
An Improved Analysis of (Variance-Reduced) Policy Gradient and Natural Policy Gradient Methods

Yanli Liu[‡] Kaiqing Zhang[†] Tamer Başar[†] Wotao Yin[‡]

[‡]Department of Mathematics, University of California, Los Angeles

[†]Department of ECE and CSL, University of Illinois at Urbana-Champaign
{yanli, wotaoyin}@math.ucla.edu, {kzhang66, basar1}@illinois.edu

Abstract

In this paper, we revisit and improve the convergence of policy gradient (PG), natural PG (NPG) methods, and their variance-reduced variants, under general smooth policy parametrizations. More specifically, with the Fisher information matrix of the policy being positive definite: i) we show that a state-of-the-art variance-reduced PG method, which has only been shown to converge to stationary points, converges to the globally optimal value up to some inherent function approximation error due to policy parametrization; ii) we show that NPG enjoys a lower sample complexity; iii) we propose SRVR-NPG, which incorporates variance-reduction into the NPG update. Our improvements follow from an observation that the convergence of (variance-reduced) PG and NPG methods can improve each other: the stationary convergence analysis of PG can be applied to NPG as well, and the global convergence analysis of NPG can help to establish the global convergence of (variance-reduced) PG methods. Our analysis carefully integrates the advantages of these two lines of works. Thanks to this improvement, we have also made variance-reduction for NPG possible, with both global convergence and an efficient finite-sample complexity.

1 Introduction

Policy gradient (PG) methods, or more generally direct policy search methods, have long been recognized as one of the foundations of reinforcement learning (RL) [1]. Specifically, PG methods directly search for the optimal policy parameter that maximizes the long-term return in Markov decision processes (MDPs), following the policy gradient ascent direction [2, 3]. This search direction can be more efficient using a preconditioning matrix, e.g., using the natural PG direction [4]. These methods have achieved tremendous empirical successes recently, especially boosted by the power of (deep) neural networks for policy parametrization [5, 6, 7, 8]. These successes are primarily attributed to the fact that PG methods naturally incorporate *function approximation* for policy parametrization, in order to handle massive and even continuous state-action spaces.

In practice, the policy gradients are usually estimated via samples using Monte-Carlo rollouts and bootstrapping [2, 9]. Such stochastic PG methods notoriously suffer from very high variances, which not only destabilize but also slow down the convergence. Several conventional approaches have been advocated to reduce the variance of PG methods, e.g., by adding a baseline [3, 10], or by using function approximation for estimating the value function, namely, developing actor-critic algorithms [11, 12, 13]. More recently, motivated by the advances of variance-reduction techniques in stochastic optimization [14, 15, 16, 17], there have been surging interests in developing *variance-reduced* PG methods [18, 19, 20, 21, 22], which are shown to be faster.

In contrast to the empirical successes of PG methods, their theoretical convergence guarantees, especially *non-asymptotic global* convergence guarantees, have not been addressed satisfactorily

until very recently [23, 24, 25, 26, 27]. By *non-asymptotic global* convergence, here we mean the convergence behavior of PG methods from any initialization, and the quality of the point they converge to (usually enjoys global optimality up to some compatible function approximation error due to policy parametrization), after a finite number of iterations/samples. These recent prominent guarantees are normally beyond the folklore *first-order* stationary-point convergence¹, as expected from a *stochastic nonconvex optimization* perspective of solving RL with PG methods. Special landscapes of the RL objective, though nonconvex, have enabled the convergence to even globally optimal values. On the other hand, none of the aforementioned variance-reduced PG methods [18, 19, 20, 21, 22] have been shown to enjoy these desired global convergence properties. It remains unclear whether these methods can converge to beyond first-order stationary policies.

Motivated by these advances and the questions that remain to be answered, we aim in this paper to improve the convergence of PG and natural PG (NPG) methods, and their variance-reduced variants, under general smooth policy parametrizations. Our contributions are summarized as follows.

Contributions. With a focus on the conventional Monte-Carlo-based PG methods, we propose a general framework for analyzing their *global convergence*. Our contribution is three-fold: first, we establish the global convergence up to compatible function approximation errors due to policy parametrization, for a variance-reduced PG method SRVR-PG [21]; second, we improve the global convergence of NPG methods established in [27], from $\mathcal{O}(\varepsilon^{-4})$ to $\mathcal{O}(\varepsilon^{-3})$; third, we propose a new variance-reduced algorithm based on NPG, and establish its global convergence with an efficient sample-complexity. These improvements are based on a framework that integrates the advantages of previous analyses on (variance reduced) PG and NPG, and rely on a (mild) assumption that the Fisher information matrix induced by the policy parametrization is positive definite (see Assumption 2.1). A comparison of previous results and our improvements is laid out in Table 1.

Related Work.

Global Convergence of (Natural) PG. Recently, there has been a surging research interest in investigating the global convergence of PG and NPG methods, which is beyond the folklore convergence to first-order stationary policies. In the special case with linear dynamics and quadratic reward, [23] shows that PG methods with random search converge to the globally optimal policy with linear rates. In [24], with a simple reward-reshaping, PG methods have been shown to converge to the second-order stationary-point policies. [26] shows that for finite-MDPs and several control tasks, the nonconvex RL objective has no suboptimal local minima. [25] prove that (natural) PG methods converge to the globally optimal value when overparametrized neural networks are used for function approximation. [27] provides a fairly general characterization of global convergence for these methods, and a basic sample complexity result for sample-based NPG updates. It is also worth noting that trust-region policy optimization (TRPO) [5], as a variant of NPG, also enjoys global convergence with overparametrized neural networks [28], and for regularized MDPs [29]. Very recently, for actor-critic algorithms, a series of non-asymptotic convergence results have also been established [30, 31, 32, 33], with global convergence guarantees when natural PG/PPO are used in the actor step.

Variance-Reduction (VR) for PG. Conventional approaches to reduce the high variance in PG methods include using (natural) actor-critic algorithms [11, 12, 13], and adding baselines [3, 10]. The idea of variance reduction (VR) is first proposed to accelerate stochastic minimization. VR algorithms such as SVRG [14, 15, 16], SAGA [17], SARAH [34], and Spider [35] achieve acceleration over SGD in both convex and nonconvex settings. SVRG is also accelerated by applying a positive definite preconditioner that captures the curvature of the objective [36]. Inspired by these successes in stochastic optimization, VR is also incorporated into PG methods [18], with empirical validations for acceleration, and analyzed rigorously in [19]. Then, [20] improves the sample complexity of SVRPG, and [21] proposes a new SRVR-PG method that uses recursively updated semi-stochastic policy gradient, which leads to an improved sample complexity of $\mathcal{O}(\varepsilon^{-1.5})$ over previous works. More recently, [22] proposes a new STORM-PG method, which blends momentum in the update and matches the sample complexity of in [21], and [37] applies the idea of SARAH and considers a more general setting with regularization. Finally, heavy-ball type of momentum has also been applied to PG methods [38]. We highlight that all these sample complexity results are for first-order stationary-point convergence (which might have arbitrarily bad performance: see (2.2)), in contrast

¹That is, finding a parameter θ such that $\|\nabla J(\theta)\|^2 \leq \varepsilon$, where J is the expected return.

to the more desired global convergence guarantees (up to some function approximation errors that can be small) that we are interested in.

