
Balancing Learning Speed and Stability in Policy Gradient via Adaptive Exploration

Matteo Papini
Politecnico di Milano

Andrea Battistello
Politecnico di Milano

Marcello Restelli
Politecnico di Milano

Abstract

In many Reinforcement Learning (RL) applications, the goal is to find an optimal deterministic policy. However, most RL algorithms require the policy to be stochastic in order to avoid instabilities and perform a sufficient amount of exploration. Adjusting the level of stochasticity during the learning process is non-trivial, as it is difficult to assess whether the costs of random exploration will be repaid in the long run, and to contain the risk of instability. We study this problem in the context of policy gradients (PG) with Gaussian policies. Using tools from the safe PG literature, we design a surrogate objective for the policy variance that captures the effects this parameter has on the learning speed and on the quality of the final solution. Furthermore, we provide a way to optimize this objective which guarantees a stable improvement of the original performance measure. We evaluate the proposed methods on simulated continuous control tasks.

1 Introduction

Reinforcement learning (RL, Sutton and Barto, 2018) is an approach to adaptive intelligence that employs a reward signal to train an autonomous agent on a general task through direct interaction with an unknown environment. The results recently achieved by RL in challenging games (Mnih et al., 2015; Silver et al., 2017; OpenAI, 2018) are astounding. However, in order to apply RL to real-world scenarios (e.g., robotics, autonomous driving, finance), we have to tackle further challenges. Unlike games, problems involving

physical systems are more naturally modeled as continuous control tasks. For this reason, we will focus on policy gradient (PG, Sutton et al., 1999; Deisenroth et al., 2013), an RL technique that employs stochastic gradient ascent to optimize parametric controllers. PG is particularly suitable for continuous tasks due to its robustness to noise, convergence properties, and versatility in policy design (Peters et al., 2005). Another perk of games is that they are easily simulated. Simulations require a reliable model of the environment, which is often not available. Learning online on a physical system (like a robot) sharpens the need for *stable* learning algorithms, as large deviations from known policies may yield unsafe behavior or unrecoverable costs.

Although we normally look for a deterministic controller, PG is only able to stably improve *stochastic* policies. A notable exception is Deterministic Policy Gradient (DPG, Silver et al., 2014; Lillicrap et al., 2016), which optimizes a deterministic policy while collecting data with a noisy version of it. Even in this case, the stochastic nature of the behavioral policy is necessary to maintain a sufficient level of exploration and avoid the local optima. This introduces an inevitable trade-off: policy stochasticity facilitates, stabilizes, and speeds up the learning of other policy parameters (Ahmed et al., 2019). On the other hand, random behavior yields worse *online* performance, and may be unsafe in some applications. A natural solution to this problem is to adapt the *amount of exploration* during the learning process. Unfortunately, this is highly non-trivial, as it falls under the infamous exploration-exploitation dilemma. If exploration is abandoned too soon, the agent may never know all the relevant aspects of the task and get stuck in suboptimal behavior. If the transition to deterministic behavior is delayed too much, the learning process may become unnecessarily long and expensive. This problem has been thoroughly studied in the Multi-Armed Bandit (MAB) literature (Bubeck and Cesa-Bianchi, 2012; Lattimore and Szepesvári, 2019). Adaptive exploration has a long history also in RL, mostly limited to the tabular setting (Kearns and Singh, 2002;

Brafman and Tennenholtz, 2002; Strehl et al., 2009; Jaksch et al., 2010; Lattimore and Hutter, 2014; Dann and Brunskill, 2015; Dann et al., 2017; Jin et al., 2018; Ok et al., 2018), with extensions to continuous states (Ortner and Ryabko, 2012; Lakshmanan et al., 2015; Bellemare et al., 2016). In Policy Gradient methods, it is common to learn the amount of stochasticity via gradient ascent on the performance measure, like any other policy parameter (Duan et al., 2016). As randomness typically erodes performance, this *greedy* approach can yield premature convergence to quasi-deterministic policies, causing learning instability and getting stuck to local optima. The current trend is to augment the traditional performance objective (Sutton et al., 1999) with an entropy bonus (Haarnoja et al., 2017, 2018; Nachum et al., 2018; Shani et al., 2018) which favors more stochastic policies. However, the exploration-exploitation trade-off is still unsolved, as one has to decide how much weight to give to the entropy bonus.

