

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

Yuhao Ding, Junzi Zhang, 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].

Recently, 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. In the exact gradient setting, 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, these advantages have mostly been established for the true gradient setting and it is not fully understood whether any geometric property can be exploited to accelerate convergence to global optimality in inexact gradient settings. In the inexact gradient settings, it is proven in [11] that the NPG with the entropy regularization has a sample complexity of  $\tilde{O}(\frac{1}{\epsilon^2})$  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 open whether a global optimality convergence result and a low sample complexity can be obtained for the PG with entropy regularization in the practical 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. 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

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 the Amazon Advertising (work done prior to joining or outside of Amazon), Palo Alto, CA 94301 USA (e-mail: junziz@amazon.com).

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

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 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  $\tilde{O}(\frac{1}{\epsilon^2})$  for the proposed algorithm under the softmax parameterization. Our result is the first to achieve the sample complexity of  $\tilde{O}(\frac{1}{\epsilon^2})$  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

Stochastic policy gradient estimators with the original one-step entropy regularization has been proposed and adopted in [1, 2, 4]. For trajectory-level entropy regularization, an exact (visitation measure-based) policy gradient formula has been derived in [14] and later re-derived in the soft-max policy parametrization setting in [9], while stochastic policy gradient estimators have not been formally proposed or studied in the literature. The only exceptions are [3, 8], where [8] provides stochastic policy gradient estimators with a related but different trajectory-level KL regularization term, while [3] utilizes a reparametrization approach to reduce policy stochasticity to a fixed generating distribution. In contrast, we focus on the more widely used trajectory-level entropy regularization and classical likelihood-ratio-based estimators.

The theoretical understanding of policy-based methods has received considerable attention recently [9, 10, 11, 15, 16, 17, 18, 19, 20, 21]. Several techniques have been developed to improve standard PG and achieve a linear convergence rate, such as adding entropy regularization [9, 10, 11, 15], exploiting natural geometries based on Bregman divergences leading to NPG or policy mirror descent [10, 11, 16], and using a geometry-aware normalized PG (GNPG) approach to exploit the non-uniformity of the value function [22]. 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 [23, 24]. 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 notions  $\mathbb{E}_\xi[\cdot]$  and  $\mathbb{E}[\cdot]$  refer to the expectation over the random variable  $\xi$  and over all of the randomness. The notion  $\text{Var}[\cdot]$  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  $l_1$ -norm,  $l_2$ -norm and  $l_\infty$ -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  $a = \mathcal{O}(\mathcal{L}(x))$  means that  $a \leq C \cdot \mathcal{L}(x)$  for some constant  $C > 0$  that is independent of  $x$ . Similarly,  $a = \tilde{\mathcal{O}}(\mathcal{L}(x))$  indicates that the previous inequality may also depend on the function  $\log(\mathcal{L}(x))$ , where  $C > 0$  is again independent of  $x$ . We use  $\text{Geom}(x)$  to denote a geometric distribution with the parameter  $x$ .

## II. PRELIMINARIES

**Markov decision processes.** RL is generally modeled as a discounted Markov decision process (MDP) defined by a tuple  $(\mathcal{S}, \mathcal{A}; P; r; \gamma)$ . Here,  $\mathcal{S}$  and  $\mathcal{A}$  denote the finite state and action spaces;  $P(s'|s; a)$  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  $r(s; a) \in [0, r]$  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  $\theta(a|s)$  associated to the parameter  $\theta \in \mathbb{R}^d$ . Let  $\mathcal{D} = \{s_0; a_0; s_1; a_1; \dots\}$  denote the data of a sampled trajectory under policy  $\theta$  with the probability distribution over the trajectory as  $\mathcal{P}(\mathcal{D} | \theta) := \mathbb{P}(s_0) \prod_{h=1}^{\infty} \mathbb{P}(s_{h+1} | s_h; a_h) \theta(a_h | s_h)$ ; where  $\mathbb{P}(s)$  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

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

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, namely,  $\max_{\theta \in \mathbb{R}^d} V^\pi(\theta) := \mathbb{E}_{s_0 \sim \rho} [V^\pi(s_0)]$ . For the notational convenience, we will denote  $V^\pi(\theta)$  by the shorthand notation  $V^\theta(\cdot)$ .

**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  $P(S_h = s; a_h = a | s_0; \pi)$  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  $\pi$  for the performance measure is exploratory [9, 18], i.e.,  $\pi(s) > 0$  adequately covers the entire state distribution:

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

In practice, when the above assumption is not satisfied, we can optimize under another initial distribution  $\pi_0$ , i.e., the gradient is taken with respect to the optimization measure  $\pi$ , where  $\pi_0$  is usually chosen as an exploratory initial distribution that adequately covers the state distribution of some optimal policy. It is shown in [15] that the difficulty of the exploration problem faced by PG algorithms can be captured through the distribution mismatch coefficient defined as  $\left\| \frac{d}{\mu} \right\|_\infty$ , where  $\frac{d}{\mu}$  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  $\theta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ ,  $\theta(a|s)$  is chosen to be  $\frac{\exp(\theta_{sa})}{\sum_{a \in \mathcal{A}} \exp(\theta_{sa})}$ . 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  $\log \theta(a|s)$  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:  $\|\nabla \log \theta(a|s)\|_2 \leq 2$ ;  $\|\nabla^2 \log \theta(a|s)\|_2 \leq 1$ :

**RL with entropy regularization.** Entropy is a commonly used regularization in RL to promote exploration and discourage premature convergence to suboptimal policies [5, 8, 25]. It is far less aggressive in penalizing small probabilities, in comparison to other common regularizations such as log barrier functions [15]. 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

$$V_\lambda^\pi(\cdot) = V^\pi(\cdot) + H(\cdot; \lambda); \quad (1)$$

where  $\lambda \geq 0$  determines the strength of the penalty and  $H(\cdot; \lambda)$  stands for the discounted entropy defined as

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

Equivalently,  $V_\lambda^\pi(\cdot)$  can be viewed as the weighted value function of  $\pi$  by adjusting the instantaneous reward to be policy-dependent regularized version as  $r^\lambda(s; a) := r(s; a) - \log \theta(a|s)$ , 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  $a \in \mathcal{A}$ )

$$Q_\lambda^\pi(s; a) = r(s; a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} [V_\lambda^\pi(s')]; \\ V_\lambda^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} [-\log \theta(a|s) + Q_\lambda^\pi(s; a)];$$

**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:

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

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 [26]. The vanilla PG with exact gradient information and the entropy regularization is summarized in Algorithm 1.

---

#### Algorithm 1 Exact PG method

---

- 1: **Inputs:**  $\{\rho, \gamma\}_{t=1}^T, 1-\gamma$ .
  - 2: **for**  $t = 1; 2; \dots; T - 1$  **do**
  - 3:      $Q_{t+1} = Q_t + \gamma \nabla V_\lambda^{\theta_t}(\cdot)$ .
  - 4: **end for**
  - 5: **Outputs:**  $Q_T$ .
- 

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 (Q1):*  $V_\lambda^\theta(s) \leq \frac{\bar{r} + \lambda \log|\mathcal{A}|}{1-\gamma}$  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 Q2):* The entropy regularized PG with respect to  $\theta$  is

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

where

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

Furthermore, the entropy regularized PG is bounded, i.e.,  $\|\nabla V_\lambda^\theta(\cdot)\| \leq G$  for all  $\theta \in \mathcal{S}$  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  $\nabla V_\lambda^\theta(\cdot)$  is Lipschitz continuous.

*Lemma 3 (Lemmas 7 and 14 Q3):* [The PG  $\nabla V_\lambda^\theta(\cdot)$  is Lipschitz continuous with some constant  $L > 0$ , i.e.,

where the value of the Lipschitz constant is defined as

Challenges for designing entropy regularized PG estimators. Existing works either consider one-step entropy regularization [2, 28], 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  $r \log \pi$  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  $\nabla_{\theta} V^{\pi}$ , we need to first draw a state-action pair  $(s; a)$  from the distribution  $\pi^{\theta}$ ;  $\bullet$  and then obtain an unbiased estimate of the action-value function  $Q^{\pi}(s; a)$ . For the standard discounted finite-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  $r^{\pi}(s; a) \log \pi^{\theta}(s; a)$  could be unbounded when  $\pi^{\theta}(s; a) = 0$ . 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 the condition (i) is satisfied, 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  $Q^{\pi}(s; a)$  is evaluated, where the horizon  $H \sim \text{Geom}(1 - \beta)$ . It can be shown that  $(s_H; a_H) \sim \pi^{\theta}(s; a)$ . Moreover, given  $(s_H; a_H)$ , we will perform Monte-Carlo rollouts for another trajectory with the horizon  $H^{\text{roll}} \sim \text{Geom}(1 - \beta)$  independent of  $H$ , and estimate the advantage function value  $\hat{A}^{\pi}(s; a)$  along the trajectory  $(s_0^{\text{roll}}; a_0^{\text{roll}}; \dots; s_{H^{\text{roll}}}^{\text{roll}}; a_{H^{\text{roll}}}^{\text{roll}})$  with  $s_0^{\text{roll}} = s; a_0^{\text{roll}} = a$  as follows:

$$\hat{A}^{\pi}(s; a) = r^{\pi}(s_0^{\text{roll}}; a_0^{\text{roll}}) + \sum_{t=1}^{H^{\text{roll}}} \beta^{t-1} \left( r^{\pi}(s_t^{\text{roll}}; a_t^{\text{roll}}) - \log \pi^{\theta}(s_t^{\text{roll}}; a_t^{\text{roll}}) \right) \quad (3)$$

The subroutines of sampling one pair  $(s; a)$  from  $\pi^{\theta}$ ;  $\bullet$ , estimating  $\hat{A}^{\pi}(s; a)$ , and estimating  $\hat{Q}^{\pi}(s; a)$  are summarized as

Algorithm 2 Sam-SA: Sample from  $\pi^{\theta}$ ;  $\bullet$

---

```

1: Inputs:  $\pi^{\theta}$ ;  $\bullet$ 
2: Draw  $H \sim \text{Geom}(1 - \beta)$ .
3: Draw  $s_0$  and  $a_0 \sim \pi^{\theta}(s_0; a_0)$ .
4: for  $h = 1; 2; \dots; H - 1$  do
5:   Simulate the next state  $s_{h+1} \sim P^{\pi}(s_h; a_h)$  and action  $a_{h+1} \sim \pi^{\theta}(s_{h+1})$ .
6: end for
7: Outputs:  $s_H; a_H$ .

```

---

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

---

```

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

```

---

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

$$\nabla_{\theta} V^{\pi} = \mathbb{E} \left[ \sum_{h=0}^{\infty} \beta^h \left( r^{\pi}(s_{h+1}; a_{h+1}) - \log \pi^{\theta}(s_{h+1}; a_{h+1}) \right) \nabla_{\theta} Q^{\pi}(s_h; a_h) \right] \quad (4)$$

where  $(s_H; a_H) \sim \text{Sam-SA}(\pi^{\theta}; \bullet)$  and  $\hat{Q}$  is defined in (3). The following lemma shows that the stochastic PG is an unbiased estimator of  $\nabla_{\theta} V^{\pi}$ .

Lemma 4: For  $\nabla_{\theta} V^{\pi}$  defined in (4), we have  $\mathbb{E} \nabla_{\theta} V^{\pi} = \nabla_{\theta} V^{\pi}$ .