NPG [27]	NPG [25]	TRPO [28]	TRPO [29]
$\mathcal{O}(\varepsilon^{-4})$	$\mathcal{O}(T_{TD}\varepsilon^{-2})$ ¹	$\mathcal{O}(\varepsilon^{-8})$	$\mathcal{O}(\varepsilon^{-4})$

NPG (2.8)	PG (2.4)	SRVR-PG (Algorithm 2)	SRVR-NPG (Algorithm 1)
$\mathcal{O}(\varepsilon^{-3})$	$\mathcal{O}(\sigma^2\varepsilon^{-4})$	$\mathcal{O}((W + \sigma^2)\varepsilon^{-3})$	$\mathcal{O}((W + \sigma^2)\varepsilon^{-2.5} + \varepsilon^{-3})$

Table 1: Comparison of sample complexities of several methods to reach global optimality up to some compatible function approximation error (see (2.9)). Our results are listed in the second table (See App. A for their derivations). We compare the number of trajectories to reach ε -optimality in expectation, up to some inherent error due to the function approximation for policy parametrization (see (2.3)). σ^2 is an upper bound for the variance of gradient estimator (see Assumption 4.1), and W is an upper bound for the variance of importance weight (see Assumption 4.3).

2 Preliminaries

We first introduce some preliminaries regarding both the MDPs and policy gradient methods.

2.1 Markov Decision Processes

Consider a discounted Markov decision process defined by a tuple $(\mathcal{S}, \mathcal{A}, \mathbb{P}, R, \gamma)$, where \mathcal{S} and \mathcal{A} denote the state and action spaces of the agent, $\mathbb{P}(s' | s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$ is the Markov kernel that determines the transition probability from (s, a) to state s' , $\gamma \in (0, 1)$ is the discount factor, and $r : \mathcal{S} \times \mathcal{A} \rightarrow [-R, R]$ is the reward function of s and a .

At each time t , the agent executes an action $a_t \in \mathcal{A}$ given the current state $s_t \in \mathcal{S}$, following a possibly stochastic policy $\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$, i.e., $a_t \sim \pi(\cdot | s_t)$. Then, given the state-action pair (s_t, a_t) , the agent observes a reward $r_t = r(s_t, a_t)$. Thus, under any policy π , one can define the *state-action value* function $Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ as

$$Q^\pi(s, a) := \mathbb{E}_{a_t \sim \pi(\cdot | s_t), s_{t+1} \sim \mathbb{P}(\cdot | s_t, a_t)} \left(\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a \right).$$

One can also define the *state-value* function $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$, and the *advantage* function $A^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, under policy π , as $V^\pi(s) := \mathbb{E}_{a \sim \pi(\cdot | s)} [Q^\pi(s, a)]$ and $A^\pi(s, a) := Q^\pi(s, a) - V^\pi(s)$, respectively. Suppose that the initial state s_0 is drawn from some distribution ρ . Then, the goal of the agent is to find the optimal policy that maximizes the expected discounted return, namely,

$$\max_{\pi} J(\pi) := \mathbb{E}_{s_0 \sim \rho} [V^\pi(s_0)]. \quad (2.1)$$

In practice, both the state and action spaces \mathcal{S} and \mathcal{A} can be very large. Thus, the policy π is usually parametrized as π_θ for some parameter $\theta \in \mathbb{R}^d$, using, for example, deep neural networks. As such, the goal of the agent is to maximize $J(\pi_\theta)$ in the space of the parameter θ , which naturally induces an optimization problem. Such a problem is in general nonconvex [24, 27], making it challenging to find the globally optimal policy.

For notational convenience, let us denote $J(\pi_\theta)$ by $J(\theta)$. Many of the previous works focus on establishing stationary convergence of policy gradient methods. That is, finding a θ that satisfies

$$\|\nabla J(\theta)\|^2 \leq \varepsilon. \quad (2.2)$$

Obviously, such a θ may not lead to a large $J(\theta)$. Instead, we are interested in finding a θ such that

$$J^* - J(\theta) \leq \mathcal{O}(\sqrt{\varepsilon_{\text{bias}}}) + \varepsilon, \quad (2.3)$$

where $J^* = \max_{\pi} J(\pi)$, and the $\mathcal{O}(\sqrt{\varepsilon_{\text{bias}}})$ term reflects the inherent error related to the possibly limited expressive power of the policy parametrization π_θ (see Assumption 4.4 for the definition).

¹In [25], T_{TD} iterations of temporal difference updates are needed at each iteration, T_{TD} can be large for wide neural networks. See App. A for details.

2.2 (Natural) Policy Gradient Methods

To solve the optimization problem (2.1), one standard way is via the policy gradient (PG) method [3]. Specifically, let $\tau_i = \{s_0^i, a_0^i, s_1^i, \dots\}$ denote the data of a sampled trajectory under policy π_θ . Then, a stochastic PG ascent update is given as

$$\theta^{k+1} = \theta^k + \eta \cdot \frac{1}{N} \sum_{i=1}^N g(\tau_i | \theta^k), \quad (2.4)$$

where $\eta > 0$ is a stepsize, N is the number of trajectories, and $g(\tau_i | \theta^k)$ estimates $\nabla J(\theta^k)$ using the trajectory τ_i . Common unbiased estimators of PG include REINFORCE [2], using the policy gradient theorem [39], and GPOMDP [9]. The commonly used GPOMDP estimator will be given by

$$g(\tau_i | \theta) = \sum_{h=0}^{\infty} \left(\sum_{t=0}^h \nabla_\theta \log \pi_\theta(a_t^i | s_t^i) \right) (\gamma^h r(s_h^i, a_h^i)), \quad (2.5)$$

where $\nabla_\theta \log \pi_\theta(a_t^i | s_t^i)$ is the *score function*. If the expectation of this infinite sum exists, then (2.5) becomes an unbiased estimate of the policy gradient of the objective $J(\theta)$ defined in (2.1). This unbiasedness is established in App. B for completeness.

