

# Adaptive Step-Size for Policy Gradient Methods

**Authored by:**

Marcello Restelli  
Matteo Pirotta  
Luca Bascetta

## **Abstract**

In the last decade, policy gradient methods have significantly grown in popularity in the reinforcement-learning field. In particular, they have been largely employed in motor control and robotic applications, thanks to their ability to cope with continuous state and action domains and partial observable problems. Policy gradient researches have been mainly focused on the identification of effective gradient directions and the proposal of efficient estimation algorithms. Nonetheless, the performance of policy gradient methods is determined not only by the gradient direction, since convergence properties are strongly influenced by the choice of the step size: small values imply slow convergence rate, while large values may lead to oscillations or even divergence of the policy parameters. Step-size value is usually chosen by hand tuning and still little attention has been paid to its automatic selection. In this paper, we propose to determine the learning rate by maximizing a lower bound to the expected performance gain. Focusing on Gaussian policies, we derive a lower bound that is second-order polynomial of the step size, and we show how a simplified version of such lower bound can be maximized when the gradient is estimated from trajectory samples. The properties of the proposed approach are empirically evaluated in a linear-quadratic regulator problem.

## **1 Paper Body**

Policy gradient methods have established as the most effective reinforcement-learning techniques in robotic applications. Such methods perform a policy search to maximize the expected return of a policy in a parameterized policy class. The reasons for their success are many. Compared to several traditional reinforcement-learning approaches, policy gradients scale well to high-dimensional continuous state and action problems, and no changes to the algorithms are needed to face uncertainty in the state due to limited and noisy sensors. Furthermore, policy representation can be properly designed for the given task, thus allowing to incorporate domain knowledge into the algorithm useful to speed up the learning process and to prevent the unexpected execution of dangerous policies

that may harm the system. Finally, they are guaranteed to converge to locally optimal policies. Thanks to these advantages, from the 1990s policy gradient methods have been widely used to learn complex control tasks [1]. The research in these years has focused on obtaining good model-free estimators of the policy gradient using data generated during the task execution. The oldest policy gradient approaches are finite-difference methods [2], that estimate gradient direction by resolving a regression problem based on the performance evaluation of policies associated to different small perturbations of the current parameterization. Finite-difference methods have some advantages: they are easy to implement, do not need assumptions on the differentiability of the policy w.r.t. the policy parameters, and are efficient in deterministic settings. On the other hand, when used on real systems, the choice of parameter perturbations may be difficult and critical for system safeness. Furthermore, the presence of uncertainties may significantly slow down the convergence rate. Such drawbacks have been overcome by likelihood ratio methods [3, 4, 5], since they do not need to generate policy parameters variations and quickly converge even in highly stochastic systems. Several

studies have addressed the problem to find minimum variance estimators by the computation of optimal baselines [6]. To further improve the efficiency of policy gradient methods, natural gradient approaches (where the steepest ascent is computed w.r.t. the Fisher information metric) have been considered [7, 8]. Natural gradients still converge to locally optimal policies, are independent from the policy parameterization, need less data to attain good gradient estimate, and are less affected by plateaus. Once an accurate estimate of the gradient direction is obtained, policy parameters are updated by:  $\theta_{t+1} = \theta_t + \alpha_t \hat{g}_t$ , where  $\alpha_t \in \mathbb{R}_+$  is the step size in the direction of the gradient. Although, given an unbiased gradient estimate, convergence to a local optimum can be guaranteed under mild conditions over the learning rate values [9], their choice may significantly affect the convergence speed or the behavior during the transient. Updating the policy with large step sizes may lead to policy oscillations or even divergence [10], while trying to avoid such phenomena by using small learning rates determines a growth in the number of iterations that is unbearable in most real-world applications. In general unconstrained programming, the optimal step size for gradient ascent methods is determined through line-search algorithms [11], that require to try different values for the learning rate and evaluate the function value in the corresponding updated points. Such an approach is unfeasible for policy gradient methods, since it would require to perform a large number of policy evaluations. Despite these difficulties, up to now, little attention has been paid to the study of step size computation for policy gradient algorithms. Nonetheless, some policy search methods based on expectation-maximization have been recently proposed; such methods have properties similar to the ones of policy gradients, but the policy update does not require to tune the step size [12, 13]. In this paper, we propose a new approach to compute the step size in policy gradient methods that guarantees an improvement at each step, thus avoiding oscillation and divergence issues. Starting from a lower bound to the difference of performance between two poli-

cies, in Section 3 we derive a lower bound in the case where the new policy is obtained from the old one by changing its parameters along the gradient direction. Such a new bound is a (polynomial) function of the step size, that, for positive values of the step size, presents a single, positive maximum ( i.e., it guarantees improvement) which can be computed in closed form. In Section 4, we show how the bound simplifies to a quadratic function of the step size when Gaussian policies are considered, and Section 5 studies how the bound needs to be changed in approximated settings (e.g., model-free case) where the policy gradient needs to be estimated directly from experience.