The next lemma shows that the proposed PG estimator  $\nabla_{\theta} V^{\pi}$  has a bounded variance even if it is unbounded when it approaches a deterministic policy.

Lemma 5: For  $\nabla_{\theta} V^{\pi}$  defined in (4), we have  $\text{Var} \nabla_{\theta} V^{\pi} \leq B^2$ ; where  $B^2 = \frac{8}{1 - \beta^2} \frac{\check{r}^2 \hat{\beta} \log \beta S^2}{1 - \beta^2}$ .

## C. Sampling the trajectory-based PG

Compared to the unbiased PG with a random horizon (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  $\nabla_{\theta} V^{\pi}$  can also be written as

$$\nabla_{\theta} V^{\pi} = \mathbb{E} \left[ \sum_{t=0}^{\infty} \beta^t \left( r^{\pi}(s_{t+1}; a_{t+1}) - \log \pi^{\theta}(s_{t+1}; a_{t+1}) \right) \nabla_{\theta} Q^{\pi}(s_t; a_t) \right];$$

where the expectation is taken over the trajectory distribution, i.e.,  $\pi^{\theta}(S)$ .

Since the distribution  $\pi^{\theta}(S)$  is unknown,  $\nabla_{\theta} V^{\pi}$  needs to be estimated from samples. The trajectory-based estimators include REINFORCE [28], 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  $\tilde{s}_0; a_0; s_1; \dots; s_{H-1}; a_{H-1}; s_H$  denote the truncation of the full trajectory of length  $H$ . Then, with the commonly used truncated GPOMDP, the truncated PG estimator for  $V^*$  can be written as:

$$\hat{V}^H = \sum_{h=0}^{H-1} \gamma^h \log \pi(a_h | s_h; \theta) - \sum_{h=0}^{H-1} \gamma^h r_h(a_h; s_h; \theta) + \log \pi(a_H | s_H; \theta) \quad (5)$$

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  $\hat{V}^H$  defined in (5), we have

$$|E[\hat{V}^H] - V^*| \leq \frac{2\gamma \log S}{1-\gamma} \leq \frac{1}{H}$$

From Lemma 6, we can observe that the bias is proportional to  $1/H$  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  $\hat{V}^H$  has a bounded variance even if it may be unbounded when approaches a deterministic policy.

Lemma 7: For  $\hat{V}^H$  defined in (5), we have

$$\text{Var}[\hat{V}^H] \leq B \frac{2\gamma^2 \log S^2}{1-\gamma^4}$$

#### D. Batched PG algorithms

In practice, we can sample and compute a batch of independently and identically distributed PG estimators  $\hat{V}_i^B$  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 away from zero if the exact PG is available. by iteratively running gradient-ascent-based algorithms, i.e.,  $\theta_{t+1} = \theta_t + \frac{1}{B} P_{i=1}^B \hat{V}_i^B$ ; where  $\alpha$  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:  $\gamma, \epsilon, B; T; \theta_1$ 
2: for  $t = 1; 2; \dots; T$  do
3:   for  $i = 1; 2; \dots; B$  do
4:      $s_{H_t}^i; a_{H_t}^i \sim \text{Sample}(S^A; \theta_t)$ 
5:      $\hat{Q}_{t,i}^B = \text{Est-EntQ}(s_{H_t}^i; a_{H_t}^i; t; \gamma)$ 
6:   end for
7:    $\theta_{t+1} = \theta_t + \frac{1}{B} \sum_{i=1}^B \hat{Q}_{t,i}^B \odot \log \pi(a_{H_t}^i | s_{H_t}^i; \theta_t)$ 
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 parameterization. However, our results in this section (and also the stationary point convergence result in Section 4C below)

can be easily extended to the general parameterizations as long as  $\log \pi(a | s; \theta)$  and  $\log \pi(a | s; \theta)$  are bounded for all  $s; a \in S \times A$ .

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 (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). 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  $V^*$  in (1) satisfies a non-uniform Łojasiewicz inequality as follows:

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

$$|V^* - V| \leq \frac{2}{\epsilon} C^2 \epsilon^{\frac{1}{d}} |V^* - V|$$

where

$$C = \frac{2}{\epsilon} \min_s \max_{s;a} \log \pi(a | s; \theta) + \frac{1}{\epsilon}$$

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  $\alpha = \frac{2}{\epsilon}$  for the entropy regularized objective, it holds that  $\inf_{\theta \in C_1} \min_{s;a} \log \pi(a | s; \theta) \geq \epsilon$ .

Remark 2: Note that by Algorithm 1  $\inf_{\theta \in C_1} \min_{s;a} \log \pi(a | s; \theta)$  is only dependent on the initialization and step-size (apart from problem dependent constants). Hence hereafter we denote  $c_1 = \inf_{\theta \in C_1} \min_{s;a} \log \pi(a | s; \theta)$ .

With Lemmas 3, 8 and 9, it is shown in Theorem 6 that the convergence rate for the entropy regularized PG is  $\frac{1}{\epsilon} \frac{1}{\epsilon^d}$  where the value of  $\epsilon$  depends on  $\inf_{\theta \in C_1} \min_{s;a} \log \pi(a | s; \theta) \geq \epsilon$  and  $\theta_{t-1}$  is generated by Algorithm 1. With a bad initialization  $\theta_1$ ,  $\min_{s;a} \log \pi(a | s; \theta_1)$  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 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  $V^{\lambda}(\theta) = \log \sum_{i=1}^2 e^{\theta_i r_i}$ .

Fig. 1. Landscape of  $V^{\lambda}(\theta) = \log \sum_{i=1}^2 e^{\theta_i r_i}$ .

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  $\theta_t^a$  is generated by Algorithm 4 for the entropy regularized objective with the step-sizes satisfying  $\sum_{t=1}^{\infty} \alpha_t = \infty$ ;  $\sum_{t=1}^{\infty} \alpha_t^2 < \infty$  and  $\alpha_t \leq \frac{2}{L}$  for all  $t = 1, 2, \dots$ . It holds that  $\lim_{t \rightarrow \infty} \mathbb{E} \| \nabla V^{\lambda}(\theta_t^a) \|^2 = 0$  with probability 1.

This result follows from classic results for the Robbins-Monro algorithm [4, 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 sizes 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 1 (see Appendix B) for the objective function  $V^{\lambda}(\theta)$  and the update rule  $\theta_{t+1}^a = \theta_t^a + \alpha_t u_t$ , where  $u_t = \nabla V^{\lambda}(\theta_t^a)$  and  $w_t = \nabla V^{\lambda}(\theta_t^a) - \nabla V^{\lambda}(\theta^*)$ .

1) From Lemma 3, we know that Condition 1 in 1 is satisfied with  $L = \frac{8r_{\max}^2}{\lambda} \frac{4 \log \frac{S}{\delta}}{\lambda}$ .

- 2) Condition 1 in 1 is satisfied by the definition of  $V^{\lambda}(\theta)$ .
- 3) Condition 1 in 1 is satisfied with  $c_1 = 1$  and  $c_2 = 1$ .
- 4) From Lemma 4 and 5, we know that Condition 4 in 1 is satisfied with  $A = \frac{8}{\lambda} \frac{r_{\max}^2 \log \frac{S}{\delta}}{\lambda}$ .
- 5) Condition 1 in 1 is satisfied by the definition of  $V^{\lambda}(\theta)$ .

In addition, it results from Lemma 1 we know that the entropy-regularized value function  $V^{\lambda}(\theta)$  is bounded. Thus, by Proposition 1, we must have  $\lim_{t \rightarrow \infty} \mathbb{E} \| \nabla V^{\lambda}(\theta_t^a) \|^2 = 0$  with probability 1. This completes the proof.

j

However, since the entropy-regularized value function is not coercive in  $\theta$  and it may be the case that the gradient  $\nabla V^{\lambda}(\theta)$  diminishing to 0 corresponds to  $\theta_t$  going to infinity instead of converging to a stationary point. In addition, the existing results [2, 35, 36] on the almost surely stationary point convergence rely on the assumption that the trajectories of the process are bounded, i.e.,  $\sup_{t \geq 0} \| Y_t \|^2 < \infty$  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, 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  $D^{\lambda}(\theta) = V^{\lambda}(\theta) - V^{\lambda}(\theta^*)$  denote the sub-optimality gap. Since the optimal policy  $\theta^*$  is unique [11], there must exist a continuum of optimal solutions

$$\mathbb{P}_{a \in \mathcal{A}} \frac{\exp \left\{ \sum_{s=1}^a \theta_s^* \right\}}{\sum_{s \in \mathcal{A}} \exp \left\{ \sum_{s=1}^a \theta_s^* \right\}} = \theta^* \quad \forall a \in \mathcal{A}; \theta_s^* \in \mathcal{A}; a > A.$$

In addition, we use  $\theta$  and  $\theta^{\dagger}$  interchangeably to denote the optimal policy of the entropy-regularized RL. Let  $\theta_t^{\dagger}$  denote the iterates of the algorithm with the exact PG (Algorithm 1) with  $\alpha_t = \frac{1}{2L}$  starting from the initial point  $\theta_1$ . For the soft-max parameterization, we have  $\theta_t^{\dagger} = \frac{1}{\sum_{s \in \mathcal{A}} e^{\theta_t^{\dagger}(s)}}$  for all  $s \in \mathcal{A}$ ;  $a > S$ , where  $C_s = \frac{1}{\sum_{s \in \mathcal{A}} e^{\theta_t^{\dagger}(s)}}$  are some constants. Then, we have

$$\min_{s \in \mathcal{A}} Z_t(s) \leq Z_t(\theta_t^{\dagger}) \leq \log \frac{1}{Z_t(\theta_t^{\dagger})} \leq \max_{s \in \mathcal{A}} Z_t(s); \quad \text{for all } t = 1, 2, \dots$$

Furthermore, by Lemma 9, we can define  $Y_{log}(\theta_t^{\dagger}) = \log \frac{1}{Z_t(\theta_t^{\dagger})}$ ; where  $c_{\log} = A_0$  is defined in Remark 2. Note that  $Y_{log}(\theta_t^{\dagger})$  is only dependent on  $\theta_t^{\dagger}$  and  $\mathcal{A}$  (apart from problem dependent constants), and  $Z_t(\theta_t^{\dagger}) \leq \log \frac{1}{Z_t(\theta_t^{\dagger})} \leq B$  for any  $\theta_t^{\dagger}$  and  $B = \frac{1}{2L}$ .

In addition, with a fair degree of hindsight and for some  $\epsilon_0$ , we define the stopping time for the iterates  $\theta_{t-1}^T$  as

$$\min_{t \geq t_0} t \text{Umin } Y_t \leq \epsilon_2 A \epsilon_0^{-1}; \quad (6)$$

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

$$G^0 \geq R^{\epsilon_2} \min_{t \geq t_0} Y_t \leq \epsilon_2 B \epsilon_0^{-1};$$

Finally, we define  $t^* = \min_{t \geq t_0} t \leq \epsilon_2$ . We are now ready to present the main result.