In practice, a *truncated* version of GPOMDP is used to approximate the infinite sum in (2.5), as

$$g(\tau_i^H | \theta) = \sum_{h=0}^{H-1} \left(\sum_{t=0}^h \nabla_\theta \log \pi_\theta(a_t^i | s_t^i) \right) (\gamma^h r(s_h^i, a_h^i)), \quad (2.6)$$

where $\tau_i^H = \{s_0^i, a_0^i, s_1^i, \dots, s_{H-1}^i, a_{H-1}^i, s_H^i\}$ is a truncation of the full trajectory τ_i of length H . (2.6) is thus a biased stochastic estimate of $\nabla J(\theta)$, with the bias being negligible for a large enough H . For notational simplicity, we denote the H -horizon trajectory distribution induced by the initial state distribution ρ and policy π_θ as $p_\rho^H(\cdot | \theta)$, that is,

$$p_\rho^H(\tau^H | \theta) = \rho(s_0) \prod_{h=0}^{H-1} \pi_\theta(a_h | s_h) \mathbb{P}(s_{h+1} | a_h, s_h).$$

Hereafter, unless otherwise stated, we refer to this H -horizon trajectory simply as *trajectory*, drawn from $p_\rho^H(\cdot | \theta)$.

As a significant variant of PG, NPG [4] also incorporates a preconditioning matrix $F_\rho(\theta)$, leading to the following update

$$F_\rho(\theta) = \mathbb{E}_{s \sim d_\rho^{\pi_\theta}} [F_s(\theta)], \quad \theta^{k+1} = \theta^k + \eta \cdot F_\rho^\dagger(\theta^k) \nabla J(\theta^k), \quad (2.7)$$

where $F_s(\theta) = \mathbb{E}_{a \sim \pi_\theta(\cdot | s)} [\nabla_\theta \log \pi_\theta(a | s) \nabla_\theta \log \pi_\theta(a | s)^\top]$ is the Fisher information matrix of $\pi_\theta(\cdot | s) \in \mathcal{P}(\mathcal{A})$, $F_\rho^\dagger(\theta^k)$ is the Moore-Penrose pseudoinverse of $F_\rho(\theta^k)$, and $d_\rho^{\pi_\theta} \in \mathcal{P}(\mathcal{S})$ is the state visitation measure induced by policy π_θ and initial distribution ρ , which is defined as

$$d_\rho^{\pi_\theta}(s) := (1 - \gamma) \mathbb{E}_{s_0 \sim \rho} \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s | s_0, \pi_\theta).$$

The NPG update (2.7) can also be written as [4, 27]

$$\theta^{k+1} = \theta^k + \eta \cdot w^k, \quad \text{with } w^k \in \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} L_{\nu_\rho^{\pi_\theta}}(w; \theta), \quad (2.8)$$

where $L_{\nu_\rho^{\pi_\theta}}(w; \theta)$ is the *compatible function approximation error* defined by

$$L_{\nu_\rho^{\pi_\theta}}(w; \theta) = \mathbb{E}_{(s,a) \sim \nu_\rho^{\pi_\theta}} \left[(A^{\pi_\theta}(s, a) - (1 - \gamma) w^\top \nabla_\theta \log \pi_\theta(a | s))^2 \right]. \quad (2.9)$$

Here, $\nu_\rho^{\pi_\theta}(s, a) = d_\rho^{\pi_\theta}(s) \pi(a | s)$ is the *state-action* visitation measure induced by π_θ and initial state distribution ρ , which can also be written as

$$\nu_\rho^{\pi_\theta}(s, a) := (1 - \gamma) \mathbb{E}_{s_0 \sim \rho} \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s, a_t = a | s_0, \pi_\theta). \quad (2.10)$$

For convenience, we will denote $\nu_\rho^{\pi_\theta}$ by ν^{π_θ} hereafter. In other words, the NPG update direction w^k is given by the minimizer of a stochastic optimization problem. In practice, one obtains an approximate NPG update direction w^k by SGD (see Procedure 1).

Regarding the NPG update (2.8), we make the following standing assumption on the Fisher information matrix induced by π_θ and ρ .

Assumption 2.1. For all $\theta \in \mathbb{R}^d$, the Fisher information matrix induced by policy π_θ and initial state distribution ρ satisfies

$$F_\rho(\theta) = \mathbb{E}_{(s,a) \sim \nu_\rho^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(a | s) \nabla_\theta \log \pi_\theta(a | s)^\top] \succcurlyeq \mu_F \cdot I_d$$

for some constant $\mu_F > 0$.

Assumption 2.1 essentially states that $F_\rho(\theta)$ behaves well as a preconditioner in the NPG update (2.8). This is a common (and minimal) requirement for the convergence of preconditioned algorithms in both convex and nonconvex settings in the optimization realm, for example, the quasi-Newton algorithms [40, 41, 42, 43], and their stochastic variants [44, 45, 46, 47, 36]. In the RL realm, one common example of policy parametrizations that can satisfy this assumption is the Gaussian policy [2, 48, 19, 21], where $\pi_\theta(\cdot | s) = \mathcal{N}(\mu_\theta(s), \Sigma)$ with mean parametrized linearly as $\mu_\theta(s) = \phi(s)^\top \theta$, where $\phi(s)$ denotes some feature matrix of proper dimensions, θ is the coefficient vector, and $\Sigma \succ 0$ is some fixed covariance matrix. In this case, the Fisher information matrix at each s becomes $\phi(s) \Sigma^{-1} \phi(s)^\top$, independent of θ , and is uniformly lower bounded (positive definite sense) if $\phi(s)$ is full-row-rank, namely, the features expanded by θ are linearly independent, which is a common requirement for linear function approximation settings [49, 50, 51]. See App. B.2 for more detailed justifications, as well as discussions on more general policy parametrizations.

In the pioneering NPG work [4], $F(\theta)$ is directly assumed to be positive definite. So is in the follow-up works on natural actor-critic algorithms [12, 13]. In fact, this way, $F(\theta)$ will define a valid Riemannian metric on the parameter space, which has been used for interpreting the desired convergence properties of natural gradient methods [52, 53]. In a recent version of [27], a relevant assumption (specifically, Assumption 6.5, item 3) is made to establish the global convergence of NPG, in which it is assumed that $\lambda_{\min}(F_\rho(\theta))$ is not too small compared with the Fisher information matrix induced by a fixed comparator policy. This can be implied by our Assumption 2.1. To sum up, the positive definiteness on the Fisher preconditioning matrix is common and not very restrictive.

In Sec. 4, we shall see that under Assumption 2.1, the stationary convergence of NPG can be analyzed, and NPG enjoys a better sample complexity of $\mathcal{O}(\varepsilon^{-3})$ in terms of its global convergence, compared with the existing sample complexity of $\mathcal{O}(\varepsilon^{-4})$ in [27]. In addition, interestingly, PG and its variance-reduced version SRVR-PG also enjoy global convergence, although the Fisher information matrix does not appear explicitly in their updates.