In this paper, we propose an alternative approach: we design a separate optimization objective for the policy stochasticity that also accounts for the long-term effects of random exploration. We do so for the special case of Gaussian policies, which are, however, the most used in practice (Duan et al., 2016). In this framework, the amount of exploration can be controlled via the policy variance parameters¹. To do so, we use insights from the *safe policy gradient* literature (Kakade and Langford, 2002; Pirotta et al., 2013, 2015; Papini et al., 2017, 2019), where the effects of policy variance on performance improvement have been already exposed, but not fully exploited. We hence propose the Meta-Exploring Policy Gradient (MEPG) algorithm, that optimizes the variance parameters via gradient ascent on the new surrogate objective, while concurrently optimizing the other parameters as in vanilla PG.

Although built upon insights from the safe PG literature, the MEPG algorithm is heuristic and comes with no formal guarantees. By developing an adaptive meta-parameter schedule, we devise the Stably Exploring Policy Gradient (SEPG) algorithm, a variant of MEPG with guarantees of monotonic improvement of the original performance objective. To do so, we generalize existing improvement guarantees (Papini et al., 2019) for Gaussian policies to the previously uncharted case of adaptive policy variance. These improvement guarantees come at the cost of worsening learning speed and sample complexity, which can be

critical in RL applications, where collecting samples is costly and time-consuming. Although the scope of this work is mostly theoretical, for the sake of applicability, we also relax the classical monotonic improvement constraint (Kakade and Langford, 2002) to a more general *bounded worsening* constraint. This allows the practitioner to specify how much performance he would accept to lose (and with which probability) at any policy update, introducing a way to trade-off the stability of learning with the time required.

The paper is structured as follows: in Section 2, we provide an essential background on PG and review the existing performance improvement guarantees for this framework. In Section 3, we present our novel exploration objective and the MEPG algorithm. In Section 4.1, we extend existing improvement guarantees for Gaussian policies to the adaptive-variance case, providing safe exact-gradient updates. In section 4.2, we generalize these results to the more realistic stochastic-gradient case, and present the SEPG algorithm. Finally, in Section 5, we evaluate the proposed algorithms on simulated control tasks. Proofs of all the formal statements are given in Appendix A.

2 Preliminaries

In this section, we provide an essential background on policy gradient methods, including existing performance-improvement guarantees.

2.1 Policy Gradient Fundamentals

A continuous Markov Decision Process (MDP, Puterman, 1994) $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \rho \rangle$ is defined by a continuous state space $\mathcal{S} \subseteq \mathbb{R}^d$; a continuous action space \mathcal{A} ; a Markovian transition kernel \mathcal{P} , where $\mathcal{P}(s'|s, a)$ is the transition density from state s to s' under action a ; a reward function \mathcal{R} , where $\mathcal{R}(s, a) \in [-R_{\max}, R_{\max}]$ is the reward for state-action pair (s, a) and R_{\max} is the maximum absolute-value reward; a discount factor $\gamma \in [0, 1)$; and an initial-state distribution ρ on \mathcal{S} . An agent’s behavior is modeled as a policy π , where $\pi(\cdot|s)$ is the density function over \mathcal{A} in state s . We study episodic MDPs with indefinite horizon. In practice, we consider episodes of length H , the effective horizon of the task. A trajectory τ is a sequence of states and actions $(s_0, a_0, s_1, a_1, \dots, s_{H-1}, a_{H-1})$ observed by following a stationary policy, where $s_0 \sim \rho$ and $s_{h+1} \sim \mathcal{P}(\cdot|s_h, a_h)$. The policy induces a measure p_π over trajectories. We denote with $\mathcal{R}(\tau)$ the total discounted reward provided by trajectory τ : $\mathcal{R}(\tau) = \sum_{h=0}^{H-1} \gamma^h \mathcal{R}(s_h, a_h)$. Policies can be ranked based on their expected total reward $J(\pi) = \mathbb{E}_{\tau \sim p_\pi} [\mathcal{R}(\tau)]$. Solving the MDP means finding an optimal policy $\pi^* \in \arg \max_{\pi} \{J(\pi)\}$.

¹Extensions to other classes of policies are possible, but require to identify explicit scale parameters e.g., the *temperature* for Softmax policies or the *variance* for (reparametrized) beta policies.