**Theorem 1:** Consider an arbitrary tolerance level  $\epsilon_0$  and a small enough tolerance level  $\epsilon_2$ . For every initial point  $\theta_1$ , if  $\tau_1$  is generated by Algorithm 5 with

$$\begin{aligned} T_1 &\leq C \frac{6D^2 \epsilon_0^{-1}}{C_0^2} \frac{8L}{C_0 \ln 2}; \quad T_2 \leq C \frac{t_0}{6} \quad t_0; \quad T \leq T_1 + T_2; \\ B_1 &\leq C \max\left\{ \frac{30}{C_0^2}, \frac{6T_1 \log T_1}{L} \right\}; \quad B_2 \leq C \frac{2 \ln^2 T_2}{6C} \quad t_0; \\ t &\leq B \min\left\{ \frac{\log T_1}{T_1 L}, \frac{8}{C_0}, \frac{1}{2L} \right\} \quad \text{for } t \leq T_1; \\ t &\leq \frac{1}{t - T_1} \quad \text{for } t \leq T_1 \end{aligned}$$

where

$$t_0 \leq \min\left\{ \frac{C}{3}, \frac{C}{2} \right\} \frac{\epsilon_0}{6 \ln 2} \frac{\epsilon_0^{-1}}{C_0^2} \exp\left\{ \frac{r}{C_0} \right\}; \quad (7)$$

$$t_0 \leq C \frac{3}{2} \quad (8)$$

$$C \leq \frac{2}{\epsilon_2} \min_{s,a} \epsilon_0^{-1} \epsilon_2 \min_{s,a} \epsilon_0^{-2} \epsilon_2 \epsilon_0^{-1} A_0; \quad (9)$$

$$C_0 \leq \frac{2}{\epsilon_2} \min_{s,a} \epsilon_0^{-1} \epsilon_2 \min_{s,a} \epsilon_0^{-2} \epsilon_2 \epsilon_0^{-1} A_0; \quad (10)$$

then we have  $P^D \geq \tau_1 \cdot B \cdot C \epsilon_0^{-1}$ . In total, it requires  $\Theta(\epsilon_0^{-2})$  samples to obtain an-optimal policy with high probability.

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

---

```

1: Inputs:  $\epsilon_0, \epsilon_2, \theta_1, B_1, B_2, T_1, T_2, \tau_1, \tau_1^T$ .
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:  $\tau_1$ 

```

---

## 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  $\epsilon_0^{-1-2\epsilon_2}$  dependency on the targeting tolerance. That being said, the polynomial dependency on  $\epsilon_0$  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 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 [15, 38]), as well as the “softmax gravity well” induced by the softmax parameterization which has also been observed in the exact gradient setting [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  $\theta_t$  remains in the region where it is uniformly bounded from below for all  $t$ . However, even with large batches, it is generally difficult to control stochastic trajectories, which eventually leads to the polynomial dependency on  $\epsilon_0$  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 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 VI, 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 VII. We then combine the two steps to prove Theorem 1 in Section VIII.

## VI. 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  $\epsilon_0$  and  $\epsilon_2$ . For every initial point  $\theta_1$ , if  $\tau_1$  is generated by Algorithm 4

<sup>2</sup>Note that similar difficulties in generalization from exact policy gradients to stochastic policy gradients have been observed in [39] which states that “unlike the true gradient setting, geometric information cannot be easily exploited in the stochastic case for accelerating policy optimization without detrimental consequences or impractical assumptions”.

with  $t \in B \min \{ \frac{\log T_1}{T_1 L}; \frac{8}{C_0^2}; \frac{1}{2L} \}$ ,  $T_1 \leq \frac{6D^*}{C_0} \cdot \frac{8L}{C_0^2 m^2}$ , and  $B_1 \max \{ \frac{30}{C_0^2}; \frac{6}{L} \}$ ,  $T_1 \log T_1 \leq \frac{6}{L}$ , then we have  $\mathbb{P}^{\mathcal{D}^*} \{ \|\hat{x}_t - x_t\| \leq C_1 \} \geq 1 - \epsilon$ .

A. Helpful lemmas

Lemma 12: Suppose that  $\hat{x}_t$  is L-smooth. Given  $\{x_t\}_{t=1}^T \in B_{\frac{1}{2L}}$  for all  $t \in C_1$ , let  $\{x_t\}_{t=1}^T$  be generated by a general update of the form  $x_{t+1} = x_t + \eta_t u_t$  and let  $\epsilon_t = u_t \odot f'(\hat{x}_t)$ . We have

$$\|f'(\hat{x}_{t+1}) - f'(\hat{x}_t)\| \leq \frac{1}{4} \eta_t Y_t^2 + \frac{1}{2} \epsilon_t Y_t^2.$$

Proof. Since  $f'$  is L-smooth, one can write

$$\begin{aligned} f'(\hat{x}_{t+1}) - f'(\hat{x}_t) &= f'(\hat{x}_t) - \eta_t u_t + \eta_t u_t + x_{t+1} - x_t e \\ f'(\hat{x}_{t+1}) - f'(\hat{x}_t) &= f'(\hat{x}_t) - \eta_t \epsilon_t + \eta_t \epsilon_t + x_{t+1} - x_t e \\ &= \eta_t \epsilon_t + \frac{1}{2} L \|x_{t+1} - x_t\|^2 e \end{aligned}$$

$$\begin{aligned} C \frac{1}{2} \eta_t Y_t^2 + \frac{1}{2} L \|x_{t+1} - x_t\|^2 &\leq \frac{1}{2} \eta_t Y_t^2 + \frac{1}{2} L \eta_t^2 \|u_t\|^2 \\ &\leq \frac{1}{2} \eta_t Y_t^2 + \frac{1}{2} L \eta_t^2 Y_t^2 \leq \frac{1}{2} \eta_t Y_t^2 \end{aligned}$$

where the constant  $A_0$  is to be determined later. By the above inequality and the definition of  $\epsilon_t$ , we have

$$\begin{aligned} f'(\hat{x}_{t+1}) - f'(\hat{x}_t) &= \eta_t \epsilon_t + \frac{1}{2} L \|x_{t+1} - x_t\|^2 e \\ &\leq \frac{1}{2} \eta_t Y_t^2 + \frac{1}{2} L \eta_t^2 Y_t^2 \leq \frac{1}{2} \eta_t Y_t^2 \end{aligned}$$

By choosing  $\eta_t = \frac{1}{2L}$  and using the fact that  $\|x_t\| \leq \frac{1}{2L}$ , we have

$$\begin{aligned} f'(\hat{x}_{t+1}) - f'(\hat{x}_t) &\leq \frac{1}{4} \eta_t Y_t^2 + \frac{1}{2} \epsilon_t Y_t^2 \\ &\leq \frac{1}{4} \eta_t Y_t^2 + \frac{1}{2} \epsilon_t Y_t^2 \end{aligned}$$

This completes the proof.  $\square$

To prove Lemma 11, we first consider the case when  $A \leq T_1$ , where  $\epsilon$  is defined in (6). Conditioning on this event, we can use Lemma 8 to show that  $\hat{x}_t$  is linearly convergent up to some aggregated estimation error.

Lemma 13: If  $t \in B \frac{1}{2L}$ , then  $\mathbb{E} \|\hat{x}_t - x_t\| \leq B \frac{C_0}{8} \cdot \frac{1}{T_1} \cdot \frac{5}{8C_0 B_1}$ .

Proof. Let  $\epsilon_t = \frac{1}{B_1} \mathbb{P}_{i=1}^{B_1} \odot V^{t,i} \odot u_t$ , where  $u_t = \frac{1}{B_1} \mathbb{P}_{i=1}^{B_1} \odot V^{t,i} \odot u_t$  and  $\odot V^{t,i} \odot u_t$  is an unbiased estimator of  $\odot V^{t,i} \odot u_t$ . Let  $\mathcal{F}_t$  denote the sigma field generated by the randomness up to iteration  $t$ . We define  $E^t = \mathbb{E} \{ \cdot | \mathcal{F}_t \}$  as the expectation operator conditioned on the sigma field  $\mathcal{F}_t$ . Since  $\odot V^{t,i} \odot u_t$

is L-smooth due to Lemma 3, it follows from Lemma 12 in the supplementary material:

$$\begin{aligned} E^t \|\hat{x}_{t+1} - x_{t+1}\| &\leq E^t \|\hat{x}_t - x_t\| + \frac{1}{4} \eta_t Y_t^2 \\ E^t \|V^{t,i} \odot u_t\| &\leq E^t \|u_t\| + \frac{1}{4} \eta_t Y_t^2 \\ BE^t \frac{1}{8} \eta_t Y_t^2 &\leq \frac{3}{4} \epsilon_t Y_t^2 + \frac{1}{4} \eta_t Y_t^2 \\ BE^t \frac{1}{8} \eta_t Y_t^2 &\leq E^t \|u_t\| + \frac{1}{4} \eta_t Y_t^2 \\ E^t \frac{1}{8} \eta_t Y_t^2 &\leq \frac{5}{8} \epsilon_t Y_t^2 + \frac{1}{4} \eta_t Y_t^2 \\ BE^t \frac{1}{8} \eta_t Y_t^2 &\leq \frac{5}{8} \epsilon_t Y_t^2 + \frac{1}{4} \eta_t Y_t^2 \end{aligned}$$

for every  $t \in B \frac{1}{2L}$ , where the second inequality uses the fact that  $u_t$  is an unbiased estimator of  $\odot V^{t,i} \odot u_t$  and the last inequality is due to Lemma 8. We now consider two cases:

- Case 1: Assume that  $A \leq t$ , which implies that  $t > G^0$  and  $C^0 \leq C \leq C^0$ . Then, we have  $E \|\hat{x}_t - x_t\| \leq B \frac{C_0}{8} \cdot \frac{1}{T_1} \cdot \frac{5}{8} E \epsilon_t Y_t^2 \leq \frac{1}{8} \cdot \frac{C_0}{8} \cdot \frac{1}{T_1} \cdot \frac{5}{8} E \epsilon_t Y_t^2$ .
- Case 2: Assume that  $A > t$  which leads to  $E \|\hat{x}_t - x_t\| \leq \frac{1}{8} \cdot \frac{C_0}{8} \cdot \frac{1}{T_1} \cdot \frac{5}{8} E \epsilon_t Y_t^2$ .

Now combining the above two cases yields the inequality

$$\begin{aligned} E \|\hat{x}_t - x_t\| &\leq \frac{1}{8} \cdot \frac{C_0}{8} \cdot \frac{1}{T_1} \cdot \frac{5}{8} E \epsilon_t Y_t^2 \\ &\leq \frac{1}{8} \cdot \frac{C_0}{8} \cdot \frac{1}{T_1} \cdot \frac{5}{8} E \epsilon_t Y_t^2 \leq \frac{1}{8} \cdot \frac{C_0}{8} \cdot \frac{1}{T_1} \cdot \frac{5}{8} E \epsilon_t Y_t^2 \end{aligned}$$

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

$$\begin{aligned} E \|\hat{x}_t - x_t\| &\leq \frac{1}{8} \cdot \frac{C_0}{8} \cdot \frac{1}{T_1} \cdot \frac{5}{8} E \epsilon_t Y_t^2 \\ &\leq \frac{1}{8} \cdot \frac{C_0}{8} \cdot \frac{1}{T_1} \cdot \frac{5}{8} E \epsilon_t Y_t^2 \end{aligned}$$

where the last equality uses the fact that  $t$  is a stopping time and the random variable  $\epsilon_t$  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

$$\begin{aligned} E \|\hat{x}_t - x_t\| &\leq \frac{1}{8} \cdot \frac{C_0}{8} \cdot \frac{1}{T_1} \cdot \frac{5}{8} E \epsilon_t Y_t^2 \\ &\leq \frac{1}{8} \cdot \frac{C_0}{8} \cdot \frac{1}{T_1} \cdot \frac{5}{8} E \epsilon_t Y_t^2 \\ &\leq \frac{1}{8} \cdot \frac{C_0}{8} \cdot \frac{1}{T_1} \cdot \frac{5}{8} E \epsilon_t Y_t^2 \end{aligned}$$

By setting  $t = T_1$ , we obtain that  $E \|\hat{x}_{T_1} - x_{T_1}\| \leq B \frac{C_0}{8} \cdot \frac{1}{T_1} \cdot \frac{5}{8} E \epsilon_{T_1} Y_{T_1}^2$ . This completes the proof.  $\square$

We now establish that  $\|\hat{x}_t - x_t\|$  will be bounded with high probability if the large batch size is used.