3 Variance-Reduced Policy Gradient Methods

Recently, [21] proposes an algorithm called Stochastic Recursive Variance Reduced Policy Gradient (SRVR-PG, see Algorithm 2), which applies variance-reduction on PG. It achieves a sample complexity of $\mathcal{O}(\varepsilon^{-1.5})$ to find an ε -stationary point, compared with the $\mathcal{O}(\varepsilon^{-2})$ sample complexity of stochastic PG. However, it remains unclear whether SRVR-PG converges globally. In this work, we provide an affirmative answer to this question by showing that SRVR-PG has a sample complexity of $\mathcal{O}(\varepsilon^{-3})$ to find an ε -optimal policy, up to some compatible function approximation error due to policy parametrization.

We also propose a new algorithm called SRVR-NPG to incorporate variance reduction into NPG, which is described in Algorithm 1. In Sec. 4, we provide a sample complexity for its global convergence, which is comparable to our improved NPG result.

In line 8 of Algorithm 1, $g_w(\tau_j^H | \theta_{t-1}^{j+1})$ is a weighted gradient estimator given by

$$g_w(\tau_j^H | \theta_{t-1}^{j+1}) = \sum_{h=0}^{H-1} w_{0:h}(\tau_j^H | \theta_{t-1}^{j+1}, \theta_t^{j+1}) \left(\sum_{i=0}^h \nabla_\theta \log \pi_\theta(a_t^i | s_t^i) \right) (\gamma^h r(s_h^i, a_h^i)), \quad (3.1)$$

Algorithm 1 Stochastic Recursive Variance Reduced Natural Policy Gradient (SRVR-NPG)

Input: number of epochs S , epoch size m , stepsize η , batch size N , minibatch size B , truncation horizon H , initial parameter $\theta_m^0 = \theta_0 \in \mathbb{R}^d$.

```
1: for  $j \leftarrow 0, \dots, S - 1$  do
2:    $\theta_0^{j+1} = \theta_m^j$ ;
3:   Sample  $\{\tau_i^H\}_{i=1}^N$  from  $p_\rho^H(\cdot | \theta_0^{j+1})$  and calculate  $u_0^{j+1} = \frac{1}{N} \sum_{i=1}^N g(\tau_i^H | \theta_0^{j+1})$ ;
4:    $w_0^{j+1} = \text{SRVR-NPG-SGD}(\nu^{\pi_{\theta_0^{j+1}}}, \pi_{\theta_0^{j+1}}, u_0^{j+1})$ ;  $\triangleright w_0^{j+1} \approx w_{0,*}^{j+1} = F_\rho^{-1}(\theta_0^{j+1})u_0^{j+1}$ ;
5:    $\theta_1^{j+1} = \theta_0^{j+1} + \eta w_0^{j+1}$ ;
6:   for  $t \leftarrow 1, \dots, m - 1$  do
7:     Sample  $B$  trajectories  $\{\tau_j^H\}_{j=1}^B$  from  $p_\rho^H(\cdot | \theta_t^{j+1})$ ;
8:      $u_t^{j+1} = u_{t-1}^{j+1} + \frac{1}{B} \sum_{j=1}^B (g(\tau_j^H | \theta_t^{j+1}) - g_w(\tau_j^H | \theta_{t-1}^{j+1}))$ ;
9:      $w_t^{j+1} = \text{SRVR-NPG-SGD}(\nu^{\pi_{\theta_t^{j+1}}}, \pi_{\theta_t^{j+1}}, u_t^{j+1})$ ;  $\triangleright w_t^{j+1} \approx w_{t,*}^{j+1} = F_\rho^{-1}(\theta_t^{j+1})u_t^{j+1}$ ;
10:     $\theta_{t+1}^{j+1} = \theta_t^{j+1} + \eta w_t^{j+1}$ ;
11:  end for
12: end for
13: return  $\theta_{\text{out}}$  chosen uniformly from  $\{\theta\}_{j=1, \dots, S; t=0, \dots, m-1}$ .
```

where the importance weight factor $w_{0:h}(\tau_j^H | \theta_{t-1}^{j+1}, \theta_t^{j+1})$ is defined by

$$w_{0:h}(\tau_j^H | \theta_{t-1}^{j+1}, \theta_t^{j+1}) = \prod_{h'=0}^h \frac{\pi_{\theta_{t-1}^{j+1}}(a_{h'} | s_{h'})}{\pi_{\theta_t^{j+1}}(a_{h'} | s_{h'})}. \quad (3.2)$$

This importance sampling makes w_t^{j+1} an unbiased estimator of $\nabla J^H(\theta_t^{j+1})$.

In lines 4 and 8 of Algorithm 1, w_t^{j+1} is produced by SRVR-NPG-SGD (see Procedure 2), which applies SGD¹ to solve the following subproblem:

$$w_t^{j+1} \approx \underset{w}{\operatorname{argmin}} \left\{ \mathbb{E}_{(s,a) \sim \nu_t^{j+1}} \left[\left(w^T \nabla_\theta \log \pi_{\theta_t^{j+1}}(a | s) \right)^2 \right] - 2 \langle w, u_t^{j+1} \rangle \right\}, \quad (3.3)$$

where ν_t^{j+1} is the state-action visitation measure induced by $\pi_{\theta_t^{j+1}}$. The exact update direction given by (3.3) is $F_\rho^{-1}(\theta_t^{j+1})u_t^{j+1}$, and as in NPG, $F_\rho(\theta_t^{j+1})$ also serves as a preconditioner.

4 Theoretical Results

Before presenting the global convergence results, we first introduce some standard assumptions.

Assumption 4.1. The truncated GPOMDP estimator $g(\tau^H | \theta)$ defined in (2.6) satisfies $\operatorname{Var}(g(\tau^H | \theta)) := \mathbb{E}[\|g(\tau^H | \theta) - \mathbb{E}[g(\tau^H | \theta)]\|^2] \leq \sigma^2$ for any θ and $\tau^H \sim p_\rho^H(\cdot | \theta)$.

Assumption 4.2. 1. $\|\nabla_\theta \log \pi_\theta(a | s)\| \leq G$ for any θ and $(s, a) \in \mathcal{S} \times \mathcal{A}$.

2. $\|\nabla_\theta \log \pi_{\theta_1}(a | s) - \nabla_\theta \log \pi_{\theta_2}(a | s)\| \leq M\|\theta_1 - \theta_2\|$ for any θ_1, θ_2 and $(s, a) \in \mathcal{S} \times \mathcal{A}$.

Assumption 4.3. For the importance weight $w_{0:h}(\tau^H | \theta_1, \theta_2)$ (3.2), there exists $W > 0$ such that

$$\operatorname{Var}(w_{0:h}(\tau^H | \theta_1, \theta_2)) \leq W, \quad \forall \theta_1, \theta_2 \in \mathbb{R}^d, \tau^H \sim p_\rho^H(\cdot | \theta_2).$$

Assumptions 4.1, 4.2 and 4.3 are standard in the analysis of PG methods and their variance reduced variants [27, 19, 20, 21]. They can be verified for simple policy parametrizations such as Gaussian policies; see [19, 57, 58] for more justifications.