Policy gradient (PG, Sutton et al., 1999; Peters and Schaal, 2008) methods restrict this optimization problem to a class of parametric policies $\Pi_\Theta = \{\pi_\theta : \theta \in \Theta \subseteq \mathbb{R}^m\}$, so that π_θ is differentiable w.r.t. θ . We denote the performance of a parametric policy π_θ with $J(\theta)$. A *locally* optimal policy can be found via gradient ascent on the performance measure:

$$\theta_{t+1} \leftarrow \theta_t + \alpha_t \nabla_\theta J(\theta_t),$$

where $\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim p_\theta} [\nabla_\theta \log p_\theta(\tau) \mathcal{R}(\tau)]$, (1)

where t denotes the current iteration, p_θ is short for p_{π_θ} , and $(\alpha_t)_{t=1}^\infty$ is a sequence of positive step sizes. In practice, $\nabla_\theta J$ is not available, but can be estimated from a batch of trajectories $\mathcal{D}_N = \{\tau_1, \dots, \tau_N\}$. The GPOMDP (Baxter and Bartlett, 2001) algorithm (a refinement of REINFORCE, Williams, 1992) provides an unbiased gradient estimator:

$$\hat{\nabla}_\theta^N J(\theta) = \frac{1}{N} \sum_{n=1}^N \sum_{h=0}^{H-1} \left(\sum_{i=0}^h \nabla_\theta \log \pi_\theta(a_i^n | s_i^n) \right) \cdot (\gamma^h \mathcal{R}(s_h^n, a_h^n) - b), \quad (2)$$

where b is a baseline used to reduce variance. Any baseline that does not depend on actions preserves the unbiasedness of the estimator.² We employ the variance-minimizing baselines provided by Peters and Schaal (2008).

A widely used (Duan et al., 2016) policy class is the Gaussian³:

$$\pi_\theta(a|s) = \frac{1}{\sqrt{2\pi}\sigma_\omega} \exp \left\{ -\frac{1}{2} \left(\frac{a - \mu_\nu(s)}{\sigma_\omega} \right)^2 \right\}, \quad (3)$$

also denoted with $\mathcal{N}(a|\mu_\nu(s), \sigma_\omega^2)$, where the action space is $\mathcal{A} = \mathbb{R}$, μ_ν is the state-dependent mean and $\sigma_\omega^2 > 0$ is the variance (σ_ω is the standard deviation). The policy parameters consist of a vector of mean parameters $\nu \in \Upsilon \subseteq \mathbb{R}^m$ and a variance parameter $\omega \in \Omega \subseteq \mathbb{R}$, i.e., $\theta \equiv [\nu^T | \omega]^T$ and $\Theta \equiv \Upsilon \times \Omega \subseteq \mathbb{R}^{m+1}$. We focus on the following, common (Rajeswaran et al., 2017) parametrization:

$$\mu_\nu(s) = \nu^T \phi(s), \quad \sigma_\omega = e^\omega, \quad (4)$$

where $\phi(\cdot)$ is a vector of m state-features⁴. We also assume that the state features are bounded in Euclidean norm, i.e., $\sup_{s \in \mathcal{S}} \|\phi(s)\| \leq \varphi$, and both Υ and Ω are convex sets.

²Some action-dependent baselines are also possible. See (Tucker et al., 2018) for a discussion.

³We consider scalar actions for simplicity. Multi-dimensional actions are discussed in Appendix D, together with heteroskedastic exploration.

⁴This is a shallow policy since we assume to have the features already available, as opposed to a deep policy, where the features are learned together with ν in an end-to-end fashion. Generalizations to deep neural policies are plausible, but beyond the scope of this paper.

Entropy regularization (Schulman et al., 2017) consists in modifying the reward to favor policy stochasticity:

$$\tilde{r}_t = (1 - \tau)r_t + \tau \mathbb{H}(\pi(\cdot|s_t)), \quad (5)$$

where $\mathbb{H}(\pi(\cdot|s_t)) = \mathbb{E}_{a \sim \pi(\cdot|s_t)} [-\log \pi(a|s_t)]$ is the entropy of the action distribution at state s_t and $\tau \in [0, 1]$ is a regularization coefficient. In the case of Gaussian policies, the entropy bonus can be computed in closed form and the GPOMDP algorithm can be easily modified to account for the new reward function (e.g., Ahmed et al., 2019). Vanilla policy gradient is recovered for $\tau = 0$.