2

## Preliminaries

A discrete-time continuous Markov decision process (MDP) is defined as a 6-tuple  $\langle S, A, P, R, \gamma, D \rangle$ , where  $S$  is the continuous state space,  $A$  is the continuous action space,  $P$  is a Markovian transition model where  $P(s_0 | s, a)$  defines the transition density between state  $s$  and  $s_0$  under action  $a$ ,  $R : S \times A \rightarrow [0, R]$  is the reward function, such that  $R(s, a)$  is the expected immediate reward for the state-action pair  $(s, a)$  and  $R$  is the maximum reward value,  $\gamma \in [0, 1)$  is the discount factor for future rewards, and  $D$  is the initial state distribution. The policy of an agent is characterized by a density distribution  $\pi(a | s)$  that specifies for each state  $s$  the density distribution over the action space  $A$ . To measure the distance between two policies we will use this norm:  $\| \pi - \pi' \|_1 = \int_S \int_A |\pi(a | s) - \pi'(a | s)| da ds$

A

that is the superior value over the state space of the total variation between the distributions over the action space of policy  $\pi_0$  and  $\pi$ . We consider infinite horizon problems where the future rewards are exponentially discounted with  $\gamma$ . For each state  $s$ , we define the utility of following a stationary policy  $\pi$  as:  $U^\pi(s) = \mathbb{E} \sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(a_t | s_t)) | s_0 = s$ .  $\pi \in \Pi$

It is known that  $V$

$\pi$

$t=0$

solves the following recursive (Bellman) equation:  $V^\pi(s) = R(s, \pi(a | s)) + \gamma \int_S P(s_0 | s, \pi(a | s)) V^\pi(s_0) ds_0$ .  $A$

$S$

2

Policies can be ranked by their expected discounted reward starting from the state distribution  $D$ :  $J^\pi = \int_S D(s) V^\pi(s) ds = \int_S D(s) (R(s, \pi(a | s)) + \gamma \int_S P(s_0 | s, \pi(a | s)) V^\pi(s_0) ds_0) ds$

$d^\pi D(s)$

$P^\pi$

$S$

$A$

$t$

where  $d^\pi = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P^t(s_0 = s | D)$  is the discounted future state distribution for a starting state distribution  $D$  [5]. Solving an MDP means to find a policy  $\pi^*$  that maximizes  $J^\pi$  the expected long-term reward:  $\pi^* = \arg \max_{\pi \in \Pi} J^\pi$

$\arg \max_{\pi} J(\pi)$ . For any MDP there exists at least one deterministic optimal policy that simultaneously maximizes  $V^{\pi}(s)$ ,  $\forall s \in \mathcal{S}$ . For control purposes, it is better to consider action values  $Q^{\pi}(s, a)$ , i.e., the value of taking action  $a$  in state  $s$  and following a policy  $\pi$  thereafter: 
$$Q^{\pi}(s, a) = R(s, a) + \gamma \sum_{s'} P(s' | s, a) V^{\pi}(s').$$

Furthermore, we define the advantage function:  $A^{\pi}(s, a) = Q^{\pi}(s, a) - V^{\pi}(s)$ , that quantifies the advantage (or disadvantage) of taking action  $a$  in state  $s$  instead of following policy  $\pi$ . In for each state  $s$ , we define the advantage of a policy  $\pi$  over policy  $\pi_0$  as  $A^{\pi}(s) = A^{\pi}(s, \pi(s)) - A^{\pi_0}(s, \pi_0(s))$  and, following [14], we define its expected value w.r.t. an initial  $\rho_0$  state distribution  $\rho_0$  as  $A^{\pi} = \sum_s \rho_0(s) A^{\pi}(s)$ . We consider the problem of finding a policy that maximizes the expected discounted reward over a class of parameterized policies  $\Pi = \{\pi_{\theta} : \theta \in \Theta\}$ , where  $\Theta$  is a compact representation of  $(a|s)$ . The exact gradient of the expected discounted reward w.r.t. the policy parameters [5] is: 
$$\nabla_{\theta} J(\theta) = \sum_s \rho_0(s) \sum_a \pi_{\theta}(a|s) \nabla_{\theta} Q^{\pi_{\theta}}(s, a).$$
 The policy parameters can be updated by following the direction of the gradient of the expected discounted reward:  $\theta_0 = \theta + \alpha \nabla_{\theta} J(\theta)$ . In the following, we will denote with  $\|\nabla_{\theta} J(\theta)\|_1$  and  $\|\nabla_{\theta} J(\theta)\|_2$  the L1 and L2 norm of the policy gradient vector, respectively.