Following the Assumption 6.5 of [27], we assume that the policy parametrization π_θ achieves a good function approximation, as measured by the *transferred compatible function approximation error*.

¹Following [27], we apply SGD [54] to make a fair comparison. One can also apply the SA algorithm [55] and AC-SA algorithm [56].

Assumption 4.4. For any $\theta \in \mathbb{R}^d$, the *transferred compatible function approximation error* satisfies

$$L_{\nu^*}(w_*^\theta; \theta) = \mathbb{E}_{(s,a) \sim \nu^*} \left[\left(A^{\pi_\theta}(s, a) - (1 - \gamma)(w_*^\theta)^\top \nabla_\theta \log \pi_\theta(a | s) \right)^2 \right] \leq \varepsilon_{\text{bias}}, \quad (4.1)$$

where $\nu^*(s, a) = d_\rho^{\pi^*}(s) \cdot \pi^*(a | s)$ is the state-action distribution induced by an optimal policy π^* that maximizes $J(\pi)$, and $w_*^\theta = \operatorname{argmin}_{w \in \mathbb{R}^d} L_{\nu_\rho^{\pi_\theta}}(w; \theta)$ is the exact NPG update direction at θ .

$\varepsilon_{\text{bias}}$ reflects the error when approximating the advantage function from the score function, it measures the capacity of the parametrization π_θ . When π_θ is the softmax parametrization, we have $\varepsilon_{\text{bias}} = 0$ [27]. When π_θ is a restricted parametrization, $\varepsilon_{\text{bias}}$ is often positive as π_θ may not contain all stochastic policies. For rich neural parametrizations, $\varepsilon_{\text{bias}}$ is very small [25].

4.1 A General Framework for Global Convergence

Inspired by the global convergence analysis of NPG in [27], we present a general framework that relates the global convergence rates of these algorithms to i) their stationary convergence rate on $J(\theta)$, and ii) the difference between their update directions and exact NPG update directions.

Proposition 4.5. Let $\{\theta^k\}_{k=1}^K$ be generated by a general update of the form

$$\theta^{k+1} = \theta^k + \eta w^k, \quad k = 0, 1, \dots, K-1.$$

Furthermore, let $w_*^k = F_\rho^{-1}(\theta^k) \nabla J(\theta^k)$ be the exact NPG update direction at θ^k . Then, we have

$$\begin{aligned} J(\pi^*) - \frac{1}{K} \sum_{k=0}^{K-1} J(\theta^k) &\leq \frac{\sqrt{\varepsilon_{\text{bias}}}}{1 - \gamma} + \frac{1}{\eta K} \mathbb{E}_{s \sim d_\rho^{\pi^*}} [\text{KL}(\pi^*(\cdot | s) \| \pi_{\theta^0}(\cdot | s))] \\ &\quad + \frac{M\eta}{2K} \sum_{k=0}^{K-1} \|w^k\|^2 + \frac{G}{K} \sum_{k=0}^{K-1} \|w^k - w_*^k\|, \end{aligned} \quad (4.2)$$

where π^* is an optimal policy that maximizes $J(\pi)$.

The detailed proof of this global convergence framework can be found in J. To obtain a high level idea, one first starts from the M -smoothness of the score function to get

$$\begin{aligned} &\mathbb{E}_{s \sim d_\rho^{\pi^*}} [\text{KL}(\pi^*(\cdot | s) \| \pi_{\theta^k}(\cdot | s)) - \text{KL}(\pi^*(\cdot | s) \| \pi_{\theta^{k+1}}(\cdot | s))] \\ &\geq \eta \mathbb{E}_{s \sim d_\rho^{\pi^*}} \mathbb{E}_{a \sim \pi^*(\cdot | s)} [\nabla_\theta \log \pi_{\theta^k}(a | s) \cdot w_*^k] \\ &\quad + \eta \mathbb{E}_{s \sim d_\rho^{\pi^*}} \mathbb{E}_{a \sim \pi^*(\cdot | s)} [\nabla_\theta \log \pi_{\theta^k}(a | s) \cdot (w^k - w_*^k)] - \frac{M\eta^2}{2} \|w^k\|^2. \end{aligned}$$

On the other hand, the renowned Performance Difference Lemma [59] tells us that

$$\mathbb{E}_{s \sim d_\rho^{\pi^*}} \mathbb{E}_{a \sim \pi^*(\cdot | s)} [A^{\pi_{\theta^k}}(s, a)] = (1 - \gamma) (J^* - J(\theta^k)).$$

To connect the advantage term $\mathbb{E}_{s \sim d_\rho^{\pi^*}} \mathbb{E}_{a \sim \pi^*(\cdot | s)} [A^{\pi_{\theta^k}}(s, a)]$ with the inner product term $\mathbb{E}_{s \sim d_\rho^{\pi^*}} \mathbb{E}_{a \sim \pi^*(\cdot | s)} [\nabla_\theta \log \pi_{\theta^k}(a | s) \cdot w_*^k]$, we invoke Assumption 4.4:

$$\mathbb{E}_{s \sim d_\rho^{\pi^*}} \mathbb{E}_{a \sim \pi^*(\cdot | s)} \left[\left(A^{\pi_\theta}(s, a) - (1 - \gamma)(w_*^\theta)^\top \nabla_\theta \log \pi_\theta(a | s) \right)^2 \right] \leq \varepsilon_{\text{bias}}, \quad \text{for any } \theta \in \mathbb{R}^d.$$

The final result follows from a telescoping sum on $k = 0, 1, \dots, K-1$.

Several remarks are in order. The first term on the right-hand side of (4.2) reflects the function approximation error due to the parametrization π_θ , and the second term is of the form $\mathcal{O}(\frac{1}{K})$. The third term depends on the stationary convergence. With Assumption 2.1, it can be shown that¹ $\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|w^k\|^2] \rightarrow 0$ for both NPG and SRVR-NPG. The proof follows from an optimization perspective and is inspired by the stationary convergence analysis of stochastic PG (see App. E).

With Assumption 2.1, we can also show that the last term of (4.2) is small. Take stochastic PG as an example; then, we have $w^k = \frac{1}{N} \sum_{i=1}^N g(\tau_i^H | \theta^k)$, and

$$\frac{1}{K} \sum_{k=0}^{K-1} \|w^k - w_*^k\| \leq \frac{1}{K} \sum_{k=0}^{K-1} \|w^k - \nabla J(\theta^k)\| + \frac{1}{K} \sum_{k=0}^{K-1} \left(1 + \frac{1}{\mu_F} \right) \|\nabla J(\theta^k)\|.$$

¹The stationary convergence of SRVR-PG has been established in [21].