2.2 Safe Policy Gradients

When learning on-line in the real world, we would like to avoid oscillations in the performance $J(\theta)$, as large deviations from previously observed behavior may compromise the safety of the system or deemed unacceptable by stakeholders. A convenient way to enforce this desideratum is through an *improvement constraint* (Thomas et al., 2015):

Definition 2.1. Given a parametric policy π_θ with current parameter θ_t , we say that update $\Delta\theta \in \mathbb{R}^{m+1}$ is *safe* w.r.t. requirement $C_t \in \mathbb{R}$ if:

$$J(\theta_{t+1}) - J(\theta_t) \geq C_t, \quad (6)$$

where $\theta_{t+1} = \theta_t + \Delta\theta$.

We talk of a *required performance improvement* when C_t is non-negative, otherwise of a *bounded worsening*. The case $C_t \equiv 0$ corresponds to the well-studied *monotonic improvement* constraint (Kakade and Langford, 2002; Pirotta et al., 2013; Papini et al., 2017). This constraint can also be used to enforce an absolute performance threshold J_{\min} , by setting $C_t = J_{\min} - J(\theta_t)$. For instance, if we want to guarantee that the performance is never worse than that of the initial policy, $J(\theta_0)$, we set $C_t = J(\theta_0) - J(\theta_t)$. This can be useful in safety-critical systems where the initial policy is designed to be safe, assuming the performance measure captures all sources of risk (García and Fernández, 2015; Amodei et al., 2016). Recent work (Papini et al., 2019) provides improvement guarantees for a general family of policies. Given positive constants ψ , κ and ξ , a policy class Π_Θ is called (ψ, κ, ξ) -smoothing if:

$$\begin{aligned} \sup_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \left[\|\nabla \log \pi_\theta(a|s)\| \right] &\leq \psi, \\ \sup_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \left[\|\nabla \log \pi_\theta(a|s)\|^2 \right] &\leq \kappa, \\ \sup_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \left[\|\nabla \nabla^T \log \pi_\theta(a|s)\| \right] &\leq \xi, \end{aligned} \quad (7)$$

for all $\theta \in \Theta$, where $\|\cdot\|$ denotes the Euclidean norm for vectors and the spectral norm for matrices. In particular, Gaussian policies *with fixed standard deviation*

(constant ω) are (ψ, κ, ξ) -smoothing with the following constants (Papini et al., 2019):

$$\psi = \frac{2\varphi}{\sqrt{2\pi}\sigma_\omega}, \quad \kappa = \xi = \frac{\varphi^2}{\sigma_\omega^2}, \quad (8)$$

where φ is the Euclidean-norm bound on state features. For a (ψ, κ, ξ) -smoothing policy, the performance improvement yielded by a policy gradient update can be lower-bounded by a function of the step size α as follows (Theorem 9 from Papini et al., 2019):

$$J(\theta_{t+1}) - J(\theta_t) \geq \alpha \|\nabla J(\theta_t)\|^2 - \alpha^2 \frac{L}{2} \|\Delta\theta\|^2, \quad (9)$$

where $\theta_{t+1} = \theta_t + \alpha \nabla J(\theta_t)$ and $L = \frac{R_{\max}}{(1-\gamma)^2} \left(\frac{2\gamma\psi^2}{1-\gamma} + \kappa + \xi \right)$. This allows to select a *safe step size* for improvement constraint C_t , i.e., a step size for which (6) is satisfied by the policy gradient update (1). In this paper, whenever multiple choices of the step size are safe, we decide to employ the *largest* safe step size. This is meant to yield faster convergence (see Section 5 for an empirical substantiation of this claim). In the fixed-variance Gaussian case, we can obtain a safe step size for the mean-parameter update (adaptation of Corollary 10 from Papini et al., 2019):

