lemma 1. Furthermore, since (21) is non-decreasing and the limit as  $N \rightarrow \infty$  exists,

by the Monotone Convergence Theorem, we can interchange the limit with the outer-expectation over  $H'$  in (18) as follows:

$$\begin{aligned} & \mathbb{E} [\hat{Q}_\lambda^\theta(s, a) \mid \theta, s, a] \quad (22) \\ &= r(s_0, a_0) + \lim_{N \rightarrow \infty} \mathbb{E}_{H'} \left[ \mathbb{E}_\tau \left[ \sum_{h=1}^N \mathbb{1}_{H' \geq h \geq 0} \gamma^{h/2} \cdot (r(s_h, a_h) \right. \right. \\ &\quad \left. \left. - \lambda \log \pi_\theta(a_h | s_h)) \mid \theta, s_0 = s, a_0 = a, H' \right] \right] \\ &= r(s_0, a_0) + \lim_{N \rightarrow \infty} \sum_{h=1}^N [\mathbb{E}_\tau [\mathbb{E}_{H'} [\mathbb{1}_{H' \geq h \geq 0}] \gamma^{h/2} \cdot (r(s_h, a_h) \\ &\quad - \lambda \log \pi_\theta(a_h | s_h)) \mid \theta, s_0 = s, a_0 = a]] \\ &= r(s_0, a_0) + \lim_{N \rightarrow \infty} \sum_{h=1}^N [\mathbb{E}_\tau [\gamma^h \cdot (r(s_h, a_h) \\ &\quad - \lambda \log \pi_\theta(a_h | s_h)) \mid \theta, s_0 = s, a_0 = a]] \\ &= r(s_0, a_0) + \mathbb{E}_\tau \left[ \sum_{h=1}^{\infty} [\gamma^h \cdot (r(s_h, a_h) \right. \\ &\quad \left. - \lambda \log \pi_\theta(a_h | s_h)) \mid \theta, s_0 = s, a_0 = a] \right] \\ &= Q_\lambda^\theta(s, a) \end{aligned}$$

where we have also used the fact that  $H'$  is drawn independently from the trajectory  $\tau$  in the first equality, the fact that  $H' \sim \text{Geom}(1 - \gamma^{1/2})$  and thus  $\mathbb{E}_{H'} [\mathbb{1}_{H' \geq h \geq 0}] = \gamma^{h/2}$  in the second equality, and the interchangeability between the limit and the summation over the trajectory  $\tau$  in the third equality. This completes the proof of the unbiasedness of  $\hat{Q}_\lambda^\theta(s, a)$ .

Now, we are ready to show unbiasedness of the stochastic gradients  $\hat{\nabla} V_\lambda^\theta(\rho)$ . It follows from Lemma 2 that

$$\begin{aligned} & \mathbb{E} [\hat{\nabla} V_\lambda^\theta(\rho) \mid \theta] \\ &= \mathbb{E}_{H, (s_H, a_H)} \left\{ \mathbb{E}_{H', (s'_{1:H'}, a'_{1:H'})} [\hat{\nabla} V_\lambda^\theta(\rho) \right. \\ &\quad \left. \mid \theta, s'_0 = s_H, a'_0 = a_H] \mid \theta \right\} \\ &= \mathbb{E}_{H, (s_H, a_H)} \left( \mathbb{E}_{H', (s'_{1:H'}, a'_{1:H'})} \left\{ \frac{1}{1 - \gamma} \nabla_\theta \log \pi_\theta(a'_0 | s'_0) \right. \right. \\ &\quad \left. \left. (\hat{Q}_\lambda^\theta(s'_0, a'_0) - \lambda \log \pi_\theta(a'_0 | s'_0)) \mid \theta, s'_0 = s_H, a'_0 = a_H, \right\} \mid \theta \right) \\ &= \mathbb{E}_{H, (s_H, a_H)} \left( \frac{1}{1 - \gamma} \nabla_\theta \log \pi_\theta(a_H | s_H) \mathbb{E}_{H', (s'_{1:H'}, a'_{1:H'})} \right. \\ &\quad \left. \{ (\hat{Q}_\lambda^\theta(s'_0, a'_0) - \lambda \log \pi_\theta(a'_0 | s'_0)) \mid \theta, s'_0 = s_H, a'_0 = a_H, \} \mid \theta \right) \\ &= \mathbb{E}_{H, (s_H, a_H)} \left\{ \frac{1}{1 - \gamma} \nabla_\theta \log \pi_\theta(a_H | s_H) \right. \\ &\quad \left. (Q_\lambda^\theta(s_H, a_H) - \lambda \log \pi_\theta(a_H | s_H)) \mid \theta \right\}. \end{aligned}$$

where we have used (22) in the last equality. By using the identity function  $\mathbb{1}_{h=H}$ , the above expression can be further written as

$$\begin{aligned} & \mathbb{E} [\hat{\nabla} V_\lambda^\theta(\rho) \mid \theta] \quad (23) \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{H, (s_H, a_H)} \left\{ \sum_{h=0}^{\infty} \mathbb{1}_{h=H} \nabla_\theta \log \pi_\theta(a_H | s_H) \right. \\ &\quad \left. (Q_\lambda^\theta(s_H, a_H) - \lambda \log \pi_\theta(a_H | s_H)) \mid \theta \right\}. \end{aligned}$$