Lemma 14: It holds that  $\mathbb{P} \{ \|\hat{x}_t - x_t\| \leq \frac{1}{B_1} \cdot \frac{1}{T_1} \cdot \frac{L}{B_1} \} \geq 1 - \epsilon$ . Proof. By the triangle inequality and the fact that the iterations of the algorithm with the exact PG are bounded by we have

$$\|\hat{x}_t - x_t\| \leq \frac{1}{B_1} \cdot \frac{1}{T_1} \cdot \frac{L}{B_1} + \frac{1}{B_1} \cdot \frac{1}{T_1} \cdot \frac{L}{B_1} + \frac{1}{B_1} \cdot \frac{1}{T_1} \cdot \frac{L}{B_1}$$



Using the update rule of the algorithm with the exact PG obtain  
 $\mathbb{E}[V_i^{t+1}]$  and the stochastic PG  $\mathbb{E}[V_i^{t+1}]$ , one  
can write

$$\begin{aligned} \mathbb{E}[V_i^{t+1}] &= \mathbb{E}[Q_{i1}^t u_i \wedge Q_{i1}^t \mathbb{E}[V_i^t]] \\ &= \mathbb{E}[BQ_{i1}^t [u_i \odot V_i^t]] \\ &= \mathbb{E}[Q_{i1}^t [u_i \odot V_i^t \odot V_i^t \odot V_i^t]] \\ &= \mathbb{E}[BQ_{i1}^t Y_{e_1} Y_{e_2} Q_{i1}^t L Z_i \quad i_2] \end{aligned}$$

By expanding  $Z_i \quad i_2$  recursively, it can be concluded that

$$\begin{aligned} \mathbb{E}[d^t] &= \mathbb{E}[BQ_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 L Z_{t-1} \quad i_2 Z_t Q_{i1}^t L Z_i \quad i_2] \\ &= \mathbb{E}[BQ_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 L Q_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 L^2 Q_{i1}^t Z_i \quad i_2 Z_t Q_{i1}^t L Z_i \quad i_2] \\ &= \mathbb{E}[Q_{i1}^t L Z_i \quad i_2 Z_t Q_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 L Q_{i1}^t Y_{e_1} Y_{e_2} \quad i_1] \\ &= \mathbb{E}[Q_{i1}^t \alpha L \quad i_1 L^2 Z_i \quad i_2 Z_t Q_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 L Q_{i1}^t Y_{e_1} Y_{e_2} \quad i_1] \\ &= \mathbb{E}[BQ_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 L Q_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 \alpha L \quad i_1 L^2 Z_i \quad i_2 Z_t Q_{i1}^t Y_{e_1} Y_{e_2} \quad i_1] \\ &= \mathbb{E}[Q_{i1}^t \alpha^2 L \quad i_1 L^3 Z_i \quad i_2 Z_t Q_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 L Q_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 \alpha L \quad i_1 L^2 Z_i \quad i_2 Z_t] \\ &= \mathbb{E}[Q_{i1}^t \alpha^3 L \quad i_1 L^3 Z_i \quad i_2 Z_t Q_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 L Q_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 \alpha^2 L \quad i_1 L^2 Z_i \quad i_2 Z_t] \\ &= \mathbb{E}[Q_{i1}^t \alpha^4 L \quad i_1 L^3 Z_i \quad i_2 Z_t Q_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 L Q_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 \alpha^3 L \quad i_1 L^2 Z_i \quad i_2 Z_t] \\ &= \mathbb{E}[Q_{i1}^t M \quad i_1 L Y_{e_1} Y_{e_2} \quad i_1] \end{aligned}$$

Then, by the de nition of in (6) and Markov inequality, we

$$\begin{aligned} \mathbb{P}[d^t > C] &\leq \mathbb{P}[BQ_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 L M \quad i_1 L Y_{e_1} Y_{e_2} \quad i_1 > C] \\ &\leq \mathbb{P}[BQ_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 L E Y_{e_1} Y_{e_2} \quad i_1 > C] \\ &\leq \mathbb{P}[BQ_{i1}^t Y_{e_1} Y_{e_2} \quad i_1 L \cdot T_1^t E Y_{e_1} Y_{e_2} \quad i_1 > C] \end{aligned}$$

where we use the fact that  $\frac{1}{4}$  for all  $t > \sim 1; 2; \dots$ . Furthermore, since  $E Y_{e_1} Y_{e_2} \leq E Y_{e_1}^2 Y_{e_2}^2 \leq \frac{1}{B_1}$ , we have  $\mathbb{P}[d^t > C] \leq \frac{B_1 \cdot T_1^t}{B_1}$ : This completes the proof.

B. Proof of Lemma 11

Proof. By combining Lemmas 13 and 14, we obtain that

$$\begin{aligned} \mathbb{P}[d^t > C] &\leq \mathbb{P}[D^t > C] \\ &\leq \mathbb{P}[AT_1; D^t > C] \\ &\leq \frac{E[AT_1; D^t] \leq \mathbb{P}[BT_1]}{C} \\ &\leq \frac{C^0 \cdot T_1^{-1} D^t \leq 5^2}{8} \\ &\leq \frac{C^0 \cdot \frac{8}{C^0} \cdot C^0 T_1^{-1} D^t \leq 5^2}{8} \\ &\leq \frac{5^2}{C^0 B_1} \frac{T_1^{-1} L \cdot T_1^{-1}}{B_1} \\ &\leq \frac{1}{2} \frac{C^0 T_1^{-1} D^t \leq 5^2}{C^0 B_1} \frac{T_1^{-1} L \cdot T_1^{-1}}{B_1} \end{aligned}$$

where the second inequality holds due to the Markov inequality, and the last inequality holds because of  $\frac{1}{m} \leq B_1 \frac{1}{2}$  for all  $m \geq C$  and  $\frac{8}{C^0} \leq C$ . By taking  $B = \min\{\frac{\log T_1}{T_1 L}; \frac{8}{C^0}; \frac{1}{2L}\}$ , we obtain

$$\begin{aligned} \mathbb{P}[d^t > C] &\leq \mathbb{P}[D^t > C] \\ &\leq \frac{1}{2} \frac{C^0 \log T_1}{8L} D^t \leq 5^2 \frac{\log T_1^{-1} \cdot T_1^{-1}}{B_1 L} \\ &\leq \frac{1}{2} \frac{C^0 \log T_1}{8L} D^t \leq 5^2 \frac{\log T_1^{-1} \cdot \frac{T_1^{-1}}{\log T_1} \log T_1}{B_1 L} \\ &\leq \frac{1}{2} \frac{C^0 \log T_1}{8L} D^t \leq 5^2 \frac{\log T_1 \cdot T_1^{-1}}{B_1 L} \\ &\leq \frac{1}{T_1} \frac{C^0 \log T_1}{8L} D^t \leq 5^2 \frac{\log T_1 \cdot T_1^{-1}}{B_1 L} \end{aligned}$$

where we have used  $\hat{D}^t \leq B e^{-\alpha t}$  in the third inequality and  $a^{ln b} \leq b^{ln a}$  in the last inequality. To guarantee  $\hat{D}^t \leq C e^{-\alpha t}$ , it suffices to have

$$T_1 \leq \frac{6D^0}{\alpha} e^{-\alpha T_1}; \quad B_1 \leq \max\left\{\frac{30}{C}, \frac{6}{L}\right\} T_1 \log T_1;$$

This completes the proof.  $\square$

### VII. 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 15:** Given a tolerance level  $\epsilon > 0$ , let  $\pi^*$  be the optimal policy of  $V^*$ . Assume further that Algorithm 4 is run for  $T_2$  iterates with a step-size sequence of the form  $\alpha_t = \frac{1}{t}$  and a batch-size sequence  $\beta_t = \frac{1}{t}$  for all  $t \in \{1, 2, \dots, T_2\}$ . If  $t_0 \leq \frac{3}{2\alpha}$ , and  $\pi_1$  is initialized in a neighborhood  $\mathcal{U}_1$  such that

$$U_1 \cap \mathcal{M} > \epsilon A \cdot \mathbb{S} \hat{D}^* \cdot B_0 \checkmark; \quad (11)$$

where  $\epsilon_0 = \min\left\{\frac{\alpha}{6 \ln 2}, \frac{\alpha}{8}\right\} \checkmark \exp\left\{-\frac{1}{\alpha}\right\}$ ; and the constant  $\checkmark > 0$ , then the event

$$\mathcal{E}_{1; T_2} \supset \min_{s,a} \hat{\pi}_t^s \leq C \epsilon \quad \forall t \in \{1, 2, \dots, T_2\} \quad (12)$$

occurs with probability at least  $1 - \epsilon$ .

#### A. Helpful lemmas

To prove Lemma 15, we first characterize the maximum amount by which  $\hat{D}^t$  can grow at each step.

**Lemma 16:** Suppose that  $\hat{\pi}_t$  is generated by Algorithm 4 with  $0 < \alpha_t \leq \frac{1}{16r}$  for all  $t \in \mathbb{N}$ . We have

$$\hat{D}^t \leq B C \frac{C^t}{4} \hat{D}^0 + \frac{t}{2} \sum_{s=1}^t \frac{1}{4} \mathbb{Y}_s \mathbb{Y}_s^2; \quad (13)$$

where  $\mathbb{Y}_t = \mathbb{E}[\mathbb{Y}_t^2]$  and  $\mathbb{Y}_t = \mathbb{E}[\mathbb{Y}_t]$ .