Lemma 2.1. *Let Π_γ be the class of Gaussian policies parametrized as in (4), but with fixed variance parameter ω . Let $\mathbf{v}_t \in \mathbb{R}^m$ and $\mathbf{v}_{t+1} = \mathbf{v}_t + \alpha_t \nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega)$. For any $C_t \leq C_t^*$, the largest step size guaranteeing $J(\mathbf{v}_{t+1}, \omega) - J(\mathbf{v}_t, \omega) \geq C_t$ is:*

$$\bar{\alpha}_t := \frac{\sigma_\omega^2}{F} \left(1 + \sqrt{1 - \frac{C_t}{C_t^*}} \right), \quad (10)$$

where $F = \frac{2\varphi^2 R_{\max}}{(1-\gamma)^2} \left(1 + \frac{2\gamma}{\pi(1-\gamma)} \right)$ and $C_t^* = \frac{\sigma_\omega^2 \|\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega)\|^2}{2F}$. \square

We have highlighted the role of the policy standard deviation. We can see how a larger σ_ω allows to take larger steps. Moreover, it increases the maximum improvement guarantee that one *can ask for* (although C_t can be selected by the user at each step, its highest feasible value C_t^* depends on the current policy variance and gradient norm). In fact, both $\bar{\alpha}_t$ and C_t^* are $\mathcal{O}(\sigma_\omega^2)$. This is due to the smoothing effect of the policy variance on the optimization landscape, in accordance with the empirical analysis from (Ahmed et al., 2019). In practice, C_t^* is too small to be of any relevance, so we are more interested in the cases when $C_t \leq 0$.

3 Adaptive Exploration

In this section, we use some insights from the safe PG literature to devise a heuristic approach to adapt the

standard deviation of a Gaussian policy during the learning process. Our desiderata are fast convergence, avoiding instabilities and not getting stuck in local optima. The algorithm we present here is heuristic. A variant with formal improvement guarantees is presented in the next section.

Consider a Gaussian policy $\pi_{\mathbf{v}, \omega}(a|s) = \mathcal{N}(a|\mu_{\mathbf{v}}(s), \sigma_\omega)$, parametrized as in (4). As mentioned above, it is common to learn the policy variance parameter via gradient ascent just like any other parameter, i.e., $\omega_{t+1} \leftarrow \omega_t + \beta_t \nabla_\omega J(\mathbf{v}_t, \omega_t)$. However, the effects of σ on the optimization landscape, exposed by Lemma 2.1, suggest to treat it with particular care, both to exploit its potential and to avoid its possible risks. In fact, adjusting the policy variance with policy gradient tends to degenerate too early into quasi-deterministic policies, getting stuck in local optima or even causing divergence issues (see Section 5). We use our understanding of the special nature of this parameter to modify GPOMDP in two ways. First of all, we make the step size *for updating the mean parameters* dependent on the policy variance, like the safe step size from Lemma 2.1. In particular, we use the following:

$$\alpha_t = \frac{\alpha \sigma_{\omega_t}^2}{\|\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)\|}, \quad (11)$$

to update the mean parameters \mathbf{v} , where $\alpha > 0$ is a hyper-parameter. This has both the effect of reducing the step size when a small σ makes the optimization landscape less smooth, preventing oscillations, and increasing it when a large σ allows it to do so, increasing the learning speed. This is not entirely unheard of, as it is exactly what a natural gradient (Kakade, 2001; Amari, 1998) would do in a *pure* Gaussian setting ($1/\sigma^2$ is the Fisher information w.r.t. the mean parameters of a Gaussian distribution, see Sehnke et al., 2008; Miyamae et al., 2010). We also divide the step size by the norm of the gradient. This is a common normalization technique (Peters et al., 2005), and is further motivated by the results of Section 4.2 on stochastic gradient updates.

As for the variance parameter ω , we treat it as a separate *meta-parameter* and we learn it in a meta-gradient fashion (Sutton, 1992; Schraudolph, 1999; Veeriah et al., 2017; Xu et al., 2018). Specifically, we employ a more far-sighted learning objective to avoid premature convergence to deterministic behavior. To do so, we look at the target performance one step in the future:

$$\begin{aligned} & J \left(\mathbf{v}_t + \alpha \sigma_{\omega_t}^2 \frac{\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)}{\|\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)\|}, \omega_t \right) \\ & \simeq J(\mathbf{v}_t, \omega_t) + \alpha \sigma_{\omega_t}^2 \|\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)\| := \mathcal{L}(\mathbf{v}_t, \omega_t), \end{aligned} \quad (12)$$