Since for the softmax parameterization  $\pi_\theta$ , we have

$$\frac{\partial \log \pi_\theta(a_H | s_H)}{\partial \theta_{s,a}} = \begin{cases} -\pi_\theta(a|s)\pi_\theta(a_H | s_H), & (s,a) \neq (s_H, a_H), \\ \pi_\theta(a_H | s_H) - \pi_\theta(a_H | s_H)\pi_\theta(a_H | s_H), & (s,a) = (s_H, a_H). \end{cases}$$

Thus, the term  $\sum_{h=0}^{\infty} \mathbb{1}_{h=H} \nabla_\theta \log \pi_\theta(a_H | s_H) (Q_\lambda^\theta(s_H, a_H) - \lambda \log \pi_\theta(a_H | s_H))$  is uniformly bounded for every  $N > 0$  and non-decreasing with respect to  $N$ , we can interchange the limit and the expectation in (23) by the Monotone Convergence Theorem to obtain

$$\begin{aligned} & \mathbb{E}[\hat{\nabla} V_\lambda^\theta(\rho) | \theta] \\ &= \sum_{h=0}^{\infty} \frac{\mathbb{P}(H=h)}{1-\gamma} \cdot \mathbb{E}_{H, (s_H, a_H)} \{ \nabla_\theta \log \pi_\theta(a_H | s_H) \\ & \quad (Q_\lambda^\theta(s_H, a_H) - \lambda \log \pi_\theta(a_H | s_H)) | \theta \} \\ &= \sum_{h=0}^{\infty} \gamma^h \cdot \mathbb{E}_{(s_h, a_h)} \{ \nabla_\theta \log \pi_\theta(a_H | s_H) \\ & \quad (Q_\lambda^\theta(s_H, a_H) - \lambda \log \pi_\theta(a_H | s_H)) \} \\ &= \sum_{h=0}^{\infty} \gamma^h \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \mathbb{P}(s_h = s, a_h = a | s_0 \sim \rho, \theta) \nabla_\theta \log \pi_\theta(a | s) \\ & \quad (Q_\lambda^\theta(s, a) - \lambda \log \pi_\theta(a | s)) \\ &= \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \nabla_\theta \log \pi_\theta(a | s) (Q_\lambda^\theta(s, a) - \lambda \log \pi_\theta(a | s)) \\ & \quad \sum_{h=0}^{\infty} \gamma^h \mathbb{P}(s_h = s, a_h = a | s_0 \sim \rho, \theta) \\ &= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_\rho^{\pi_\theta}, a \sim \pi_\theta(\cdot | s)} [\nabla_\theta \log \pi_\theta(a | s) \\ & \quad (Q_\lambda^\theta(s, a) - \lambda \log \pi_\theta(a | s))]. \end{aligned}$$

where the second equality is due to the fact that  $H \sim \text{Geom}(1-\gamma)$  and thus  $\mathbb{P}(h=H) = (1-\gamma)\gamma^h$ , and the forth equality is due to the linearity of the integral and the finiteness of the state and action spaces. This completes the proof of unbiasedness of  $\hat{\nabla} V_\lambda^\theta(\rho)$ .  $\square$

### C. Proof of Lemma 5

*Proof.* We first note that the policy gradient estimator  $\hat{\nabla} V_\lambda^\theta(\rho)$  can be decomposed as:

$$\begin{aligned} & \nabla_\theta \log \pi_\theta(a_H | s_H) (\hat{Q}_\lambda^\theta(s_H, a_H) - \lambda \log \pi_\theta(s_H, a_H)) \\ &= \nabla_\theta \log \pi_\theta(a_H | s_H) \left( \sum_{i=0}^{H'} \gamma^{i/2} (r(s'_i, a'_i) - \lambda \log \pi_\theta(a'_i | s'_i)) \right), \end{aligned}$$

where  $H \sim \text{Geom}(1-\gamma)$ ,  $H' \sim \text{Geom}(1-\gamma^{1/2})$ ,  $(s^H, a^H) \sim \nu_\rho^{\pi_\theta}(s, a)$ ,  $s'_0 = s_H$ ,  $a'_0 = a_H$ . To streamline the presentation, we introduce the following notations:

$$\begin{aligned} g_1(s_H, a_H) &= \sum_{i=0}^{H'} \gamma^{i/2} r(s'_i, a'_i), \\ g_2(s_H, a_H) &= \sum_{i=0}^{H'} \gamma^{i/2} \lambda \log \pi_\theta(a'_i | s'_i), \end{aligned}$$