**Proof.** Since  $V^*$  is  $L$ -smooth in light of Lemma 3, it follows from Lemma 12 that

$$\begin{aligned} \hat{D}^t &\leq \hat{D}^0 \\ &\leq \frac{t}{4} \sum_{s=1}^t \mathbb{Y}_s \mathbb{Y}_s^2 \\ &\leq \frac{t}{4} \sum_{s=1}^t \mathbb{Y}_s \mathbb{Y}_s \cdot \mathbb{Y}_s \leq \frac{t}{4} \sum_{s=1}^t \mathbb{Y}_s \mathbb{Y}_s^2 \\ &\leq \frac{t}{4} \sum_{s=1}^t \mathbb{Y}_s \mathbb{Y}_s \cdot \mathbb{Y}_s \leq \frac{t}{4} \sum_{s=1}^t \mathbb{Y}_s \mathbb{Y}_s^2 \\ &\leq \frac{t}{4} \sum_{s=1}^t \mathbb{Y}_s \mathbb{Y}_s \cdot \mathbb{Y}_s \leq \frac{t}{4} \sum_{s=1}^t \mathbb{Y}_s \mathbb{Y}_s^2 \\ &\leq \frac{t}{4} \sum_{s=1}^t \mathbb{Y}_s \mathbb{Y}_s \cdot \mathbb{Y}_s \leq \frac{t}{4} \sum_{s=1}^t \mathbb{Y}_s \mathbb{Y}_s^2 \end{aligned}$$

for every  $t \in \mathbb{N}$ , where the last inequality is due to Lemma 8.  $\square$

The quantity by which  $\hat{D}^t$  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 [37, 40, 41, 42], we now encode the error terms (13) as  $M_n = \sum_{t=1}^n \alpha_t \mathbb{Y}_t$  and  $S_n = \sum_{t=1}^n \alpha_t \mathbb{Y}_t^2$ . Since  $\mathbb{E}[\mathbb{Y}_t] = 0$ , we have  $\mathbb{E}[M_n] = 0$ . Therefore,  $M_n$  is a zero-mean martingale; likewise,  $\mathbb{E}[S_n] \leq C S_{n-1}$ , 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  $\hat{\pi}_t$  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  $\hat{D}^t$  remains close to  $D^*$ . We begin by introducing the cumulative mean square error  $R_n = \sum_{t=1}^n S_t$ : By construction, we have

$$\begin{aligned} R_n &\leq M_{n-1}^2 + S_{n-1} + \frac{1}{4} \mathbb{Y}_n \mathbb{Y}_n^2 \\ R_{n-1} &\leq 2M_{n-1}^2 + \frac{1}{4} \mathbb{Y}_n \mathbb{Y}_n^2 \end{aligned}$$

Hence,  $\mathbb{E}[R_n] \leq R_{n-1} + 2\mathbb{E}[M_{n-1}^2] + \frac{1}{4} \mathbb{E}[\mathbb{Y}_n \mathbb{Y}_n^2] \leq C R_{n-1}$ ; i.e.,  $R_n$  is a submartingale. With a fair degree of hindsight, we define as:

$$\mathcal{U} \supset \epsilon A \cdot \mathbb{S} \hat{D}^* \cdot B_0 \checkmark; \quad (14)$$

To condition it further, we also define the events

$$\begin{aligned} \mathcal{E}_n &= \hat{\pi}_t^s \leq C \epsilon \quad \forall t \in \{1, 2, \dots, n\} \\ \mathcal{E}_n &= \hat{D}_t \leq B_0 \quad \forall t \in \{1, 2, \dots, n\} \end{aligned}$$

By definition, we also have  $\mathcal{E}_0 = \mathcal{U}$  (because the set-building index set fork 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  $\hat{\pi}_t$  has escaped the vicinity of  $\pi^*$ .

Let the notation  $\mathbb{1}_A$  indicate the logical indicator of an event  $A$ , i.e.,  $\mathbb{1}_A = 1$  if  $A$  and  $0$  otherwise. For brevity, we write  $\mathbb{1}_n = \mathbb{1}_{\mathcal{E}_n}$  for the natural filtration of  $\mathbb{Y}_n$ . Now, we are ready to state the next lemma.

**Lemma 17:** Let  $\pi^*$  be the optimal policy. Then, for all  $n > 1$ , the following statements hold:

- 1)  $\mathbb{1}_{n-1} \leq \mathbb{1}_n$  and  $\mathbb{1}_{n-1} \leq \mathbb{1}_n$ .
- 2)  $\mathbb{1}_{n-1} \leq \mathbb{1}_n$ .
- 3) Consider the "large noise" event

$$\mathcal{E}_n = \mathbb{1}_{\mathcal{E}_n} \mathbb{1}_{\mathcal{E}_{n-1}} \mathbb{1}_{\mathcal{E}_{n-2}} \dots \mathbb{1}_{\mathcal{E}_0} \leq \checkmark \mathbb{1}_{\mathcal{U}}$$

and let  $\tilde{R}_n = \sum_{t=1}^n \mathbb{1}_{\mathcal{E}_t} S_t$  denote the cumulative error subject to the noise being "small" until time  $n$ . Then,

$$\mathbb{E}[R_n] \leq \mathbb{E}[\tilde{R}_n] + G^2 \sum_{t=1}^n \frac{1}{4B} \mathbb{P}[\mathcal{E}_t^c] \checkmark; \quad (15)$$

By convention, we write  $E_0 = g$  and  $R_0 = 0$ .

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

- ✓ For the base case  $n = 1$ , we have  $\|R_1 - R^*\|_{E_1} > U$  because  $R_1$  is initialized in  $U_1 \cap U$ . Since  $E_0 = g$ , our claim follows.
- ✓ For the inductive step, assume that  $\|R_{n-1} - R^*\|_{E_{n-1}} > U$  for some  $n \geq 2$ . To show that  $\|R_n - R^*\|_{E_n} > U$ , we fix a realization in  $E_n$  such that  $R_t \in B^c$  for all  $t = 1, 2, \dots, n$ . Since  $E_n \subseteq E_{n-1}$ , the inductive hypothesis posits that  $\|R_{n-1} - R^*\|_{E_{n-1}} > U$  for all  $t = 1, 2, \dots, n$ ; hence, it suffices to show that  $\|R_n - R_{n-1}\|_{E_n} > U$ . To that end, given that  $\|R_{n-1} - R^*\|_{E_{n-1}} > U$  for all  $t = 1, 2, \dots, n$ , the distance estimator (13) readily gives  $D_{t-1} \geq BD_{t-1}^* + \frac{1}{4} Y_{t-1}^2$  for all  $t = 1, 2, \dots, n$ . Therefore, after telescoping, we obtain

$$D_n \geq BD_n^* + \frac{1}{4} \sum_{t=1}^n Y_t^2 \geq \frac{1}{4} \sum_{t=1}^n Y_t^2 \geq \frac{1}{4} \sum_{t=1}^n Y_t^2 \geq \frac{1}{4} \sum_{t=1}^n Y_t^2$$

by the inductive hypothesis. This completes the induction. For Statement 3, we decompose  $R_n$  as

$$R_n = R_{n-1} + \frac{1}{n} \sum_{t=1}^n (R_t - R_{t-1})$$

where we have used the fact that  $R_0 = 0$  so  $R_{n-1} = \frac{1}{n-1} \sum_{t=1}^{n-1} R_t$  (recall that  $E_{n-1} \subseteq E_n$ ). Then, by the definition of  $R_n$ , we have

$$R_n - R_{n-1} = \frac{1}{n} \sum_{t=1}^n (R_t - R_{t-1}) - \frac{1}{n-1} \sum_{t=1}^{n-1} (R_t - R_{t-1})$$

and therefore

$$\|R_n - R_{n-1}\|_{E_n} \leq \frac{1}{n} \sum_{t=1}^n \|R_t - R_{t-1}\|_{E_n} \leq \frac{1}{n} \sum_{t=1}^n \|R_t - R_{t-1}\|_{E_{n-1}} \leq \frac{1}{n} \sum_{t=1}^n \|R_t - R_{t-1}\|_{E_{n-1}}$$

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

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

$$E \left[ \frac{1}{n} \sum_{t=1}^n (R_t - R_{t-1}) \mid E_{n-1} \right] = \frac{1}{n} \sum_{t=1}^n E \left[ R_t - R_{t-1} \mid E_{n-1} \right]$$

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

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

$$\frac{1}{4} \sum_{t=1}^n Y_t^2 \geq \frac{1}{4} \sum_{t=1}^n Y_t^2 \geq \frac{1}{4} \sum_{t=1}^n Y_t^2 \geq \frac{1}{4} \sum_{t=1}^n Y_t^2$$

Thus, putting together all of the above, we obtain  $\|R_n - R^*\|_{E_n} > U$ . Since  $R_{n-1} \in A^c$  if  $E_{n-1}$  occurs, we obtain  $\|R_n - R^*\|_{E_n} > U$ . This completes the proof of Statement 3.  $\square$

With the above results, we can show that the cumulative mean square error  $\|R_n - R^*\|_{E_n}$  is small with high probability at all times.

Lemma 18: Consider an arbitrary tolerance level  $\epsilon > 0$ . If Algorithm 4 is run with a step-size schedule of the form  $\eta_t = \frac{1}{t}$  where  $t_0 \in \mathbb{N}$  and a batch size schedule  $B_t \leq \frac{1}{t}$ , we have  $P(\|R_n - R^*\|_{E_n} > \epsilon) \leq \epsilon$  for all  $n \geq t_0$ . Proof. We begin by bounding the probability of the ‘‘large noise’’ event  $E_n = E_{n-1} \cap E_n$  as follows:

$$P(E_n) \leq P(\|R_n - R^*\|_{E_n} > \epsilon) \leq P(\|R_n - R^*\|_{E_n} > \epsilon) \leq P(\|R_n - R^*\|_{E_n} > \epsilon)$$

which is derived by using the fact that  $\|R_n - R^*\|_{E_n} \leq \frac{1}{n} \sum_{t=1}^n \|R_t - R^*\|_{E_n}$ . Now, by summing up (15), we conclude that  $E[\|R_n - R^*\|_{E_n}^2] \leq \frac{1}{n} \sum_{t=1}^n P(E_t) \leq \frac{1}{n} \sum_{t=1}^n \frac{1}{t^2} \leq \frac{1}{n}$ . Hence, combining the above results, we obtain the estimate

$$Q \leq \frac{1}{n} \sum_{t=1}^n \frac{1}{t^2} \leq \frac{1}{n} \sum_{t=1}^n \frac{1}{t^2} \leq \frac{1}{n}$$

where  $P_{t-1}^a \leq \frac{1}{t^2}$  and we have used the relations  $R_0 = 0$  and  $E_0 = g$  (by convention). By choosing  $t_0 \leq \frac{1}{\epsilon}$ , we ensure that  $\frac{1}{t_0} \leq \epsilon$ ; moreover, since the events  $E_t$  are disjoint for all  $t = 1, 2, \dots$ , we obtain  $P(\|R_n - R^*\|_{E_n} > \epsilon) \leq \sum_{t=1}^n P(E_t) \leq \frac{1}{n} \sum_{t=1}^n \frac{1}{t^2} \leq \frac{1}{n}$ . Hence,  $P(\|R_n - R^*\|_{E_n} > \epsilon) \leq \frac{1}{n}$  as claimed.  $\square$

Furthermore, we can show that the entropy-regularized value function  $V^* = \arg \min_{s \in S} \mathbb{E}[\sum_{t=1}^n \ell_t(s)]$  is locally quadratic around the optimal policy  $s^*$ .

Lemma 19: For every policy  $\pi$ , we have

$$D^* \leq C \frac{\min_{s \in S} \hat{s}^*}{2 \ln 2} \sum_{s \in S} \pi(s) \sum_{a \in A} \pi(a) \sum_{s \in S} \pi(s) \leq C \sum_{s \in S} \pi(s) \sum_{a \in A} \pi(a) \sum_{s \in S} \pi(s)$$

Proof. It follows from the soft sub-optimality difference lemma (Lemma 26 in [43]) that

$$V^* - V^{\pi} \leq \frac{1}{1 - \gamma} \sum_{s \in S} \pi(s) D_{KL}(\pi^* \parallel \pi) \sum_{s \in S} \pi(s) \leq \frac{1}{1 - \gamma} \sum_{s \in S} \pi(s) \sum_{a \in A} \pi(a) \sum_{s \in S} \pi(s)$$

$$C \frac{\min_{s \in S} \hat{s}^*}{2 \ln 2} \sum_{s \in S} \pi(s) \sum_{a \in A} \pi(a) \sum_{s \in S} \pi(s) \leq C \sum_{s \in S} \pi(s) \sum_{a \in A} \pi(a) \sum_{s \in S} \pi(s)$$

$$C \frac{\min_{s \in S} \hat{s}^*}{2 \ln 2} \sum_{s \in S} \pi(s) \sum_{a \in A} \pi(a) \sum_{s \in S} \pi(s) \leq C \sum_{s \in S} \pi(s) \sum_{a \in A} \pi(a) \sum_{s \in S} \pi(s)$$

where the first inequality is due to Theorem 11.6 [4] stating that

$$D_{KL}(P^* \circ S^* \circ C \frac{1}{2 \ln 2} Y^* \circ Q^* \circ Y_1^2)$$

for every two discrete distribution  $P^*$  and  $Q^*$ . Moreover, the second inequality is due to  $\|s^* - C^{-1} \cdot \hat{s}^*\|$  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 15

Proof. Since the sequence  $\epsilon_n$  is decreasing and  $\epsilon_n \leq \epsilon_{n-1}$  (by the second part of Lemma 17), Lemma 18 yields that  $P^*_{T_2} \circ C \inf_n P^*_{n-1} \circ C \inf_n P^*_{n-1} \circ C 1 - \epsilon$  provided that  $t_0$  is chosen large enough.