3

### Policy Gradient Formulation

In this section we provide a lower bound to the improvement obtained by updating the policy parameters along the gradient direction as a function of the step size. The idea is to start from the general lower bound on the performance difference between any pair of policies introduced in [15] and specialize it to the policy gradient framework. Lemma 3.1 (Continuous MDP version of Corollary 3.6 in [15]). For any pair of stationary policies corresponding to parameters  $\theta$  and  $\theta_0$  and for any starting state distribution  $\rho_0$ , the difference between the performance of policy  $\pi_{\theta_0}$  and policy  $\pi_{\theta}$  can be bounded as follows 
$$J(\theta_0) - J(\theta) \geq \sum_s \rho_0(s) A^{\pi_{\theta_0}}(s) \sum_a \pi_{\theta}(a|s) \nabla_{\theta} Q^{\pi_{\theta_0}}(s, a) - \frac{1}{2} \sum_s \rho_0(s) \sum_a \pi_{\theta}(a|s) \sum_{a'} \pi_{\theta_0}(a'|s) \|Q^{\pi_{\theta_0}}(s, a) - Q^{\pi_{\theta_0}}(s, a')\|^2$$
 where  $\|Q^{\pi_{\theta_0}}\|_{\infty}$  is the supremum norm of the  $Q^{\pi_{\theta_0}}$  function:  $\|Q^{\pi_{\theta_0}}\|_{\infty} = \sup_{s,a} Q^{\pi_{\theta_0}}(s, a)$ .

As we can notice from the above bound, to maximize the performance improvement, we need to find a new policy  $\pi_{\theta_0}$  that is associated to large average advantage  $A^{\pi_{\theta_0}}$ , but, at the same time, is not too different from the current policy  $\pi_{\theta}$ . Policy gradient approaches provide search directions characterized by increasing advantage values and, through the step size value, allow to control the difference between the new policy and the target one. Exploiting a lower bound to the first order Taylor's expansion, we can bound the difference between the current policy and the new policy, whose parameters are adjusted along the gradient direction, as a function of the step size  $\alpha$ . Lemma 3.2. Let the update of the policy parameters be  $\theta_0 = \theta + \alpha \nabla_{\theta} J(\theta)$ . Then 
$$J(\theta_0) \geq J(\theta) + \alpha \sum_s \rho_0(s) \sum_a \pi_{\theta}(a|s) \nabla_{\theta} Q^{\pi_{\theta}}(s, a) - \frac{\alpha^2}{2} \sum_s \rho_0(s) \sum_a \pi_{\theta}(a|s) \sum_{a'} \pi_{\theta}(a'|s) \|\nabla_{\theta} Q^{\pi_{\theta}}(s, a) - \nabla_{\theta} Q^{\pi_{\theta}}(s, a')\|^2$$

T

$J(\pi_{\theta}(\mathbf{s}), \pi_{\theta}(\mathbf{s})) = J(\pi_{\theta}(\mathbf{s}), \pi_{\theta}(\mathbf{s})) + J(\pi_{\theta}(\mathbf{s}), \pi_{\theta}(\mathbf{s}))$   
where  $J(\pi_{\theta}(\mathbf{s}), \pi_{\theta}(\mathbf{s})) = J(\pi_{\theta}(\mathbf{s}), \pi_{\theta}(\mathbf{s}))$ .

$\inf_{\theta} c(\theta, 1)$   
!

$\sum_{i,j=1}^m X_{ij}^2 (\pi_{\theta}(\mathbf{s}), \pi_{\theta}(\mathbf{s})) + c(\theta, 1) + I(i=j)$

By combining the two previous lemmas, it is possible to derive the policy performance improvement obtained following the gradient direction. Theorem 3.3. Let the update of the parameters be  $\theta_0 = \theta + \eta \nabla J(\theta)$ . Then for any stationary policy  $\pi_{\theta}(\mathbf{s})$  and any starting state distribution  $\nu$ , the difference in performance between  $\pi_{\theta_0}$  and  $\pi_{\theta}$  is lower bounded by:  $J(\pi_{\theta_0}, \nu) - J(\pi_{\theta}, \nu) \geq \eta \sum_{i,j=1}^m X_{ij}^2 (\pi_{\theta}(\mathbf{s}), \pi_{\theta}(\mathbf{s})) + c(\theta, 1) + I(i=j)$