Then, the policy gradient estimator  $\hat{\nabla} V_\lambda^\theta(\rho)$  can be decomposed as:

$$\hat{\nabla} V_\lambda^\theta(\rho) = \frac{1}{1-\gamma} \nabla_\theta \log \pi_\theta(a_H | s_H) (g_1(s_H, a_H) - g_2(s_H, a_H)).$$

By the definition of the variance and the Cauchy-Schwarz inequality, we have

$$\begin{aligned} & \text{Var}(\hat{\nabla} V_\lambda^\theta(\rho)) \\ &= \frac{\|\nabla_\theta \log \pi_\theta(a_H | s_H)\|^2}{(1-\gamma)^2} \text{Var}(g_1(s_H, a_H) - g_2(s_H, a_H)) \\ &\leq \frac{2\|\nabla_\theta \log \pi_\theta(a_H | s_H)\|^2}{(1-\gamma)^2} (\text{Var}(g_1(s_H, a_H)) \\ & \quad + \text{Var}(g_2(s_H, a_H))) \\ &\leq \frac{8}{(1-\gamma)^2} (\text{Var}(g_1(s_H, a_H)) + \text{Var}(g_2(s_H, a_H))), \end{aligned}$$

where the last inequality follows from  $\|\nabla_\theta \log \pi_\theta(a_H | s_H)\| \leq 2$  [25]. Since  $g_1(s, a)$  is uniformly bounded, i.e.,  $\|g_1(s, a)\| \leq \frac{\bar{r}}{1-\gamma^{1/2}}$  for all  $s \in \mathcal{S}, a \in \mathcal{A}$ , we must have

$$\text{Var}(g_1(s_H, a_H)) \leq \frac{\bar{r}^2}{(1-\gamma^{1/2})^2}.$$

Then, it remains to prove the bounded variance of  $g_2$ . Firstly, it can be seen that

$$\begin{aligned} \|g_2\|^2 &\leq \lambda^2 \left( \sum_{i=0}^{H'} \gamma^{i/2} \log \pi_\theta(a'_i | s'_i) \right)^2 \\ &= \lambda^2 \left( \sum_{i=0}^{H'} \gamma^{i/4} \gamma^{i/4} \log \pi_\theta(a'_i | s'_i) \right)^2 \\ &\leq \lambda^2 \left( \sum_{i=0}^{H'} \gamma^{i/2} \right) \left( \sum_{i=0}^{H'} \gamma^{i/2} (\log \pi_\theta(a'_i | s'_i))^2 \right) \\ &\leq \frac{\lambda^2}{1-\gamma^{1/2}} \left( \sum_{i=0}^{H'} \gamma^{i/2} (\log \pi_\theta(a'_i | s'_i))^2 \right), \end{aligned}$$

where the second inequality is due to the Cauchy-Schwarz inequality. By fixing the state action pair  $(s_H, a_H)$  and the horizon  $H'$  for now and taking expectation of  $g_2$  only over the sample trajectory  $\tau' = \{s'_0, a'_0, \dots, s'_H, a'_H\}$ , it holds that

$$\begin{aligned} & \mathbb{E}_{\tau' \sim p(\tau' | \theta)} [\|g_2\|^2] \\ &\leq \frac{\lambda^2}{1-\gamma^{1/2}} \sum_{i=0}^{H'} \gamma^{i/2} \mathbb{E}_{\tau' \sim p(\tau' | \theta)} [(\log \pi_\theta(a'_i | s'_i))^2]. \end{aligned} \quad (24)$$

Since the realizations of  $a'_i$  and  $s'_i$  do not depend on the randomness in  $s'_{i+1}, a'_{i+1}, \dots, s'_H$ , we have

$$\begin{aligned} & \mathbb{E}_{\tau' \sim p(\tau' | \theta)} [(\log \pi_\theta(a'_i | s'_i))^2] \\ &= \mathbb{E}_{s'_1 \sim p(\cdot | a'_0, s'_0) \dots a'_{H-1} \sim \pi_\theta(\cdot | s'_{H-1}), s'_H \sim p(\cdot | s'_{H-1}, a'_{H-1})} [(\log \pi_\theta(a'_i | s'_i))^2] \\ &= \mathbb{E}_{s'_1 \sim p(\cdot | a'_0, s'_0) \dots a'_i \sim \pi_\theta(\cdot | s'_i)} [(\log \pi_\theta(a'_i | s'_i))^2] \\ &= \mathbb{E}_{s'_1 \sim p(\cdot | a'_0, s'_0) \dots s'_i \sim p(\cdot | a'_{i-1}, s'_{i-1})} \left[ \sum_{a'_i \in \mathcal{A}} \pi_\theta(a'_i | s'_i) (\log \pi_\theta(a'_i | s'_i))^2 \right]. \end{aligned}$$