Now, it remains to show that  $\|D^*_{T_2} - B\| \leq \frac{\epsilon}{2}$  for all  $t = 1, 2, \dots, T_2$ . By Lemma 19, we have

$$\begin{aligned} & \mathbb{E} \left[ \frac{S^*_{t,a} \circ S^*_{t,a} \circ S^*_{t,a}}{2D^*_{t,a} \circ \ln 2} \right] \leq \mathbb{E} \left[ \frac{S^*_{t,a} \circ S^*_{t,a} \circ S^*_{t,a}}{2^2 \epsilon_0 \circ \ln 2} \right] \\ & \leq \frac{6}{\min_{s,a} \hat{s}^*} \leq B \exp^{-\frac{r}{1}} \cdot B \min_{s,a} \hat{s}^* \end{aligned}$$

where the second inequality is due to the condition that the event  $T_2$  occurs, the third inequality is due to  $\epsilon_0 \leq B$ , the fourth inequality is due to the definition of  $\epsilon_0$ , and the last inequality is due to Theorem 1 [4] where it holds that  $\log \frac{S^*_{t,a} \circ S^*_{t,a} \circ S^*_{t,a}}{\hat{s}^*_{t,a} \circ \hat{s}^*_{t,a} \circ \hat{s}^*_{t,a}} \leq C \frac{r}{1} \cdot \hat{s}^*_{t,a} \circ \hat{s}^*_{t,a} \circ \hat{s}^*_{t,a} > S \circ A$ .

Now, it can be easily verified that  $\mathbb{E} \left[ \frac{S^*_{t,a} \circ S^*_{t,a} \circ S^*_{t,a}}{\min_{s,a} \hat{s}^*} \right] \leq C \frac{r}{1} \cdot \hat{s}^*_{t,a} \circ \hat{s}^*_{t,a} \circ \hat{s}^*_{t,a}$ . For every  $t = 1, 2, \dots, T_2$ , let  $s; a = \operatorname{argmin}_{s,a} \hat{s}^*_{t,a} \circ S^*_{t,a} \circ S^*_{t,a}$ . One can write

$$\min_{s;a} \hat{s}^*_{t,a} \circ S^*_{t,a} \circ S^*_{t,a} \leq C \frac{r}{1} \cdot \hat{s}^*_{t,a} \circ \hat{s}^*_{t,a} \circ \hat{s}^*_{t,a} \leq \min_{s;a} \hat{s}^*_{t,a} \circ S^*_{t,a} \circ S^*_{t,a} \leq C \frac{r}{1} \cdot \min_{s;a} \hat{s}^*_{t,a} \circ S^*_{t,a} \circ S^*_{t,a}$$

where the last inequality is due to  $\hat{s}^*_{t,a} \circ S^*_{t,a} \circ S^*_{t,a} \leq C \min_{s;a} \hat{s}^*_{t,a} \circ S^*_{t,a} \circ S^*_{t,a}$  for every  $s > S$  and  $a > A$ . Thus, we obtain  $\|D^*_{T_2} - B\| \leq \frac{\epsilon}{2}$ . This completes the proof.  $\square$

VIII. PROOF OF THEOREM 1

From Lemma 11, we conclude that, with a large batch size, the iteration will converge to a neighborhood of the optimal solution with high probability. From Lemma 15, 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. Proof. From Lemma 11, we can conclude that  $\|D^*_{T_1} - B\| \leq \frac{\epsilon}{2}$  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  $\|D^*_{T_1} - B\| \leq \frac{\epsilon}{2}$ . It follows from Lemma 16 that

$$\mathbb{E} \left[ \frac{D^*_{t+1} \circ 1_{T_1}}{4} \right] \leq \mathbb{E} \left[ \frac{D^*_{t+1} \circ 1_{T_1}}{2} \right] \leq \mathbb{E} \left[ \frac{D^*_{t+1} \circ 1_{T_1}}{4} \right] \leq \frac{t}{4} \mathbb{E} \left[ Y_{t+1}^2 \right]$$

for all  $t \leq T_1$ , where  $t \leq T_1$  and  $t \leq T_1$  is defined in (12). When the event  $T_1$  occurs, we have  $\|D^*_{T_1} - B\| \leq \frac{\epsilon}{2}$ , where  $B$  is defined in (9). By taking the expectation, we can obtain

$$\begin{aligned} \mathbb{E} \left[ \frac{D^*_{t+1} \circ 1_{T_1}}{2} \right] & \leq \mathbb{E} \left[ \frac{D^*_{t+1} \circ 1_{T_1}}{4} \right] \leq \frac{t}{4} \mathbb{E} \left[ Y_{t+1}^2 \right] \\ \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{2} \right] & \leq \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \leq \frac{t}{4} \mathbb{E} \left[ Y_{t+1}^2 \right] \\ \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] & \leq \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \leq \frac{t}{4B} \end{aligned}$$

where the first equality is due to the fact that  $T_1$  is deterministic conditioning on  $F_t$ , the second equality is due to the unbiasedness of conditioning on  $F_t$ , and the first inequality is due to (17). Therefore,  $\mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \leq \frac{t}{4B}$ . Arguing inductively yields that

$$\begin{aligned} \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] & \leq \frac{t}{4} \mathbb{E} \left[ Y_{t+1}^2 \right] \\ \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] & \leq \frac{t}{4} \mathbb{E} \left[ Y_{t+1}^2 \right] \leq \frac{t}{4} \mathbb{E} \left[ Y_{t+1}^2 \right] \\ \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] & \leq \frac{t}{4} \mathbb{E} \left[ Y_{t+1}^2 \right] \leq \frac{t}{4} \mathbb{E} \left[ Y_{t+1}^2 \right] \end{aligned}$$

By taking  $t = \frac{4}{C \cdot \epsilon_0}$ , we obtain that

$$\mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \leq \frac{t}{4} \mathbb{E} \left[ Y_{t+1}^2 \right] \leq \frac{t}{4} \mathbb{E} \left[ Y_{t+1}^2 \right] \leq \frac{t}{4} \mathbb{E} \left[ Y_{t+1}^2 \right]$$

By the law of total probability and the Markov inequality, we obtain that

$$\begin{aligned} P^* \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] & \leq P^* \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \leq P^* \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \\ P^* \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] & \leq P^* \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \leq P^* \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \\ P^* \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] & \leq P^* \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \leq P^* \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \\ \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] & \leq \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \leq \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \\ \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] & \leq \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \leq \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \\ \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] & \leq \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \leq \mathbb{E} \left[ \frac{D^*_{T_1} \circ 1_{T_1}}{4} \right] \end{aligned}$$

where the second inequality follows from Lemma 15. To guarantee that  $\|D^*_{T_1} - B\| \leq \frac{\epsilon}{2}$ , it suffices to have  $\frac{t}{6} \mathbb{E} \left[ Y_{t+1}^2 \right] \leq \frac{\epsilon}{2}$ . This completes the proof.

IX. 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  $\tilde{O}(\frac{1}{\epsilon})$ , 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.

#### 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," *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," *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: Integrating 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," *arXiv preprint arXiv: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," *arXiv preprint arXiv:2102.00135*, 2021.
- [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] 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.
- [15] A. Agarwal, S. M. Kakade, J. D. Lee, and G. Mahajan, "On the theory of policy gradient methods: Optimality, approximation, and distribution shift," *arXiv preprint arXiv:1908.00261*, 2019.
- [16] L. Xiao, "On the convergence rates of policy gradient methods," *arXiv preprint arXiv:2201.07443*, 2022.
- [17] 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.
- [18] J. Bhandari and D. Russo, "Global optimality guarantees for policy gradient methods," *arXiv preprint arXiv:1906.01786*, 2019.
- [19] J. Zhang, C. Ni, Z. Yu, C. Szepesvari, and M. Wang, "On the convergence and sample efficiency of variance-reduced policy gradient method," *arXiv preprint arXiv:2102.08607*, 2021.
- [20] 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.
- [21] 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.
- [22] 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.
- [23] 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.
- [24] Y. Ding, J. Zhang, and J. Lavaei, "On the global convergence of momentum-based policy gradient," *arXiv preprint arXiv:2110.10116*, 2021.
- [25] B. Eysenbach and S. Levine, "If MaxEnt RL is the answer, what is the question?," *arXiv preprint arXiv:1910.01913*, 2019.

[26] R. S. Sutton and A. G. Barto, Reinforcement learning: An introduction MIT press, 2018.

[27] S. Cayci, N. He, and R. Srikant, “Linear convergence of entropy-regularized natural policy gradient with linear function approximation,” arXiv preprint arXiv:2106.04096 2021.

[28] R. J. Williams, “Simple statistical gradient-following algorithms for connectionist reinforcement learning,” Machine learning vol. 8, no. 3-4, pp. 229–256, 1992.

[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 pseudotrajectories 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,” arXiv preprint arXiv:2006.11144 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] Y.-G. Hsieh, F. Lutzeler, J. Malick, and P. Mertikopoulos, “On the convergence of single-call stochastic extragradient methods,” arXiv preprint arXiv:1908.08465 2019.

[41] —, “Explore aggressively, update conservatively: Stochastic extragradient methods with variable stepsize scaling,” arXiv preprint arXiv:2003.10162 2020.

[42] 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.

[43] J. Mei, C. Xiao, C. Szepesvari, and D. Schuurmans, “On the global convergence rates of softmax policy gradient methods,” International Conference on Machine Learning PMLR, 2020, pp. 6820–6829.

[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,” arXiv preprint arXiv:1702.08892 2017.

[46] Y. Ding, J. Zhang, and J. Lavaei, “On the global convergence of momentum-based policy gradient,” arXiv preprint arXiv:2110.10116 2021.

[47] R. Yuan, R. M. Gower, and A. Lazaric, “A general sample complexity analysis of vanilla policy gradient,” arXiv preprint arXiv:2107.11433 2021.

[48] 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.