$\sum_{i,j=1}^m X_{ij}^2 (\pi_{\theta}(\mathbf{s}), \pi_{\theta}(\mathbf{s})) + c(\theta, 1) + I(i=j)$

$\sum_{i,j=1}^m X_{ij}^2 (\pi_{\theta}(\mathbf{s}), \pi_{\theta}(\mathbf{s})) + c(\theta, 1) + I(i=j)$

$\sum_{i,j=1}^m X_{ij}^2 (\pi_{\theta}(\mathbf{s}), \pi_{\theta}(\mathbf{s})) + c(\theta, 1) + I(i=j)$

$\sum_{i,j=1}^m X_{ij}^2 (\pi_{\theta}(\mathbf{s}), \pi_{\theta}(\mathbf{s})) + c(\theta, 1) + I(i=j)$

The above bound is a forth order polynomial of the step size, whose stationary points, being the roots of a third order polynomial  $ax^3 + bx^2 + cx + d$ , can be expressed in closed form. It is worth to notice that, for positive values of  $\eta$ , the bound presents a single stationary point that corresponds to a local maximum. In fact, since  $a, b \geq 0$  and  $d \geq 0$ , the Descartes' rule of signs gives the existence and uniqueness of the real positive root. In the following section, we will show, in the case of Gaussian policies, how the bound in Theorem 3.3 can be reduced to a second order polynomial in  $\eta$ , thus obtaining a simpler closed-form solution for optimal (w.r.t. the bound) step size.

4

The Gaussian Policy Model

In this section we consider the Gaussian policy model with fixed standard deviation  $\sigma$  and the mean is a linear combination of the state feature vector  $\phi(\mathbf{s})$  using a parameter vector  $\theta$  of size  $m$ :

$\pi_{\theta}(\mathbf{s}) = \frac{1}{\sigma^m} \exp\left(-\frac{1}{2\sigma^2} \phi(\mathbf{s})^T \theta\right)$  In the case of Gaussian policies, each second order derivative of policy  $\pi_{\theta}$  can be easily bounded. Lemma 4.1. For any Gaussian policy  $\pi_{\theta}(\mathbf{s}) = \frac{1}{\sigma^m} \exp\left(-\frac{1}{2\sigma^2} \phi(\mathbf{s})^T \theta\right)$ , the second order derivative of the policy can be bounded as follows:

$\frac{\partial^2 \pi_{\theta}(\mathbf{s})}{\partial \theta_i \partial \theta_j} \leq \frac{1}{\sigma^m} \phi(\mathbf{s})_i \phi(\mathbf{s})_j$

This result allows to restate Lemma 3.2 in the case of Gaussian policies:  $J(\pi_{\theta_0}, \nu) - J(\pi_{\theta}, \nu) \geq \eta \sum_{i,j=1}^m X_{ij}^2 (\pi_{\theta}(\mathbf{s}), \pi_{\theta}(\mathbf{s})) + c(\theta, 1) + I(i=j)$

$\sum_{i,j=1}^m X_{ij}^2 (\pi_{\theta}(\mathbf{s}), \pi_{\theta}(\mathbf{s})) + c(\theta, 1) + I(i=j)$

In the following we will assume that features  $\phi$  are uniformly bounded: Assumption 4.1. All the basis functions are uniformly bounded by  $M$ :  $\phi(\mathbf{s})_i \leq M$ ,  $\forall \mathbf{s} \in \mathcal{S}$ ,  $i = 1, \dots, m$ . Exploiting Pinsker's inequality [16]

(which upper bounds the total variation between two distributions with their Kullback-Liebler divergence), it is possible to provide the following upper bound to the supremum norm between two Gaussian policies. Lemma 4.2. For any pair of stationary policies  $\pi$  and  $\pi_0$ , so that  $\pi_0 = \pi + \eta \nabla J(\pi)$ , supremum norm of their difference can be upper bounded as follows:  $\|\pi - \pi_0\|_1 \leq 4\eta \|\nabla J(\pi)\|_1$ .

$$\|\pi - \pi_0\|_1 \leq 4\eta \|\nabla J(\pi)\|_1$$

By plugging the results of Lemmas 4.1 and 4.2 into Equation (1) we can obtain a lower bound to the performance difference between a Gaussian policy  $\pi$  and another policy along the gradient direction that is quadratic in the step size  $\eta$ . Theorem 4.3. For any starting state distribution  $\mu$ , and any pair of stationary Gaussian policies  $\pi \in \mathcal{N}(\mu, \Sigma)$  and  $\pi_0 \in \mathcal{N}(\mu, \Sigma)$ , so that  $\pi_0 = \pi + \eta \nabla J(\pi)$  and under Assumption 4.1, the difference between the performance of  $\pi_0$  and the one of  $\pi$  can be lower bounded as follows:

$$J(\pi_0) - J(\pi) \geq \eta \nabla J(\pi)^T \mu - \frac{\eta^2}{2} \nabla^2 J(\pi) \mu$$