By checking the optimality conditions for the optimization problem

$$\max \sum_{i=1}^n x_i (\log x_i)^2 \quad \text{such that} \quad \sum_{i=1}^n x_i = 1, \quad (25)$$

it can be concluded that the maximizer for the constrained problem (25) is  $x_1 = x_2 = \dots = x_n = \frac{1}{n}$  and the maximum solution is  $(\log n)^2$ .

Thus, we have  $\sum_{a_h \in \mathcal{A}} \pi_\theta(a_h | s_h) (\log \pi_\theta(a_h^i | s_h^i))^2 \leq (\log |\mathcal{A}|)^2$  and

$$\begin{aligned} & \mathbb{E}_{\tau' \sim p(\tau' | \theta)} \left[ (\log \pi_\theta(a_h^i | s_h^i))^2 \right] \\ & \leq \mathbb{E}_{s_0 \sim \rho, a_0 \sim \pi_\theta(\cdot | s_0), s_1 \sim p(\cdot | a_0, s_0) \dots s_h \sim p(\cdot | a_{h-1}, s_{h-1})} \left[ (\log |\mathcal{A}|)^2 \right] \\ & = (\log |\mathcal{A}|)^2. \end{aligned}$$

By substituting the above inequality into (24), we obtain that

$$\begin{aligned} \mathbb{E}_{\tau' \sim p(\tau' | \theta)} \left[ \|g_2\|^2 \right] & \leq \frac{\lambda^2}{1 - \gamma^{1/2}} \sum_{i=0}^{H'} \gamma^{i/2} \mathbb{E}_{\tau' \sim p(\tau' | \theta)} \left[ (\log \pi_\theta(a_i^i | s_i^i))^2 \right] \\ & \leq \frac{(\lambda \log |\mathcal{A}|)^2}{1 - \gamma^{1/2}} \sum_{i=0}^{H'} \gamma^{i/2} \\ & \leq \frac{(\lambda \log |\mathcal{A}|)^2}{(1 - \gamma^{1/2})^2}, \end{aligned}$$

for every  $H' > 0$ . By taking expectation of  $g_2$  over the state action pair  $(s_H, a_H)$  and the horizon  $H'$ , it yields that

$$\mathbb{E} \left[ \|g_2\|^2 \right] \leq \frac{(\lambda \log |\mathcal{A}|)^2}{(1 - \gamma^{1/2})^2},$$

which further implies that  $\text{Var} \left[ \|g_2\|^2 \right] \leq \frac{(\lambda \log |\mathcal{A}|)^2}{(1 - \gamma^{1/2})^2}$ . This completes the proof.  $\square$

#### D. Proof of Lemma 6

*Proof.* To simplify the notation, we define  $\tilde{r}_{j,\theta}^\lambda := r_j(s_j, a_j) - \lambda \log \pi_\theta(a_j | s_j)$ . By definition, we have

$$\begin{aligned} & \mathbb{E}[\hat{\nabla} V_\lambda^{\theta, H}(\rho)] - \nabla V_\lambda^\theta(\rho) \\ & = \mathbb{E} \left[ \sum_{h=0}^{H-1} \nabla \log \pi_\theta(a_h | s_h) \left( \sum_{j=H}^{\infty} \gamma^j \tilde{r}_{j,\theta}^\lambda \right) \right] \\ & \quad + \sum_{h=H}^{\infty} \nabla \log \pi_\theta(a_h | s_h) \left( \sum_{j=h}^{\infty} \gamma^j \tilde{r}_{j,\theta}^\lambda \right) \end{aligned}$$

Then, by the Cauchy-Schwarz inequality and the triangle inequality, we obtain

$$\begin{aligned} & \left\| \mathbb{E}[\hat{\nabla} V_\lambda^{\theta, H}(\rho)] - \nabla V_\lambda^\theta(\rho) \right\|_2 \\ & \leq \mathbb{E} \left[ \sum_{h=0}^{H-1} \left\| \nabla \log \pi_\theta(a_h | s_h) \right\| \left( \sum_{j=H}^{\infty} \gamma^j \tilde{r}_{j,\theta}^\lambda \right) \right] \\ & \quad + \mathbb{E} \left[ \sum_{h=H}^{\infty} \left\| \nabla \log \pi_\theta(a_h | s_h) \right\| \left( \sum_{j=h}^{\infty} \gamma^j \tilde{r}_{j,\theta}^\lambda \right) \right]. \end{aligned}$$