When H and N are large enough, w^k is a low-variance estimator of $\nabla J^H(\theta^k)$, and $\nabla J^H(\theta^k)$ is close to $\nabla J(\theta^k)$, this makes the first term above small. The second term also goes to 0 as θ^k approaches stationarity.

4.2 Global Convergence Results

By applying Proposition 4.5 on the PG, NPG, SRVR-PG, and SRVR-NPG updates and analyzing their stationary convergence, we obtain their global convergence rates. In the following, we only keep the dependences on σ^2 (the variance of the gradient estimator), W (variance of importance weight), $\frac{1}{1-\gamma}$ (the effective horizon) and ε (target accuracy). The specific choice of the parameters and sample complexities, as well as the proof, can be found in the appendix.

Theorem 4.6. In the stochastic PG (2.4) with the truncated GPOMDP estimator (2.6), take $\eta = \frac{1}{4L_J}$, $K = \mathcal{O}\left(\frac{1}{(1-\gamma)^2\varepsilon^2}\right)$, $N = \mathcal{O}\left(\frac{\sigma^2}{\varepsilon^2}\right)$, and $H = \mathcal{O}\left(\log\left(\frac{1}{(1-\gamma)\varepsilon}\right)\right)$. Then, we have

$$J(\pi^*) - \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[J(\theta^k)] \leq \frac{\sqrt{\varepsilon_{\text{bias}}}}{1-\gamma} + \varepsilon.$$

In total, stochastic PG samples $\mathcal{O}\left(\frac{\sigma^2}{(1-\gamma)^2\varepsilon^4}\right)$ trajectories.

Remark 4.7. $L_J = \frac{MR}{(1-\gamma)^2}$ is the Lipschitz constant of ∇J , see Lemma B.1 for details.

Remark 4.8. Theorem 4.6 improves the result of [27, Thm. 6.11] from (impractical) full gradients to sample-based stochastic gradients.

Theorem 4.9. In the NPG update (2.8), let us apply $\mathcal{O}\left(\frac{1}{(1-\gamma)^4\varepsilon^2}\right)$ iterations of SGD as in Procedure 1 to obtain an update direction. In addition, take $\eta = \frac{\mu_F^2}{4G^2L_J}$ and $K = \mathcal{O}\left(\frac{1}{(1-\gamma)^2\varepsilon}\right)$. Then,

$$J^* - \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[J(\theta^k)] \leq \frac{\sqrt{\varepsilon_{\text{bias}}}}{1-\gamma} + \varepsilon.$$

In total, NPG samples $\mathcal{O}\left(\frac{1}{(1-\gamma)^6\varepsilon^3}\right)$ trajectories.

Remark 4.10. Compared with [27, Coro. 6.10], Theorem 4.9 improves the sample complexity of NPG by $\mathcal{O}(\varepsilon^{-1})$. This is because our stationary convergence analysis on NPG allows for a constant stepsize η , while [27, Coro. 6.10] applies a stepsize of $\eta = \mathcal{O}(1/\sqrt{K})$. It is worth noting that the $\mathcal{O}(\sqrt{\varepsilon_{\text{bias}}})$ term is the same as in [27], and we also apply the average SGD [54] to solve the NPG subproblem (2.8).

Theorem 4.11. In SRVR-PG (Algorithm 2), take $\eta = \frac{1}{8L_J}$, $S = \mathcal{O}\left(\frac{1}{(1-\gamma)^{2.5}\varepsilon}\right)$, $m = \mathcal{O}\left(\frac{(1-\gamma)^{0.5}}{\varepsilon}\right)$, $B = \mathcal{O}\left(\frac{W}{(1-\gamma)^{0.5}\varepsilon}\right)$, $N = \mathcal{O}\left(\frac{\sigma^2}{\varepsilon}\right)$, and $H = \mathcal{O}\left(\log\left(\frac{1}{(1-\gamma)\varepsilon}\right)\right)$. Then, we have

$$J^* - \frac{1}{Sm} \sum_{s=0}^{S-1} \sum_{t=0}^{m-1} \mathbb{E}[J(\theta_t^{j+1})] \leq \frac{\sqrt{\varepsilon_{\text{bias}}}}{1-\gamma} + \varepsilon.$$

In total, SRVR-PG samples $\mathcal{O}\left(\frac{W+\sigma^2}{(1-\gamma)^{2.5}\varepsilon^3}\right)$ trajectories.

Remark 4.12. Theorem 4.11 establishes the global convergence of SRVR-PG proposed in [21], where only stationary convergence is shown. Also, compared with stochastic PG, SRVR-PG enjoys a better sample complexity thanks to its faster stationary convergence.

Theorem 4.13. In SRVR-NPG (Algorithm 1), let us apply $\mathcal{O}\left(\frac{1}{(1-\gamma)^4\varepsilon^2}\right)$ iterations of SGD as in Procedure 2 to obtain an update direction. In addition, take $\eta = \frac{\mu_F}{16L_J}$, $S = \mathcal{O}\left(\frac{1}{(1-\gamma)^{2.5}\varepsilon^{0.5}}\right)$, $m = \mathcal{O}\left(\frac{(1-\gamma)^{0.5}}{\varepsilon^{0.5}}\right)$, $B = \mathcal{O}\left(\frac{W}{(1-\gamma)^{0.5}\varepsilon^{1.5}}\right)$, $N = \mathcal{O}\left(\frac{\sigma^2}{\varepsilon^2}\right)$, and $H = \mathcal{O}\left(\log\left(\frac{1}{(1-\gamma)\varepsilon}\right)\right)$. Then,

$$J^* - \frac{1}{Sm} \sum_{s=0}^{S-1} \sum_{t=0}^{m-1} \mathbb{E}[J(\theta_t^{j+1})] \leq \frac{\sqrt{\varepsilon_{\text{bias}}}}{1-\gamma} + \varepsilon.$$

In total, SRVR-NPG samples $\mathcal{O}\left(\frac{W+\sigma^2}{(1-\gamma)^{2.5}\varepsilon^{2.5}} + \frac{1}{(1-\gamma)^6\varepsilon^3}\right)$ trajectories.

Remark 4.14. Compared with SRVR-PG, our SRVR-NPG has a better dependence on W and σ^2 , which could be large in practice (especially W). The current sample complexity of SRVR-NPG is not better than our (improved) result of NPG since, in our analysis, the advantage of variance reduction is offset by the cost of solving the subproblems.

5 Numerical Experiments

In this section, we compare the numerical performances of stochastic PG, NPG, SRVR-PG, and SRVR-NPG. Specifically, we test on benchmark reinforcement learning environments Cartpole and Mountain Car. Our implementation is based on the implementation of SRVPG¹ and SRVR-PG², and can be found in the supplementary material.