Since the linear coefficient is positive and the quadratic one is negative, the bound in Theorem 4.3 has a single maximum attained for some positive value of  $\eta$ . Corollary 4.4. The performance lower bound provided in Theorem 4.3 is maximized by choosing the following step size:  $\eta^* = \frac{\nabla J(\pi)^T \mu}{\nabla^2 J(\pi) \mu}$ .

$$\eta^* = \frac{\nabla J(\pi)^T \mu}{\nabla^2 J(\pi) \mu}$$

that guarantees the following policy performance improvement  $J(\pi_0) - J(\pi) \geq \frac{1}{2} \nabla J(\pi)^T \mu$ .

5

#### Approximate Framework

The solution for the tuning of the step size presented in the previous section depends on some constants (e.g., discount factor and the variance of the Gaussian policy) and requires to be able to compute some quantities (e.g., the policy gradient and the supremum value of the  $Q$ -function). In many real-world applications such quantities cannot be computed (e.g., when the state-transition model is unknown or too large for exact methods) and need to be estimated from experience samples. In this section, we study how the step size can be chosen when the gradient is estimated through sample trajectories to guarantee a performance improvement in high probability. For sake of easiness, we consider a simplified version of the bound in Theorem 4.3, in order to obtain a bound where the only element that needs to be estimated is the policy gradient  $\nabla J(\pi)$ . Corollary 5.1. For any starting state distribution  $\mu$ , and any pair of stationary Gaussian policies  $\pi \in \mathcal{N}(\mu, \Sigma)$  and  $\pi_0 \in \mathcal{N}(\mu, \Sigma)$ , so that  $\pi_0 = \pi + \eta \nabla J(\pi)$  and under Assumption 4.1, the difference between the performance of  $\pi_0$  and  $\pi$  is lower bounded by:

$$J(\pi_0) - J(\pi) \geq \eta \nabla J(\pi)^T \mu - \frac{\eta^2}{2} \nabla^2 J(\pi) \mu$$

value:  $\frac{1}{2} (1 - \gamma) \sum_{k=0}^{\infty} \gamma^k J^k(\theta) = ?$

2.  $\frac{1}{2} \sum_{k=0}^{\infty} \gamma^k \|J^k(\theta) - J^k(\theta^*)\|_2^2$ —A— RM?  $\frac{1}{2} \sum_{k=0}^{\infty} \gamma^k \|J^k(\theta) - J^k(\theta^*)\|_2^2$  Since we are assuming that the policy gradient  $J^k(\theta)$  is estimated through trajectory samples, the lower bound in Corollary 5.1 must take into consideration the associated approximation error.  $\frac{1}{2} \sum_{k=0}^{\infty} \gamma^k \|J^k(\theta) - J^k(\theta^*)\|_2^2$  of Given a set of trajectories obtained following policy  $\theta$ , we can produce an estimate  $\hat{J}^k(\theta)$  of the policy gradient and we assume to be able to produce a vector  $\mathbf{e}^k = [e_1^k, \dots, e_m^k]$ , so that the  $i$ th component of the approximation error is bounded at least with probability  $1 - \delta_i$ :

$$|e_i^k| \leq \frac{1}{\delta_i} \quad \text{P} \geq 1 - \delta_i \quad \forall k \leq m$$

Given the approximation error vector  $\mathbf{e}^k$ , we can adjust the bound in Corollary 5.1 to produce a new  $m$  bound that holds at least with probability  $(1 - \delta)$ . In particular, to preserve the inequality sign, the estimated approximation error must be used to decrease the  $L_2$ -norm of the policy gradient in the first term (the one that provides the positive contribution to the performance improvement) and to increase the  $L_1$ -norm in the penalization term. To lower bound the  $L_2$ -norm, we introduce the  $\mathbf{b}^k = J^k(\theta)$  whose components are a lower bound to the absolute value of the policy gradient vector  $\hat{J}^k(\theta)$  built on the basis of the approximation error:  $\mathbf{b}^k = J^k(\theta) - \mathbf{e}^k$ , where  $\mathbf{0}$  denotes the  $m$ -size vector with all zeros, and  $\max$  denotes the component-wise maximum.  $\mathbf{b}^k = J^k(\theta)$ : Similarly, to upper bound the  $L_1$ -norm of the policy gradient, we introduce the vector  $\mathbf{c}^k = J^k(\theta) + \mathbf{e}^k$ . Theorem 5.2. Under the same assumptions of Corollary 5.1, and provided