APPENDIX A

PROPERTIES OF STOCHASTIC POLICY GRADIENT

A. Proof of Lemma 4

Proof. We first show the unbiasedness of the Q-estimate, i.e.,

$$E \hat{Q}^s; a^s; S; s; a = Q^s; a^s \text{ for all } \hat{s}; a^s > S \setminus A \text{ and } s \in S.$$

In particular, from the definition of  $\hat{Q}^s; a^s$ , we have

$$E \hat{Q}^s; a^s; S; s; a = \sum_{h=1}^H \gamma^{h-1} \sum_{a_h} \mathbb{1}_{\{a_h = a\}} \log \hat{a}_h \mathbb{S}_h^{\bullet\bullet} S; s_0; s; a_0; a$$

$$E \hat{r}^s; a_0 = \sum_{h=1}^H \gamma^{h-1} \sum_{a_h} \mathbb{1}_{\{a_h = a\}} \log \hat{a}_h \mathbb{S}_h^{\bullet\bullet} S; s_0; s; a_0; a;$$

where we have replaced  $\mathbb{1}_{\{a_h = a\}}$  by  $\mathbb{1}_{\{a_h = a\}}$  since we use the indicator function  $\mathbb{1}$  such that the summand for  $C \setminus H^{\text{oe}}$  is null. In addition, by the law of total expectation, we have

$$E \hat{Q}^s; a^s; S; s; a = \sum_{a'} \mathbb{1}_{\{a' = a\}} E \hat{r}^s; a_0 = \sum_{a'} \mathbb{1}_{\{a' = a\}} \sum_{h=1}^H \gamma^{h-1} \sum_{a_h} \mathbb{1}_{\{a_h = a\}} \log \hat{a}_h \mathbb{S}_h^{\bullet\bullet} S; s_0; s; a_0; a; H^{\text{oe}}; \tag{18}$$

where the trajectory equal to  $\tilde{s}_0; a_0; s_1; a_1; \dots$ . The inner expectation over  $\tilde{\cdot}$  can be written as

$$E \hat{r}^s; a_0 = \sum_{a'} \mathbb{1}_{\{a' = a\}} \sum_{h=1}^H \gamma^{h-1} \sum_{a_h} \mathbb{1}_{\{a_h = a\}} \log \hat{a}_h \mathbb{S}_h^{\bullet\bullet} S; s_0; s; a_0; a; H^{\text{oe}}$$

$$= \sum_{a'} \mathbb{1}_{\{a' = a\}} \sum_{h=1}^H \gamma^{h-1} \sum_{a_h} \mathbb{1}_{\{a_h = a\}} \log \hat{a}_h \mathbb{S}_h^{\bullet\bullet} P^{\bullet} \cdot S; s_0; s; a_0; a; H^{\text{oe}}$$

$$= \sum_{a'} \mathbb{1}_{\{a' = a\}} \sum_{h=1}^H \gamma^{h-1} \sum_{a_h} \mathbb{1}_{\{a_h = a\}} \log \hat{a}_h \mathbb{S}_h^{\bullet\bullet} P^{\bullet} \cdot S; s_0; s; a_0; a; H^{\text{oe}} \tag{19}$$

By the definition of the probability over the sample trajectory  $\hat{P}^h$ , for every  $\epsilon > 0$ ;  $1; 2; \dots$ , it holds that

$$\begin{aligned} & \mathbb{P}(\hat{r}_{s_0; a_0}^h - \log \hat{a}_h \mathbb{S}_h \leq \epsilon) \leq \mathbb{P}(\hat{r}_{s_0; a_0}^h - \log \hat{a}_h \mathbb{S}_h \leq \epsilon) \\ & \mathbb{P}(\hat{r}_{s_0; a_0}^h - \log \hat{a}_h \mathbb{S}_h \leq \epsilon) \leq \mathbb{P}(\hat{r}_{s_0; a_0}^h - \log \hat{a}_h \mathbb{S}_h \leq \epsilon) \\ & \mathbb{P}(\hat{r}_{s_0; a_0}^h - \log \hat{a}_h \mathbb{S}_h \leq \epsilon) \leq \mathbb{P}(\hat{r}_{s_0; a_0}^h - \log \hat{a}_h \mathbb{S}_h \leq \epsilon) \end{aligned}$$

By  $\frac{1}{e}$ :

where the last inequality follows from  $\frac{1}{e} \leq \log x \leq \frac{1}{e}$  for  $x > 0; 1$ . Thus, for each trajectory and  $N \geq A_0$ , we have

$$\mathbb{P}(\hat{r}_{s_0; a_0}^h - \log \hat{a}_h \mathbb{S}_h \leq \frac{1}{e}) \leq \frac{1}{1 - \frac{1}{e}} \quad (20)$$

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 in (19) as follows:

$$\begin{aligned} & \mathbb{E} \hat{r}_{s_0; a_0}^h = \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \\ & \mathbb{E} \hat{r}_{s_0; a_0}^h = \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \\ & \mathbb{E} \hat{r}_{s_0; a_0}^h = \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \end{aligned}$$

In addition, for every  $N \geq A_0$ , we have

$$\begin{aligned} & \mathbb{E} \hat{r}_{s_0; a_0}^h \leq \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \\ & \mathbb{E} \hat{r}_{s_0; a_0}^h \leq \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \\ & \mathbb{E} \hat{r}_{s_0; a_0}^h \leq \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \end{aligned} \quad (21)$$

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,

the limit with the outer-expectation over  $\mathbb{P}^h$  in (18) as follows:

$$\begin{aligned} & \mathbb{E} \hat{r}_{s_0; a_0}^h = \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \\ & \mathbb{E} \hat{r}_{s_0; a_0}^h = \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \\ & \mathbb{E} \hat{r}_{s_0; a_0}^h = \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \end{aligned} \quad (22)$$

where we have also used the fact that  $\mathbb{S}_h$  is drawn independently from the trajectory in the first equality, the fact that  $\mathbb{P}^h \sim \text{Geom}(1 - \frac{1}{e})$  and thus  $\mathbb{E} \mathbb{S}_h \sim \frac{1}{1 - \frac{1}{e}}$  in the second equality, and the interchangeability between the limit and the summation over the trajectory in the third equality. This completes the proof of the unbiasedness of  $\hat{r}_{s_0; a_0}^h$ .

Now, we are ready to show unbiasedness of the stochastic gradients  $\hat{V}^h$ . It follows from Lemma 2 that

$$\begin{aligned} & \mathbb{E} \hat{V}^h = \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \\ & \mathbb{E} \hat{V}^h = \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \\ & \mathbb{E} \hat{V}^h = \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \end{aligned}$$

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

$$\begin{aligned} & \mathbb{E} \hat{V}^h = \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \\ & \mathbb{E} \hat{V}^h = \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \\ & \mathbb{E} \hat{V}^h = \mathbb{E} \left[ \frac{1}{h} \sum_{h=1}^N \log \hat{a}_h \mathbb{S}_h \right] \end{aligned} \quad (23)$$

Since for the softmax parameterization, we have

$$\frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} = \frac{\pi_{s; a} - \hat{\pi}_{s; a}}{\hat{\pi}_{s; a}}$$

Thus, the term  $\mathbb{E}_{h \sim P_{h_0}} \frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} \log \hat{\pi}_{s; a}$  is uniformly bounded for every  $N \geq 0$  and non-decreasing with respect to  $N$ , we can interchange the limit and the expectation by the Monotone Convergence Theorem to obtain

$$\begin{aligned} \mathbb{E} \frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} &= \mathbb{E}_{h \sim P_{h_0}} \frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} \log \hat{\pi}_{s; a} \\ &= \mathbb{E}_{h \sim P_{h_0}} \frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} \log \hat{\pi}_{s; a} \\ &= \mathbb{E}_{h \sim P_{h_0}} \frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} \log \hat{\pi}_{s; a} \\ &= \mathbb{E}_{h \sim P_{h_0}} \frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} \log \hat{\pi}_{s; a} \end{aligned}$$

where the second equality is due to the fact that  $\text{Geom}(1)$  is a geometric distribution and thus  $\mathbb{E}[h] = 1$ , 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  $\hat{V}$ .

### B. Proof of Lemma 5

Proof. We first note that the policy gradient estimator  $\hat{V}$  can be decomposed as:

$$\hat{V} = \mathbb{E}_{s \sim \pi} \left[ \frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} \log \hat{\pi}_{s; a} \right]$$

where  $H \sim \text{Geom}(1)$ ;  $H^{\text{oe}} \sim \text{Geom}(1-2^{-H})$ ,  $\hat{s}_H; a_H$  is a state-action pair. To streamline the presentation, we introduce the following notations:

$$g_1^{\hat{s}_H; a_H} = \mathbb{E}_{i \sim H^{\text{oe}}} \left[ r^i \hat{\pi}_{s_i^{\text{oe}}, a_i^{\text{oe}}} \right];$$

$$g_2^{\hat{s}_H; a_H} = \mathbb{E}_{i \sim H^{\text{oe}}} \left[ \log \hat{\pi}_{s_i^{\text{oe}}, a_i^{\text{oe}}} \right];$$

Then, the policy gradient estimator  $\hat{V}$  can be decomposed as:

$$\hat{V} = \frac{1}{1} \mathbb{E}_{s \sim \pi} \left[ \frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} g_1^{\hat{s}_H; a_H} - g_2^{\hat{s}_H; a_H} \right]$$

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

$$\begin{aligned} \text{Var}(\hat{V}) &\leq \mathbb{E} \left[ \left( \frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} \right)^2 \right] \text{Var}(g_1^{\hat{s}_H; a_H} - g_2^{\hat{s}_H; a_H}) \\ &\leq \mathbb{E} \left[ \left( \frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} \right)^2 \right] \left( \text{Var}(g_1^{\hat{s}_H; a_H}) + \text{Var}(g_2^{\hat{s}_H; a_H}) \right) \end{aligned}$$

where the last inequality follows from  $\mathbb{E} \left[ \left( \frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} \right)^2 \right] \leq B$  [46]. Since  $g_1^{\hat{s}_H; a_H}$  is uniformly bounded, i.e.  $|g_1^{\hat{s}_H; a_H}| \leq B$  for all  $s > S; a > A$ , we must have

$$\text{Var}(g_1^{\hat{s}_H; a_H}) \leq B \frac{r^2}{1-2^{-2}}$$

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

$$\begin{aligned} \mathbb{E} \left[ \left( \frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} \right)^2 \right] &\leq \mathbb{E}_{i \sim H^{\text{oe}}} \left[ \left( \frac{\partial \log \hat{\pi}_{s_i^{\text{oe}}, a_i^{\text{oe}}}}{\partial \theta_{s_i^{\text{oe}}, a_i^{\text{oe}}}} \right)^2 \right] \\ &\leq \mathbb{E}_{i \sim H^{\text{oe}}} \left[ \left( \frac{\partial \log \hat{\pi}_{s_i^{\text{oe}}, a_i^{\text{oe}}}}{\partial \theta_{s_i^{\text{oe}}, a_i^{\text{oe}}}} \right)^2 \right] \\ &\leq \mathbb{E}_{i \sim H^{\text{oe}}} \left[ \left( \frac{\partial \log \hat{\pi}_{s_i^{\text{oe}}, a_i^{\text{oe}}}}{\partial \theta_{s_i^{\text{oe}}, a_i^{\text{oe}}}} \right)^2 \right] \\ &\leq \mathbb{E}_{i \sim H^{\text{oe}}} \left[ \left( \frac{\partial \log \hat{\pi}_{s_i^{\text{oe}}, a_i^{\text{oe}}}}{\partial \theta_{s_i^{\text{oe}}, a_i^{\text{oe}}}} \right)^2 \right] \end{aligned}$$

where the second inequality is due to the Cauchy-Schwarz inequality. By fixing the state-action pairs  $\hat{s}_H; a_H$  and the horizon  $H^{\text{oe}}$  for now and taking expectation only over the sample trajectory  $s_0^{\text{oe}}, a_0^{\text{oe}}, \dots, s_H^{\text{oe}}, a_H^{\text{oe}}$ , it holds that

$$\mathbb{E}_{s \sim \pi} \left[ \left( \frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} \right)^2 \right] \leq \mathbb{E}_{i \sim H^{\text{oe}}} \left[ \mathbb{E}_{s_0^{\text{oe}}, a_0^{\text{oe}}, \dots, s_H^{\text{oe}}, a_H^{\text{oe}}} \left[ \left( \frac{\partial \log \hat{\pi}_{s_i^{\text{oe}}, a_i^{\text{oe}}}}{\partial \theta_{s_i^{\text{oe}}, a_i^{\text{oe}}}} \right)^2 \right] \right] \quad (24)$$