For both tasks, we apply a Gaussian policy of the form $\pi_\theta(a | s) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\mu_\theta(s)-a)^2}{2\sigma^2}\right)$ where the mean $\mu_\theta(s)$ is modeled by a neural network with Tanh as the activation function.

For the Cartpole problem, we apply a neural network of size 32×1 and a horizon of $H = 100$. In addition, each training algorithm uses 5000 trajectories in total. For the Mountain Car problem, we apply a neural network of size 64×1 and take $H = 1000$. 3000 trajectories are allowed for each algorithm. The numerical performance comparison, as well as the settings of algorithm-specific parameters, can be found in Figures 1 and 2. In App. O, we provide more implementation details.

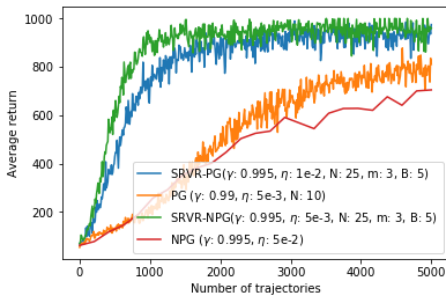


Figure 1: Numerical Performances on Cartpole. For PG, SRVR-PG and SRVR-NPG, we report the undiscounted average return averaged over 10 runs. For NPG, we report the averaged return over 40 runs. Overall, SRVR-NPG has the best performance.

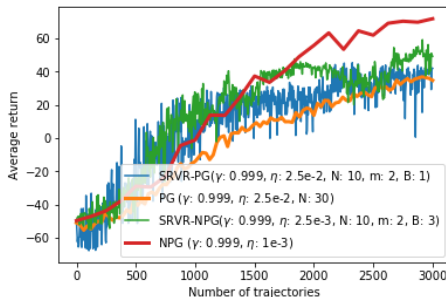


Figure 2: Numerical Performances on Mountain Car. For PG, SRVR-PG and SRVR-NPG, we report the undiscounted average return averaged over 10 runs. For NPG, we report the averaged return over 40 runs. Overall, NPG has the best performance.

6 Concluding Remarks

In this work, we have introduced a framework for analyzing the global convergence of (natural) PG methods and their variance-reduced variants, under the assumption that the Fisher information matrix is positive definite. We have established the sample complexity for the global convergence of stochastic PG and its variance-reduced variant SRVR-PG, and improved the sample complexity of NPG. In addition, we have introduced SRVR-NPG, which incorporates variance-reduction into NPG, and enjoys both global convergence guarantee and an efficient sample complexity. Our improved analysis hinges on exploiting the advantages of previous analyses on (variance reduced) PG and NPG methods, which may be of independent interest, and can be used to design faster variance-reduced NPG methods in the future.

¹<https://github.com/Dam930/rllab>

²<https://github.com/xgfelicia/SRVRPG>

Broader Impact

The results of this paper improves the performance of policy-gradient methods for reinforcement learning, as well as our understanding to the existing methods. Through reinforcement learning, our study will also benefit several research communities such as machine learning and robotics. We do not believe that the results in this work will cause any ethical issue, or put anyone at a disadvantage in our society.

Acknowledgements

Yanli Liu and Wotao Yin were partially supported by the Office of Naval Research (ONR) Grant N000141712162. Yanli Liu was also supported by UCLA Dissertation Year Fellowship. Kaiqing Zhang and Tamer Başar were supported in part by the US Army Research Laboratory (ARL) Cooperative Agreement W911NF-17-2-0196, and in part by the Office of Naval Research (ONR) MURI Grant N00014-16-1-2710.

References

- [1] Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT press, 2018.
- [2] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3-4):229–256, 1992.
- [3] Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems*, pages 1057–1063, 2000.
- [4] Sham M Kakade. A natural policy gradient. In *Advances in Neural Information Processing Systems*, pages 1531–1538, 2002.
- [5] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International Conference on Machine Learning*, pages 1889–1897, 2015.
- [6] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [7] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International Conference on Machine Learning*, pages 1928–1937, 2016.
- [8] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [9] Jonathan Baxter and Peter L Bartlett. Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research*, 15:319–350, 2001.
- [10] Cathy Wu, Aravind Rajeswaran, Yan Duan, Vikash Kumar, Alexandre M Bayen, Sham Kakade, Igor Mordatch, and Pieter Abbeel. Variance reduction for policy gradient with action-dependent factorized baselines. *arXiv preprint arXiv:1803.07246*, 2018.
- [11] Vijay R Konda and John N Tsitsiklis. Actor-critic algorithms. In *Advances in Neural Information Processing Systems*, pages 1008–1014, 2000.
- [12] Jan Peters and Stefan Schaal. Natural actor-critic. *Neurocomputing*, 71(7-9):1180–1190, 2008.
- [13] Shalabh Bhatnagar, Richard S Sutton, Mohammad Ghavamzadeh, and Mark Lee. Natural actor-critic algorithms. *Automatica*, 45(11):2471–2482, 2009.

- [14] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems*, pages 315–323, 2013.
- [15] Zeyuan Allen-Zhu and Elad Hazan. Variance reduction for faster non-convex optimization. In *International Conference on Machine Learning*, pages 699–707, 2016.
- [16] Sashank J Reddi, Ahmed Hefny, Suvrit Sra, Barnabás Póczos, and Alex Smola. Stochastic variance reduction for nonconvex optimization. In *International Conference on Machine Learning*, pages 314–323, 2016.
- [17] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems*, pages 1646–1654, 2014.
- [18] Tianbing Xu, Qiang Liu, and Jian Peng. Stochastic variance reduction for policy gradient estimation. *arXiv preprint arXiv:1710.06034*, 2017.
- [19] Matteo Papini, Damiano Binaghi, Giuseppe Canonaco, Matteo Pirodda, and Marcello Restelli. Stochastic variance-reduced policy gradient. In *Proceedings of the 35th International Conference on Machine Learning*, pages 4026–4035, 2018.
- [20] Pan Xu, Felicia Gao, and Quanquan Gu. An improved convergence analysis of stochastic variance-reduced policy gradient. *arXiv preprint arXiv:1905.12615*, 2019.
- [21] Pan Xu, Felicia Gao, and Quanquan Gu. Sample efficient policy gradient methods with recursive variance reduction. In *International Conference on Learning Representations*, 2020.
- [22] Huizhuo Yuan, Xiangru Lian, Ji Liu, and Yuren Zhou. Stochastic recursive momentum for policy gradient methods. *arXiv preprint arXiv:2003.04302*, 2020.
- [23] Maryam Fazel, Rong Ge, Sham M Kakade, and Mehran Mesbahi. Global convergence of policy gradient methods for the linear quadratic regulator. In *International Conference on Machine Learning*, 2018.
- [24] Kaiqing Zhang, Alec Koppel, Hao Zhu, and Tamer Başar. Global convergence of policy gradient methods to (almost) locally optimal policies. *arXiv preprint arXiv:1906.08383*, 2019.
- [25] Lingxiao Wang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural policy gradient methods: Global optimality and rates of convergence. *arXiv preprint arXiv:1909.01150*, 2019.
- [26] Jalaj Bhandari and Daniel Russo. Global optimality guarantees for policy gradient methods. *arXiv preprint arXiv:1906.01786*, 2019.
- [27] Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *arXiv preprint arXiv:1908.00261v5*, 2019.
- [28] Boyi Liu, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural trust region/proximal policy optimization attains globally optimal policy. In *Advances in Neural Information Processing Systems*, pages 10564–10575, 2019.
- [29] Lior Shani, Yonathan Efroni, and Shie Mannor. Adaptive trust region policy optimization: Global convergence and faster rates for regularized MDPs. *arXiv preprint arXiv:1909.02769*, 2019.
- [30] Tengyu Xu, Zhe Wang, and Yingbin Liang. Improving sample complexity bounds for actor-critic algorithms. *arXiv preprint arXiv:2004.12956*, 2020.
- [31] Tengyu Xu, Zhe Wang, and Yingbin Liang. Non-asymptotic convergence analysis of two time-scale (natural) actor-critic algorithms. *arXiv preprint arXiv:2005.03557*, 2020.
- [32] Yue Wu, Weitong Zhang, Pan Xu, and Quanquan Gu. A finite time analysis of two time-scale actor critic methods. *arXiv preprint arXiv:2005.01350*, 2020.