that it is available a

$\mathbf{b}^k$  policy gradient estimate  $\hat{J}^k(\theta)$ , so that  $\mathbf{P} \geq 1 - \delta_i \quad \forall i \leq m$ , the difference  $m$

between the performance of  $\theta_0$  and  $\theta$  can be lower bounded at least with probability  $(1 - \delta)$ :

$$\frac{1}{2} \sum_{k=0}^{\infty} \gamma^k \|J^k(\theta) - J^k(\theta^*)\|_2^2 - \frac{1}{2} \sum_{k=0}^{\infty} \gamma^k \|\mathbf{e}^k\|_2^2$$

that is maximized by the following step size value:

$$\frac{1}{2} \sum_{k=0}^{\infty} \gamma^k \|J^k(\theta) - J^k(\theta^*)\|_2^2 = \frac{1}{2} \sum_{k=0}^{\infty} \gamma^k \|\mathbf{e}^k\|_2^2$$

$$\frac{1}{2} \sum_{k=0}^{\infty} \gamma^k \|J^k(\theta) - J^k(\theta^*)\|_2^2 = \frac{1}{2} \sum_{k=0}^{\infty} \gamma^k \|\mathbf{e}^k\|_2^2$$

In the following, we will discuss how the approximation error of the policy gradient can be bounded. Among the several methods that have been proposed over the years, we focus on two well-understood policy-gradient estimation approaches: REINFORCE [3] and G(PO)MDP [4]/policy gradient theorem (PGT) [5]. 5.1

Approximation with REINFORCE gradient estimator

The REINFORCE approach [3] is the main exponent of the likelihood ratio family. The episodic REINFORCE gradient estimator is given by: 
$$\hat{g}_i = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^H (R_{n,k} - b) \nabla \log \pi(a_n; s_n, \theta) \quad (5.1)$$

where  $N$  is the number of  $H$ -step trajectories generated from a system by rollouts and  $b$  is a baseline that can be chosen arbitrary, but usually with the goal of minimizing the variance of the gradient estimator. The main drawback of REINFORCE is its variance, that is strongly affected by the length of the trajectory horizon  $H$ . The goal is to determine the number of trajectories  $N$  in order to obtain the desired accuracy of the gradient estimate. To achieve this, we exploit the upper bound to the variance of the episodic REINFORCE gradient estimator introduced in [17] for Gaussian policies. Lemma 5.3 (Adapted from Theorem 2 in [17]). Given a Gaussian policy  $\pi(a|s, \theta)$ ,

$N \geq \frac{2}{\epsilon^2} \frac{1}{\gamma^2} \frac{1}{\sigma^2} \frac{1}{\rho}$ , under the assumption of uniformly bounded rewards and basis functions (Assumption 4.1), we have the following upper bound to the variance of the  $i$ -th component of the episodic REINFORCE gradient estimate  $\hat{g}_i$

$$\text{Var}(\hat{g}_i) \leq \frac{2}{N} \frac{1}{\gamma^2} \frac{1}{\sigma^2} \frac{1}{\rho} \quad (5.2)$$

The result in the previous Lemma combined with the Chebyshev's inequality allows to provide a high probability upper bound to the gradient approximation error using the episodic REINFORCE gradient estimator.

Theorem 5.4. Given a Gaussian policy  $\pi(a|s, \theta)$ ,  $N \geq \frac{2}{\epsilon^2} \frac{1}{\gamma^2} \frac{1}{\sigma^2} \frac{1}{\rho}$ , under the assumption of uniformly bounded rewards and basis functions (Assumption 4.1), using the following number of  $H$ -step trajectories:  $N \geq \frac{2}{\epsilon^2} \frac{1}{\gamma^2} \frac{1}{\sigma^2} \frac{1}{\rho}$ , the REINFORCE is such that with probability  $1 - \delta$ : the gradient estimate  $\hat{g}_i$

$$b - \epsilon \leq \hat{g}_i \leq b + \epsilon \quad (5.2)$$

Approximation with G(PO)MDP/PGT gradient estimator

Although the REINFORCE method is guaranteed to converge at the true gradient at the fastest possible pace, its large variance can be problematic in practice. Advances in the likelihood ratio gradient estimators have produced new approaches that significantly reduce the variance of the estimate. Focusing on the class of "vanilla" gradient estimator, two main approaches have been proposed: policy gradient theorem (PGT) [5] and G(PO)MDP [4]. In [6], the authors show that, while the algorithms  $\hat{g}_i^{\text{G(PO)MDP}}$  and  $\hat{g}_i^{\text{PGT}}$  look different, their gradient estimate are equal, i.e., for reason, we can limit our attention to the PGT formulation: 
$$\hat{g}_i^{\text{PGT}} = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^H \nabla \log \pi(a_n; s_n, \theta) \quad (5.3)$$

$$l=k$$

$$b_{nl}$$

where  $R$  have the objective to reduce the variance of the gradient estimate. Following the procedure used to bound the approximation error of REINFORCE, we need an upper bound to the variance of the gradient estimate of PGT that is provided by the following lemma (whose proof is similar to the one used in [17] for the REINFORCE case).