Since the realizations of  $s_i^{\text{oe}}$  and  $a_i^{\text{oe}}$  do not depend on the randomness in  $s_1^{\text{oe}}, a_1^{\text{oe}}, \dots, s_H^{\text{oe}}, a_H^{\text{oe}}$ , we have

$$\begin{aligned} \mathbb{E}_{s \sim \pi} \left[ \left( \frac{\partial \log \hat{\pi}_{s; a}}{\partial \theta_{s; a}} \right)^2 \right] &\leq \mathbb{E}_{s_0^{\text{oe}}, a_0^{\text{oe}}, \dots, s_H^{\text{oe}}, a_H^{\text{oe}}} \left[ \mathbb{E}_{i \sim H^{\text{oe}}} \left[ \left( \frac{\partial \log \hat{\pi}_{s_i^{\text{oe}}, a_i^{\text{oe}}}}{\partial \theta_{s_i^{\text{oe}}, a_i^{\text{oe}}}} \right)^2 \right] \right] \\ &\leq \mathbb{E}_{s_0^{\text{oe}}, a_0^{\text{oe}}, \dots, s_H^{\text{oe}}, a_H^{\text{oe}}} \left[ \mathbb{E}_{i \sim H^{\text{oe}}} \left[ \left( \frac{\partial \log \hat{\pi}_{s_i^{\text{oe}}, a_i^{\text{oe}}}}{\partial \theta_{s_i^{\text{oe}}, a_i^{\text{oe}}}} \right)^2 \right] \right] \end{aligned}$$



By checking the optimality conditions for the optimization problem (25), it holds that

$$\max_{x_i} \sum_{i=1}^n x_i \log x_i \quad \text{such that} \quad \sum_{i=1}^n x_i = 1; \quad (25)$$

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  $\log n$ .

Thus, we have  $P_{a_h} \log a_h s_h \leq \log a_h s_h$  and  $\log a_h s_h$

$$E_{p^*} \log a_h s_h \leq \log a_h s_h$$

$$E_{s_0; a_0} \log a_h s_h \leq \log a_h s_h$$

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

$$E_{p^*} Y_2^2 \leq B \frac{2}{1-2} Q_{i=0}^{H^{\infty}} E_{p^*} \log a_h s_h$$

$$B \frac{\log a_h s_h}{1-2} Q_{i=0}^{H^{\infty}}$$

$$B \frac{\log a_h s_h}{1-2 \cdot 2};$$

for every  $H^{\infty} > 0$ . By taking expectation over the state action pair  $(s_H; a_H)$  and the horizon  $H^{\infty}$ , it yields that

$$E Y_2^2 \leq B \frac{\log a_h s_h}{1-2 \cdot 2};$$

which further implies that  $\text{Var} Y_2^2 \leq B \frac{\log a_h s_h}{1-2 \cdot 2}$ . This completes the proof.

C. Proof of Lemma 6

Proof. To simplify the notation, we denote  $r_j = a_j s_j$ . By definition, we have

$$E \sum_{h=0}^H \log a_h s_h \leq \sum_{j=1}^J r_j \log a_j s_j$$

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

$$[E \sum_{h=0}^H \log a_h s_h]^2 \leq \sum_{j=1}^J r_j \log a_j s_j$$

$$E \sum_{h=0}^H \log a_h s_h \leq \sum_{j=1}^J r_j \log a_j s_j$$

$$E \sum_{h=0}^H \log a_h s_h \leq \sum_{j=1}^J r_j \log a_j s_j \quad (26)$$

$$E \sum_{h=0}^H \log a_h s_h \leq \sum_{j=1}^J r_j \log a_j s_j \quad (27)$$

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

$$E \sum_{h=0}^H \log a_h s_h \leq \sum_{j=1}^J r_j \log a_j s_j$$

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 in (26) as follows:

$$E \sum_{h=0}^H \log a_h s_h \leq \sum_{j=1}^J r_j \log a_j s_j$$

Due to  $P_a \log a_h s_h \leq \log a_h s_h$  the term in (26) can be upper bounded as

$$E \sum_{h=0}^H \log a_h s_h \leq \sum_{j=1}^J r_j \log a_j s_j$$

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

$$E \sum_{h=0}^H \log a_h s_h \leq \sum_{j=1}^J r_j \log a_j s_j$$

Combining the above two inequalities, we have

$$[E \sum_{h=0}^H \log a_h s_h]^2 \leq \sum_{j=1}^J r_j \log a_j s_j$$

This completes the proof.

D. Proof of Lemma 7

Proof. For the simplicity of the notation, we first define:

$$g_1^H = \sum_{h=0}^H \log a_h s_h \quad (28)$$

$$g_2^H = \sum_{h=0}^H \log a_h s_h \quad (29)$$

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

$$\text{Var} \left[ \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i \right] \leq \sum_{h=0}^{H-1} \text{Var} \left[ \log \hat{a}_h^i \hat{s}_h^i \right] \quad (30)$$

$$\leq \sum_{h=0}^{H-1} \left( \text{E} \left[ \log \hat{a}_h^i \hat{s}_h^i \right]^2 - \left( \text{E} \left[ \log \hat{a}_h^i \hat{s}_h^i \right] \right)^2 \right) \quad (31)$$

$$\leq \sum_{h=0}^{H-1} \text{Var} \left[ \log \hat{a}_h^i \hat{s}_h^i \right] \quad (32)$$

As shown in Lemma 4.2 of [7], the fact that  $\log \hat{a}_h^i \hat{s}_h^i \in \mathcal{B}^2$  for all  $h \in \mathbb{R}^{\text{SSSS}}$  directly implies that  $\text{Var} \left[ \log \hat{a}_h^i \hat{s}_h^i \right] \leq B \frac{4r^2}{1-4}$  for all  $h \in \mathbb{R}^{\text{SSSS}}$ .

Then, it remains to prove the bounded variance of  $Y_2$ . Firstly, it can be observed that

$$Y_2 = \sum_{h=0}^{H-1} \sum_{j=0}^{h-1} \log \hat{a}_j^i \hat{s}_j^i - \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i$$

$$\leq \sum_{h=0}^{H-1} \sum_{j=0}^{h-1} \log \hat{a}_j^i \hat{s}_j^i - \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i$$

$$\leq \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i$$

where the first inequality is due to the triangle inequality and the second inequality is due to  $\log \hat{a}_j^i \hat{s}_j^i \in \mathcal{B}^2$ . Then, by taking the square of  $Y_2$ , we obtain

$$Y_2^2 \leq \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i$$

$$\leq \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i$$

$$\leq \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i$$

$$\leq \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i$$

$$\leq \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i$$

$$\leq \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i$$

$$\leq \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i$$

where the second inequality is due to the Cauchy-Schwarz inequality and the last inequality is due to  $\log \hat{a}_h^i \hat{s}_h^i \in \mathcal{B}^2$ ,  $P_{h,0}^a \leq h \frac{1}{1-2}$  and  $P_{h,0}^a \leq h \frac{1}{1}$ .

By taking expectation of  $Y_2$  over the sample trajectory, it holds that

$$\text{E} \left[ \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i \right] \leq \sum_{h=0}^{H-1} \text{E} \left[ \log \hat{a}_h^i \hat{s}_h^i \right] \quad (33)$$

Since the realizations of  $\hat{a}_h^i$  and  $\hat{s}_h^i$  do not depend on the randomness in  $\mathbb{S}_{h-1}, \mathbb{a}_{h-1}, \dots, \mathbb{S}_H$ , we have

$$\text{E} \left[ \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i \right] = \sum_{h=0}^{H-1} \text{E} \left[ \log \hat{a}_h^i \hat{s}_h^i \right]$$

$$= \sum_{h=0}^{H-1} \text{E} \left[ \log \hat{a}_h^i \hat{s}_h^i \right]$$

$$= \sum_{h=0}^{H-1} \text{E} \left[ \log \hat{a}_h^i \hat{s}_h^i \right]$$

$$= \sum_{h=0}^{H-1} \text{E} \left[ \log \hat{a}_h^i \hat{s}_h^i \right]$$

Since the maximizer for the constrained problem

$$\max_{x_1, \dots, x_n} \sum_{i=1}^n x_i \log x_i \quad \text{such that } \sum_{i=1}^n x_i = 1;$$

is  $x_1 = x_2 = \dots = x_n = \frac{1}{n}$  and the maximum solution is  $\log n^{-2}$ . Thus, we have  $\text{E} \left[ \log \hat{a}_h^i \hat{s}_h^i \right] \leq B \log \mathcal{A} S^2$  and

$$\text{E} \left[ \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i \right] \leq \sum_{h=0}^{H-1} B \log \mathcal{A} S^2$$

$$= B \log \mathcal{A} S^2 \quad (34)$$

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

$$\text{Var} \left[ \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i \right] \leq \sum_{h=0}^{H-1} \text{E} \left[ \log \hat{a}_h^i \hat{s}_h^i \right]^2$$

$$\leq \sum_{h=0}^{H-1} B^2 \log^2 \mathcal{A} S^2$$

$$\leq \sum_{h=0}^{H-1} B^2 \log^2 \mathcal{A} S^2$$

$$\leq \sum_{h=0}^{H-1} B^2 \log^2 \mathcal{A} S^2$$

Finally, by substituting  $\text{Var} \left[ \log \hat{a}_h^i \hat{s}_h^i \right]$ ,  $\text{Var} \left[ \log \hat{a}_h^i \hat{s}_h^i \right]$  and  $\text{Var} \left[ \log \hat{a}_h^i \hat{s}_h^i \right]$  into (30), it holds that

$$\text{Var} \left[ \sum_{h=0}^{H-1} \log \hat{a}_h^i \hat{s}_h^i \right] \leq B \frac{12r^2 + 24 \log^2 \mathcal{A} S^2}{1-4}$$

This completes the proof.  $\square$

## APPENDIX B

### HELPFUL RESULTS FOR THE PROOF OF LEMMA 10

Proposition 1 (Proposition 3 in [8]): Consider the problem  $\max_{x \in \mathbb{R}^d} f^*(x)$ , where  $\mathbb{R}^d$  denotes the  $d$ -dimensional Euclidean space. Let  $x_t^a$  be a sequence generated by the iterative method  $x_{t+1} = x_t + \eta_t u_t + w_t$ , where  $\eta_t$  is a deterministic positive step-size,  $u_t$  is an update direction, and  $w_t$  is a random noise term. Let  $F_t$  be an increasing sequence of  $\sigma$ -fields. Assume that

- 1)  $f^*$  is a continuously differentiable function and there exists a constant  $L$  such that

$$\| \nabla f^*(x) - \nabla f^*(y) \| \leq L \| x - y \|; \quad \forall x, y \in \mathbb{R}^d;$$

- 2)  $x_t$  and  $u_t$  are  $F_t$  measurable.