- [33] Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang. A two-timescale framework for bilevel optimization: Complexity analysis and application to actor-critic. *arXiv preprint arXiv:2007.05170*, 2020.
- [34] Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. Sarah: A novel method for machine learning problems using stochastic recursive gradient. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2613–2621. JMLR. org, 2017.
- [35] Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. In *Advances in Neural Information Processing Systems*, pages 689–699, 2018.
- [36] Yanli Liu, Fei Feng, and Wotao Yin. Acceleration of svrg and katyusha x by inexact preconditioning. In *International Conference on Machine Learning*, pages 4003–4012, 2019.
- [37] Nhan H Pham, Lam M Nguyen, Dzung T Phan, Phuong Ha Nguyen, Marten van Dijk, and Quoc Tran-Dinh. A hybrid stochastic policy gradient algorithm for reinforcement learning. *arXiv preprint arXiv:2003.00430*, 2020.
- [38] Feihu Huang, Shangqian Gao, Jian Pei, and Heng Huang. Momentum-based policy gradient methods. In *Thirty-seventh International Conference on Machine Learning (ICML 2020)*, 2020.
- [39] Richard S Sutton, Andrew G Barto, and Ronald J Williams. Reinforcement learning is direct adaptive optimal control. *IEEE Control Systems Magazine*, 12(2):19–22, 1992.
- [40] Charles George Broyden. The convergence of a class of double-rank minimization algorithms 1. general considerations. *IMA Journal of Applied Mathematics*, 6(1):76–90, 1970.
- [41] Roger Fletcher. A new approach to variable metric algorithms. *The computer journal*, 13(3):317–322, 1970.
- [42] Donald Goldfarb. A family of variable-metric methods derived by variational means. *Mathematics of computation*, 24(109):23–26, 1970.
- [43] David F Shanno. Conditioning of quasi-newton methods for function minimization. *Mathematics of computation*, 24(111):647–656, 1970.
- [44] Richard H Byrd, Samantha L Hansen, Jorge Nocedal, and Yoram Singer. A stochastic quasi-newton method for large-scale optimization. *SIAM Journal on Optimization*, 26(2):1008–1031, 2016.
- [45] Philipp Moritz, Robert Nishihara, and Michael Jordan. A linearly-convergent stochastic l-bfgs algorithm. In *Artificial Intelligence and Statistics*, pages 249–258, 2016.
- [46] Robert Gower, Donald Goldfarb, and Peter Richtárik. Stochastic block bfgs: Squeezing more curvature out of data. In *International Conference on Machine Learning*, pages 1869–1878, 2016.
- [47] Xiao Wang, Shiqian Ma, Donald Goldfarb, and Wei Liu. Stochastic quasi-newton methods for nonconvex stochastic optimization. *SIAM Journal on Optimization*, 27(2):927–956, 2017.
- [48] Yan Duan, Xi Chen, Rein Houthooft, John Schulman, and Pieter Abbeel. Benchmarking deep reinforcement learning for continuous control. In *International Conference on Machine Learning*, pages 1329–1338, 2016.
- [49] John N Tsitsiklis and Benjamin Van Roy. Analysis of temporal-difference learning with function approximation. In *Advances in Neural Information Processing Systems*, pages 1075–1081, 1997.
- [50] Francisco S Melo, Sean P Meyn, and M Isabel Ribeiro. An analysis of reinforcement learning with function approximation. In *International Conference on Machine Learning*, pages 664–671, 2008.

- [51] Richard S Sutton, Hamid Reza Maei, Doina Precup, Shalabh Bhatnagar, David Silver, Csaba Szepesvári, and Eric Wiewiora. Fast gradient-descent methods for temporal-difference learning with linear function approximation. In *International Conference on Machine Learning*, pages 993–1000, 2009.
- [52] Shun-Ichi Amari. Natural gradient works efficiently in learning. *Neural computation*, 10(2):251–276, 1998.
- [53] James Martens. New insights and perspectives on the natural gradient method. *arXiv preprint arXiv:1412.1193*, 2014.
- [54] Francis Bach and Eric Moulines. Non-strongly-convex smooth stochastic approximation with convergence rate $O(1/n)$. In *Advances in Neural Information Processing Systems*, pages 773–781, 2013.
- [55] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009.
- [56] Saeed Ghadimi and Guanghui Lan. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization i: A generic algorithmic framework. *SIAM Journal on Optimization*, 22(4):1469–1492, 2012.
- [57] Matteo Pirotta, Marcello Restelli, and Luca Bascetta. Adaptive step-size for policy gradient methods. In *Advances in Neural Information Processing Systems*, pages 1394–1402, 2013.
- [58] Corinna Cortes, Yishay Mansour, and Mehryar Mohri. Learning bounds for importance weighting. In *Advances in Neural Information Processing Systems*, pages 442–450, 2010.
- [59] Sham M Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *International Conference on Machine Learning*, volume 2, pages 267–274, 2002.
- [60] Anirban DasGupta. The exponential family and statistical applications. In *Probability for Statistics and Machine Learning*, pages 583–612. Springer, 2011.
- [61] Solomon Kullback. *Information Theory and Statistics*. Courier Corporation, 1997.
- [62] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- [63] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.