Lemma 5.5. Given a Gaussian policy  $\pi(a|s, \theta) \in \mathcal{N}(\mu(s, \theta), \Sigma)$ , under the assumption of uniformly bounded rewards and basis functions (Assumption 4.1), we have the following upper bound on the variance of the  $i$ -th component of the PGT gradient estimate  $\hat{g}_i$ :

$$\begin{aligned} \text{Var}(\hat{g}_i) &\leq \frac{1}{N} \sum_{t=1}^N \text{Tr}(\Sigma^{-1} \mathbb{E}[\nabla^2 \log \pi(a_t|s_t, \theta_t)] \Sigma) \\ &\leq \frac{1}{N} \sum_{t=1}^N \text{Tr}(\Sigma^{-1} \mathbb{E}[\nabla^2 \log \pi(a_t|s_t, \theta_t)] \Sigma) \\ &\leq \frac{1}{N} \sum_{t=1}^N \text{Tr}(\Sigma^{-1} \mathbb{E}[\nabla^2 \log \pi(a_t|s_t, \theta_t)] \Sigma) \end{aligned}$$

As expected, since the variance of the gradient estimate obtained with PGT is smaller than the one with REINFORCE, also the upper bound of the PGT variance is smaller than REINFORCE one. In particular, while the variance with REINFORCE grows linearly with the time horizon, using PGT the dependence on the time horizon is significantly smaller. Finally, we can derive the upper bound for the approximation error of the gradient estimated of PGT.

Theorem 5.6. Given a Gaussian policy  $\pi(a|s, \theta) \in \mathcal{N}(\mu(s, \theta), \Sigma)$ , under the assumption of uniformly bounded rewards and basis functions (Assumption 4.1), using the following number of  $H$ -step trajectories:

$N = \frac{1}{\epsilon^2} \left( \frac{1}{\eta^2} + \frac{1}{\epsilon^2} \right)$ , the PGT gradient estimate  $\hat{g}_i$  generated by PGT is such that with probability  $1 - \delta$ : the

$\eta$	$N$
0.50	17138 1675
0.75	8669 697
1.00	5120 499
1.25	3348
1.50	23651 2342
1.75	17516 1714
2.00	13480 1287
5.00	21888 2163
7.50	9740 849

Table 1: Convergence speed in exact LQG scenario with  $\gamma = 0.95$ . The table reports the number of iterations required by the exact gradient approach, starting from  $\theta = 0$ , to learn the optimal policy parameter  $\theta^* = 0.6037$  with an accuracy of 0.01, for different step-size values. Three different set of experiments are shown: constant step size, decreasing step size, and the step size proposed in Corollary 4.4. The table contains itmax when no convergence happens in

30, 000 iterations, and  $\beta$  when the algorithm diverges ( $\beta \geq 1$  or  $\beta \leq 0$ ). Best performances are reported in boldface. 10, 000 RF PGT

it 822 29, 761

$\beta$  0.0030 0.2176

Number of trajectories 100, 000 it  $\beta$  51, 731 0.3068 63, 985 0.4013

500, 000 it  $\beta$  75, 345 0.4088 83, 983 0.4558

Table 2: Convergence speed in approximate LQG scenario with  $\gamma = 0.9$ . The table reports, starting from  $\beta = 0$  and fixed  $\gamma = 1$ , the number of iterations performed before the proposed step size  $\beta$  becomes 0 and the last value of the policy parameter. Results are shown for different number of trajectories (of 20 steps each) used in the gradient estimation by REINFORCE and PGT.