Since  $\|\nabla \log \pi_\theta(a | s)\|_2 \leq 2$  for all  $\theta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ , it holds that

$$\begin{aligned} & \mathbb{E}[\hat{\nabla} V_\lambda^{\theta, H}(\rho)] - \nabla V_\lambda^\theta(\rho)_2 \\ & \leq 2 \mathbb{E} \left[ \sum_{h=0}^{H-1} \left( \sum_{j=H}^{\infty} \gamma^j \tilde{r}_{j,\theta}^\lambda \right) \right] \end{aligned} \quad (26)$$

$$+ 2 \mathbb{E} \left[ \sum_{h=H}^{\infty} \left( \sum_{j=h}^{\infty} \gamma^j \tilde{r}_{j,\theta}^\lambda \right) \right]. \quad (27)$$

For the term in (26), we can rewrite it as

$$2 \mathbb{E}_{\tau^\infty} \left[ \sum_{h=0}^{H-1} \left( \sum_{j=H}^{\infty} \gamma^j \tilde{r}_{j,\theta}^\lambda \right) \right] = \sum_{\tau^\infty} \left[ \sum_{h=0}^{H-1} \left( \sum_{j=H}^{\infty} \gamma^j \tilde{r}_{j,\theta}^\lambda \right) \right] \cdot \mathbb{P}(\tau^\infty)$$

Since left-hand side of the above equation is non-decreasing and the limit as  $N \rightarrow \infty$  exists, by the Monotone Convergence Theorem, we can interchange the limit with the summation over the trajectory  $\tau^\infty$  in (26) as follows:

$$\begin{aligned} & 2 \mathbb{E}_{\tau^\infty} \left[ \sum_{h=0}^{H-1} \left( \sum_{j=H}^{\infty} \gamma^j \tilde{r}_{j,\theta}^\lambda \right) \right] \\ & = 2 \sum_{h=0}^{H-1} \left( \sum_{j=H}^{\infty} \gamma^j \mathbb{E}_{\tau^\infty} [r_j(s_j, a_j) - \lambda \log \pi_\theta(a_j | s_j)] \right). \end{aligned}$$

Due to  $-\sum_a \pi(a | s) \cdot \log \pi(a | s) \leq \log |\mathcal{A}|$ , the term in (26) can be upper bounded as

$$\begin{aligned} & 2 \mathbb{E}_{\tau^\infty} \left[ \sum_{h=0}^{H-1} \left( \sum_{j=H}^{\infty} \gamma^j \tilde{r}_{j,\theta}^\lambda \right) \right] \leq 2(\bar{r} + \lambda \log |\mathcal{A}|) \sum_{h=0}^{H-1} \left( \sum_{j=H}^{\infty} \gamma^j \right) \\ & \leq \frac{2(\bar{r} + \lambda \log |\mathcal{A}|) H \gamma^H}{1 - \gamma}. \end{aligned}$$

Similarly, we can interchange the limit with the summation over the trajectory  $\tau^\infty$  in (27) and upper bound it as

$$\begin{aligned} & 2 \mathbb{E} \left[ \sum_{h=H}^{\infty} \left( \sum_{j=h}^{\infty} \gamma^j \tilde{r}_{j,\theta}^\lambda \right) \right] \leq 2(\bar{r} + \lambda \log |\mathcal{A}|) \sum_{h=H}^{\infty} \sum_{j=h}^{\infty} \gamma^j \\ & \leq \frac{2(\bar{r} + \lambda \log |\mathcal{A}|) \gamma^H}{(1 - \gamma)^2}. \end{aligned}$$

Combining the above two inequalities, we have

$$\begin{aligned} & \left\| \mathbb{E}[\hat{\nabla} V_\lambda^{\theta, H}(\rho)] - \nabla V_\lambda^\theta(\rho) \right\|_2 \\ & \leq \frac{2(\bar{r} + \lambda \log |\mathcal{A}|) \gamma^H}{(1 - \gamma)} \left( H + \frac{1}{1 - \gamma} \right). \end{aligned}$$

This completes the proof.  $\square$

#### E. Proof of Lemma 7

*Proof. Step 1: Decomposition of the variance:* For the simplicity of the notation, we first define:

$$g_1(\tau^H | \theta, \rho) = \sum_{h=0}^{H-1} \left( \sum_{j=0}^h \nabla \log \pi_\theta(a_j | s_j) \right) \gamma^h r_h(s_h, a_h) \quad (28)$$

$$g_2(\tau^H | \theta, \rho) = \lambda \sum_{h=0}^{H-1} \left( \sum_{j=0}^h \nabla \log \pi_\theta(a_j | s_j) \right) (-\gamma^h \log \pi_\theta(a_h | s_h)). \quad (29)$$