6

## Numerical Simulations and Discussion

In this section we show results related to some numerical simulations of policy gradient in the linear-quadratic Gaussian regulation (LQG) problem as formulated in [6]. The LQG problem is

characterized by a transition model  $s_{t+1} = N(s_t + a_t, \Sigma)$ , Gaussian policy  $a_t = N(\mu(s_t), \Sigma_a)$  and quadratic reward  $r_t = -0.5(s_t^T Q s_t + a_t^T R a_t)$ . The range of state and action spaces is bounded to the interval  $[-2, 2]$  and the initial state is drawn uniformly at random. This scenario is particularly instructive since it allows to exactly compute all terms involved in the bounds. We first present results in the exact scenario and then we move toward the approximated one. Table 1 shows how the number of iterations required to learn a near-optimal value of the policy parameter changes according to the standard deviation of the Gaussian policy and the step-size value. As expected, very small values of the step size allow to avoid divergence, but the learning process needs many iterations to reach a good performance (this can be observed both when the step size is kept constant and when it decreases). On the other hand, larger step-size values may lead to divergence. In this example, the higher the policy variance, the lower is the step size value that allows to avoid divergence, since, in LQG, higher policy variance implies larger policy gradient values. Using the step size  $\beta$  from Corollary 4.4 the policy gradient algorithm avoids divergence (since it guarantees an improvement at each iteration), and the speed of convergence is strongly affected by the variance of the Gaussian policy. In general, when the policy are nearly deterministic (small variance in the Gaussian case), small changes in the parameters lead to large distances between the policies, thus negatively affecting the lower bound in Equation 1. As we can notice from the expression of  $\beta$  in Corollary 4.4, considering policies with high variance (that might be a problem in real-world applications) allows to safely take larger step size, thus speeding up the learning process. Nonetheless, increasing the variance over some threshold (making policies nearly random) produces very bad policies, so that changing the policy parameter has a small impact on the performance, and as a result slows down the learning process. How to identify an optimal variance value is an interesting future research direction. Table 2 provides numerical results in the approximated settings, showing the effect of varying the number of trajectories used to estimate the gradient by

REINFORCE and PGT. Increasing the number of trajectories reduces the uncertainty on the gradient estimates, thus allowing to use larger step sizes and reaching better performances. Furthermore, the smaller variance of PGT w.r.t. REINFORCE allows the former to achieve better performances. However, even with a large number of trajectories, the approximated errors are still quite large preventing to reach very high performance. For this reason, future studies will try to derive tighter bounds. Further developments include extending these results to other policy models (e.g., Gibbs policies) and to other policy gradient approaches (e.g., natural gradient). 8

## 2 References

- [1] Jan Peters and Stefan Schaal. Policy gradient methods for robotics. In *Intelligent Robots and Systems*, 2006 IEEE/RSJ International Conference on, pages 2219–2225. IEEE, 2006.
- [2] James C Spall. Multivariate stochastic approximation using a simultaneous perturbation gradient approximation. *Automatic Control, IEEE Transactions on*, 37(3):332–341, 1992.
- [3] Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3-4):229–256, May 1992.
- [4] Jonathan Baxter and Peter L. Bartlett. Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research*, 15:319–350, 2001.
- [5] Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12(22), 2000.
- [6] Jan Peters and Stefan Schaal. Reinforcement learning of motor skills with policy gradients. *Neural Networks*, 21(4):682–697, 2008.
- [7] Sham Kakade. A natural policy gradient. *Advances in neural information processing systems*, 14:1531–1538, 2001.
- [8] Jan Peters and Stefan Schaal. Natural actor-critic. *Neurocomputing*, 71(7):1180–1190, 2008.
- [9] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, pages 400–407, 1951.
- [10] P. Wagner. A reinterpretation of the policy oscillation phenomenon in approximate policy iteration. *Advances in Neural Information Processing Systems*, 24, 2011.
- [11] Jorge J Morfe and David J Thuente. Line search algorithms with guaranteed sufficient decrease. *ACM Transactions on Mathematical Software (TOMS)*, 20(3):286–307, 1994.
- [12] J. Kober and J. Peters. Policy search for motor primitives in robotics. In *Advances in Neural Information Processing Systems 22 (NIPS 2008)*, Cambridge, MA: MIT Press, 2009.
- [13] Nikos Vlassis, Marc Toussaint, Georgios Kontes, and Savas Piperidis. Learning model-free robot control by a monte carlo em algorithm. *Autonomous Robots*, 27(2):123–130, 2009.
- [14] S.M. Kakade. On the sample complexity of reinforcement learning. PhD thesis, PhD thesis, University College London, 2003.
- [15] Matteo Pirodda, Marcello Restelli, Alessio Pecorino, and Daniele Calandriello. Safe policy iteration. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, volume 28, pages 307–315. JMLR Workshop and Conference Proceedings, May 2013.
- [16] S.

Pinsker. Information and Information Stability of Random Variable and Processes. HoldenDay Series in Time Series Analysis. Holden-Day, Inc., 1964. [17]  
Tingting Zhao, Hirotaka Hachiya, Gang Niu, and Masashi Sugiyama. Analysis and improvement of policy gradient estimation. Neural Networks, 26(0):118 ? 129, 2012.