By the definition of the variance and the Cauchy-Schwarz inequality, we have

$$\begin{aligned} & \text{Var}(\hat{\nabla} V_\lambda^{\theta, H}(\rho)) \\ &= \mathbb{E} \left[ (g_1(\tau^H|\theta, \rho) + g_2(\tau^H|\theta, \rho)) \right. \\ & \quad \left. - \mathbb{E}[g_1(\tau^H|\theta, \rho)] - \mathbb{E}[g_2(\tau^H|\theta, \rho)] \right]^2 \end{aligned} \quad (30)$$

$$\leq 3 \mathbb{E} \left[ (g_1(\tau^H|\theta, \rho) - \mathbb{E}[g_1(\tau^H|\theta, \rho)])^2 \right. \\ \left. + 3 \mathbb{E} \left[ (g_2(\tau^H|\theta, \rho) - \mathbb{E}[g_2(\tau^H|\theta, \rho)])^2 \right] \right] \quad (31)$$

$$= 3 \left( \text{Var}(g_1(\tau^H|\theta, \rho)) + \text{Var}(g_2(\tau^H|\theta, \rho)) \right). \quad (32)$$

**Step 2: Bounded variance of  $g_1$ :** As shown in Lemma 4.2 of [46], the fact that  $\|\nabla \log \pi_\theta(a|s)\|_2 \leq 2$  for all  $\theta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$  directly implies that  $\text{Var}(g_1(\tau^H|\theta, \rho)) \leq \frac{4\bar{r}^2}{(1-\gamma)^4}$  for all  $\theta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ .

**Step 3: Bounded variance of  $g_2$ :** Then, it remains to prove the bounded variance of  $g_2$ . Firstly, it can be observed that

$$\begin{aligned} \|g_2\| &= \lambda \left\| \sum_{h=0}^{H-1} \left( \sum_{j=0}^h \nabla \log \pi_\theta(a_j^i | s_j^i) \right) (-\gamma^h \log \pi_\theta(a_h^i | s_h^i)) \right\| \\ &\leq -\lambda \sum_{h=0}^{H-1} \left( \sum_{j=0}^h \|\nabla \log \pi_\theta(a_j^i | s_j^i)\| \right) \gamma^h \log \pi_\theta(a_h^i | s_h^i) \\ &\leq -2\lambda \sum_{h=0}^{H-1} (h+1) \gamma^h \log \pi_\theta(a_h^i | s_h^i). \end{aligned}$$

where the first inequality is due to the triangle inequality and the second inequality is due to  $\|\nabla \log \pi_\theta(a_j^i | s_j^i)\| \leq 2$ . Then, by taking the square of  $\|g_2\|$ , we obtain

$$\begin{aligned} \|g_2\|^2 &\leq 4\lambda^2 \left( \sum_{h=0}^{H-1} (h+1) \gamma^h \log \pi_\theta(a_h^i | s_h^i) \right)^2 \\ &= 4\lambda^2 \left( \sum_{h=0}^{H-1} (h+1) \sqrt{\gamma^h} \sqrt{\gamma^h} \log \pi_\theta(a_h^i | s_h^i) \right)^2 \\ &\leq 4\lambda^2 \left( \sum_{h=0}^{H-1} (h+1)^2 \gamma^h \right) \left( \sum_{h=0}^{H-1} \gamma^h (\log \pi_\theta(a_h^i | s_h^i))^2 \right) \\ &= 4\lambda^2 \left( \sum_{h=0}^{H-1} (h^2 + 2h + 1) \gamma^h \right) \left( \sum_{h=0}^{H-1} \gamma^h (\log \pi_\theta(a_h^i | s_h^i))^2 \right) \\ &\leq 4\lambda^2 \left( \frac{\gamma^2 + \gamma}{(1-\gamma)^3} + \frac{2\gamma}{(1-\gamma)^2} + \frac{1}{1-\gamma} \right) \\ & \quad \left( \sum_{h=0}^{H-1} \gamma^h (\log \pi_\theta(a_h^i | s_h^i))^2 \right) \\ &= 4\lambda^2 \left( \frac{\gamma+1}{(1-\gamma)^3} \right) \left( \sum_{h=0}^{H-1} \gamma^h (\log \pi_\theta(a_h^i | s_h^i))^2 \right) \end{aligned}$$

where the second inequality is due to the Cauchy-Schwarz inequality and the last inequality is due to  $\sum_{h=0}^{\infty} h^2 \gamma^h = \frac{\gamma^2 + \gamma}{(1-\gamma)^3}$ ,  $\sum_{h=0}^{\infty} h \gamma^h = \frac{\gamma}{(1-\gamma)^2}$  and  $\sum_{h=0}^{\infty} \gamma^h = \frac{1}{1-\gamma}$ .

By taking expectation of  $g_2$  over the sample trajectory  $\tau^H$ , it holds that

$$\begin{aligned} & \mathbb{E}_{\tau^H \sim p(\tau^H|\theta)} \left[ \|g_2\|^2 \right] \\ &\leq 4\lambda^2 \left( \frac{\gamma+1}{(1-\gamma)^3} \right) \sum_{h=0}^{H-1} \gamma^h \mathbb{E}_{\tau^H \sim p(\tau^H|\theta)} \left[ (\log \pi_\theta(a_h^i | s_h^i))^2 \right]. \end{aligned} \quad (33)$$

Since the realizations of  $a_h^i$  and  $s_h^i$  do not depend on the randomness in  $s_{h+1}, a_{h+1}, \dots, s_H$ , we have

$$\begin{aligned} & \mathbb{E}_{\tau^H \sim p(\tau^H|\theta)} \left[ (\log \pi_\theta(a_h^i | s_h^i))^2 \right] \\ &= \mathbb{E}_{s_0 \sim \rho, a_0 \sim \pi_\theta(\cdot|s_0), s_1 \sim p(\cdot|a_0, s_0) \dots a_{H-1} \sim \pi_\theta(\cdot|s_{H-1}), s_H \sim p(\cdot|s_{H-1}, a_{H-1})} \\ & \quad \left[ (\log \pi_\theta(a_h^i | s_h^i))^2 \right] \\ &= \mathbb{E}_{s_0 \sim \rho, a_0 \sim \pi_\theta(\cdot|s_0), s_1 \sim p(\cdot|a_0, s_0) \dots a_h \sim \pi_\theta(\cdot|s_h)} \left[ (\log \pi_\theta(a_h^i | s_h^i))^2 \right] \\ &= \mathbb{E}_{s_0 \sim \rho, a_0 \sim \pi_\theta(\cdot|s_0), s_1 \sim p(\cdot|a_0, s_0) \dots s_h \sim p(\cdot|a_{h-1}, s_{h-1})} \\ & \quad \left[ \sum_{a_h \in \mathcal{A}} \pi_\theta(a_h | s_h) (\log \pi_\theta(a_h^i | s_h^i))^2 \right]. \end{aligned}$$

Since the maximizer for the constrained problem

$$\max \sum_{i=1}^n x_i (\log x_i)^2 \quad \text{such that} \quad \sum_{i=1}^n x_i = 1,$$

is  $x_1 = x_2 = \dots = x_n = \frac{1}{n}$  and the maximum solution is  $(\log n)^2$ . Thus, we have  $\sum_{a_h \in \mathcal{A}} \pi_\theta(a_h | s_h) (\log \pi_\theta(a_h^i | s_h^i))^2 \leq (\log |\mathcal{A}|)^2$  and

$$\begin{aligned} & \mathbb{E}_{\tau^H \sim p(\tau^H|\theta)} \left[ (\log \pi_\theta(a_h^i | s_h^i))^2 \right] \\ &\leq \mathbb{E}_{s_0 \sim \rho, a_0 \sim \pi_\theta(\cdot|s_0), s_1 \sim p(\cdot|a_0, s_0) \dots s_H \sim p(\cdot|a_{H-1}, s_{H-1})} \left[ (\log |\mathcal{A}|)^2 \right] \\ &= (\log |\mathcal{A}|)^2. \end{aligned} \quad (34)$$

By combining (33) and (34), we have

$$\begin{aligned} \text{Var}(g_2) &\leq \mathbb{E}_{\tau^H \sim p(\tau^H|\theta)} \left[ \|g_2\|^2 \right] \\ &\leq 4\lambda^2 \left( \frac{\gamma+1}{(1-\gamma)^3} \right) \sum_{h=0}^{H-1} \gamma^h \mathbb{E}_{\tau \sim p(\tau|\theta)} \left[ (\log \pi_\theta(a_h^i | s_h^i))^2 \right] \\ &\leq 4\lambda^2 \left( \frac{\gamma+1}{(1-\gamma)^3} \right) \sum_{h=0}^{H-1} \gamma^h (\log |\mathcal{A}|)^2 \\ &\leq \frac{8\lambda^2 (\log |\mathcal{A}|)^2}{(1-\gamma)^4}. \end{aligned}$$

Finally, by substituting  $\text{Var}(g_1)$ ,  $\text{Var}(g_2)$  and  $\text{Var}(g_3)$  into (30), it holds that

$$\text{Var}(\hat{\nabla} V_\lambda^{\theta, H}(\rho)) \leq \frac{12\bar{r}^2 + 24\lambda^2 (\log |\mathcal{A}|)^2}{(1-\gamma)^4}.$$

This completes the proof.  $\square$

## APPENDIX B OTHER HELPFUL RESULTS.

*Lemma 19:* Suppose that  $f(x)$  is  $\bar{L}$ -smooth, i.e.,  $\|\nabla f(x) - \nabla f(y)\| \leq \bar{L} \|x - y\|$ . Given  $0 < \eta_t \leq \frac{1}{2\bar{L}}$  for all  $t \geq 1$ , let  $\{x_t\}_{t=1}^T$  be generated by  $x_{t+1} = x_t + \eta_t u_t$  and let  $e_t = u_t - \nabla f(x_t)$ . We have

$$f(x_{t+1}) \geq f(x_t) + \frac{\eta_t}{4} \|u_t\|_2^2 - \frac{\eta_t}{2} \|e_t\|_2^2.$$

*Proof.* Since  $f(x)$  is  $\bar{L}$ -smooth, one can write

$$\begin{aligned} & f(x_{t+1}) - f(x_t) - \langle u_t, x_{t+1} - x_t \rangle \\ &= f(x_{t+1}) - f(x_t) - \langle \nabla f(x_t), x_{t+1} - x_t \rangle \\ &\quad + \langle \sqrt{\eta_t}(\nabla f(x_t) - u_t), \frac{1}{\sqrt{\eta_t}}(x_{t+1} - x_t) \rangle \\ &\geq -\frac{2\bar{L}}{2} \|x_{t+1} - x_t\|^2 - \frac{b\eta_t}{2} \|\nabla f(x_t) - u_t\|_2^2 - \frac{1}{2b\eta_t} \|x_{t+1} - x_t\|_2^2 \\ &= \left(-\frac{1}{2b\eta_t} - \frac{2\bar{L}}{2}\right) \|x_{t+1} - x_t\|_2^2 - \frac{b\eta_t}{2} \|e_t\|_2^2, \end{aligned}$$

where the constant  $b > 0$  is to be determined later. By the above inequality and the definition of  $x_{t+1}$ , we have

$$\begin{aligned} & f(x_{t+1}) \\ &\geq f(x_t) + \langle u_t, x_{t+1} - x_t \rangle - \left(\frac{1}{2b\eta_t} + \frac{2\bar{L}}{2}\right) \|x_{t+1} - x_t\|_2^2 - \frac{b\eta_t}{2} \|e_t\|_2^2 \\ &= f(x_t) + \eta_t \|u_t\|^2 - \left(\frac{\eta_t}{2b} + \frac{2\bar{L}\eta_t^2}{2}\right) \|u_t\|_2^2 - \frac{b\eta_t}{2} \|e_t\|_2^2. \end{aligned}$$

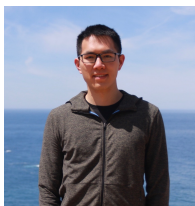
By choosing  $b = 1$  and using the fact that  $0 < \eta_t \leq \frac{1}{2\bar{L}}$ , we have

$$\begin{aligned} f(x_{t+1}) &\geq f(x_t) + \left(\frac{\eta_t}{2} - \frac{L\eta_t^2}{2}\right) \|u_t\|_2^2 - \frac{\eta_t}{2} \|e_t\|_2^2 \\ &\geq f(x_t) + \frac{\eta_t}{4} \|u_t\|_2^2 - \frac{\eta_t}{2} \|e_t\|_2^2. \end{aligned}$$

This completes the proof.  $\square$



**Hyunin Lee** is currently a Ph.D. student in Mechanical Engineering at the University of California, Berkeley. His research is focused on the reinforcement learning and the optimization theory. He obtained the B.S. degree in Mechanical Engineering from Seoul National University in 2022.



**Yuhao Ding** is currently a Ph.D. candidate in Industrial Engineering and Operations Research at the University of California, Berkeley. He has worked on different interdisciplinary problems in optimization and control theory. He obtained the B.E. degree in Aerospace Engineering from Nanjing University of Aeronautics and Astronautics in 2016, and the M.S. degree in Electrical and Computer Engineering from University of Michigan, Ann Arbor in 2018.



**Javad Lavaei** is currently an Associate Professor in the Department of Industrial Engineering and Operations Research at the University of California, Berkeley. He has worked on different interdisciplinary problems in power systems, optimization theory, control theory, and data science. He is an associate editor of the IEEE Transactions on Automatic Control, the IEEE Transactions on Smart Grid, and the IEEE Control System Letters.



**Junzi Zhang** is currently an applied scientist in Amazon.com LLC. He obtained a Ph.D. degree in Computational Mathematics from Stanford University and a B.S. degree in Applied Mathematics from School of Mathematical Sciences, Peking University. His research is focused on the design and analysis of optimization algorithms and software, and extends broadly into the fields of machine learning, causal inference and decision-making systems.