where we performed a first-order approximation. The

Algorithm 1 MEPG

```

1: Input: Initial parameters  $\mathbf{v}_0$  and  $\omega_0$ , step size  $\alpha > 0$ , meta step size  $\eta > 0$ , batch size  $N$ 
2: for  $t = 1, 2, \dots$  do
3:   Collect a batch of  $N$  trajectories with  $\pi_{\mathbf{v}_t, \omega_t}$ 
4:   Estimate  $\hat{\nabla}_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)$ ,  $\hat{\nabla}_{\omega} J(\mathbf{v}_t, \omega_t)$  and  $\hat{\nabla}_{\omega} \|\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)\|$  ▷ See Appendix B
5:    $\hat{\nabla}_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t) = \hat{\nabla}_{\omega} J(\mathbf{v}_t, \omega_t) + \alpha e^{2\omega_t} \left( 2 \|\hat{\nabla}_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)\| + \hat{\nabla}_{\omega} \|\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)\| \right)$ 
6:    $\omega_{t+1} \leftarrow \omega_t + \eta \hat{\nabla}_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t) / \|\hat{\nabla}_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t)\|$ 
7:    $\mathbf{v}_{t+1} \leftarrow \mathbf{v}_t + \alpha e^{2\omega_t} \hat{\nabla}_{\mathbf{v}} J(\mathbf{v}_t, \omega_t) / \|\hat{\nabla}_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)\|$ 
8: end for
    
```

gradient of \mathcal{L} w.r.t. ω is:

$$\nabla_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t) = \nabla_{\omega} J(\mathbf{v}_t, \omega_t) + 2\alpha\sigma_{\omega_t}^2 \|\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)\| + \alpha\sigma_{\omega_t}^2 \nabla_{\omega} \|\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)\|. \quad (13)$$

The first term of the sum is the usual policy gradient w.r.t. ω , and accounts for the immediate effect of policy stochasticity on performance. The role of the second term is to increase the step size α_t , more so if the gradient w.r.t. \mathbf{v} is large. The third term is meant to modify the policy variance to increase the gradient norm and can be seen as a way to escape local optima. The last two terms, together, account for the long-term effects of modifying the policy variance. We propose to update ω in the direction of the (normalized) meta-gradient $\nabla_{\omega} \mathcal{L}$ using a meta-step size $\eta > 0$:

$$\omega_{t+1} \leftarrow \omega_t + \eta \frac{\nabla_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t)}{\|\nabla_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t)\|}. \quad (14)$$

In practice, exact gradients are not available. The policy gradient for the mean-parameter update can be estimated with GPOMDP (2). Computing $\hat{\nabla}_{\omega} \mathcal{L}$ also requires estimating $\nabla_{\omega} \|\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)\|$, which is computationally no more expensive, but could suffer from more variance (See Appendix B). The pseudocode for the resulting algorithm, called Meta-Exploring Policy Gradient (MEPG), is provided in Algorithm 1.

4 Stable Exploration

In this section, we extend the performance improvement guarantees reported in Section 2.2 to Gaussian policies with adaptive variance and we use these theoretical results to devise a variant of Algorithm 1 with improvement guarantees.

4.1 Exact framework

Existing guarantees for fixed-variance Gaussian policies are based on the smoothing constants from (Papini et al., 2019), reported in (8). These depend (inversely) on σ_{ω}^2 , hence are no longer constant once we allow ω to vary. In particular, they tend to infinity as the policy approaches determinism. Unfortu-

nately, this is enough to invalidate the safety guarantees. A workaround would be to replace σ_{ω} with a lower bound, which can be imposed by constraining the parameter space Ω or by changing the parametrization. However, this would make the improvement bounds unnecessarily conservative, and would prevent the agent to converge to deterministic behavior in the end. For these reasons, we instead propose to update the mean and variance parameters *alternately*. First, we show that Gaussian policies are smoothing w.r.t. the variance parameter *alone*:

Lemma 4.1. *Let Π_{Ω} be the class of Gaussian policies parametrized as in (4), but with fixed mean parameter \mathbf{v} . Π_{Ω} is $\left(\frac{4}{\sqrt{2\pi e}}, 2, 2\right)$ -smoothing. \square*

This allows to devise a safe policy-gradient update for the variance parameters:

Theorem 4.2. *Let Π_{Ω} be the class of policies defined in Lemma 4.1. Let $\omega_t \in \Omega$ and $\omega_{t+1} \leftarrow \omega_t + \beta_t \nabla_{\omega} J(\mathbf{v}, \omega_t)$. For any $C_t \leq C_t^*$, the largest step-size satisfying (6) is:*

$$\bar{\beta}_t = \frac{1}{G} \left(1 + \sqrt{1 - \frac{C_t}{C_t^*}} \right), \quad (15)$$

where $C_t^* = \frac{\|\nabla_{\omega} J(\mathbf{v}, \omega_t)\|^2}{2G}$ and $G = \frac{4R_{\max}}{(1-\gamma)^2} \left(1 + \frac{4\gamma}{\pi e(1-\gamma)} \right)$. \square

Similarly to the mean parameters case, $\beta_t = \frac{1}{G}$ is the greedy-safe step size and $\beta_t = \frac{2}{G}$ is the largest step size guaranteeing monotonic improvement.

Alternately updating the mean parameter as in Lemma 2.1 and the variance parameter as in Theorem 4.2 ensures $J(\mathbf{v}_{t+1}, \omega_{t+1}) - J(\mathbf{v}_t, \omega_t) \geq C_t$ for all t .

However, Theorem 4.2 still pertains *naïve* variance updates, which suffer from all the problems discussed in Section 3. The next question is how to optimize the surrogate exploratory objective \mathcal{L} from (12) while satisfying the original constraint (6) on the performance objective J . The following Theorem provides a safe

Algorithm 2 SEPG

```

1: Input: Initial parameters  $\mathbf{v}_0$  and  $\omega_0$ , batch size  $N$ , improvement thresholds  $\{C_t\}_{t=1}^\infty$ , confidence parameter
    $\delta$ , discount factor  $\gamma$ , maximum reward  $R_{\max}$ , feature bound  $\varphi$ 
2: for  $t = 1, 2, \dots$  do
3:   Collect a batch of  $N$  trajectories with  $\pi_{\mathbf{v}_t, \omega_t}$ 
4:   Estimate  $\hat{\nabla}_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)$ ,
5:   if  $t$  is odd then
6:      $\mathbf{v}_{t+1} = \mathbf{v}_t + \tilde{\alpha}_t \nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)$  ▷ Safe step size from Section 4.2
7:   else
8:     estimate  $\hat{\nabla}_{\omega} J(\mathbf{v}_t, \omega_t)$  and  $\hat{\nabla}_{\omega} \|\hat{\nabla}_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)\|$  ▷ See Appendix B
9:      $\hat{\nabla}_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t) = \hat{\nabla}_{\omega} J(\mathbf{v}_t, \omega_t) + \tilde{\alpha}_t \left( 2 \|\hat{\nabla}_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)\| + \hat{\nabla}_{\omega} \|\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega_t)\| \right)$ 
10:     $\omega_{t+1} = \omega_t + \tilde{\eta}_t \nabla_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t)$  ▷ Safe meta-step size from Section 4.2
11:   end if
12: end for
    
```

update for a smoothing policy in the direction of a generic update vector \mathbf{x}_t ($\nabla_{\omega} \mathcal{L}$ in MEPG):

Theorem 4.3. *Let Π_{Θ} be a (ψ, κ, ξ) -smoothing policy class, $\boldsymbol{\theta}_t \in \Theta$, and $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \eta_t \mathbf{x}_t$, where $\mathbf{x}_t \in \mathbb{R}^m$ and $\eta_t \in \mathbb{R}$ is a (possibly negative) step size. Let $\lambda_t := \frac{\langle \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_t), \mathbf{x}_t \rangle}{\|\mathbf{x}_t\|}$ be the scalar projection of $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_t)$ onto \mathbf{x}_t . For any $C_t \leq C_t^*$, provided $\lambda_t \neq 0$, the largest step size guaranteeing $J(\boldsymbol{\theta}_{t+1}) - J(\boldsymbol{\theta}_t) \geq C_t$ is:*

$$\bar{\eta}_t = \frac{|\lambda_t|}{L \|\mathbf{x}_t\|} \left(\text{sign}(\lambda_t) + \sqrt{1 - \frac{C_t}{C_t^*}} \right), \quad (16)$$

where $C_t^* = \frac{\lambda_t^2}{2L}$ and $L = \frac{R_{\max}}{(1-\gamma)^2} \left(\frac{2\gamma\psi^2}{1-\gamma} + \kappa + \xi \right)$. \square

Note that a positive performance improvement up to C_t^* can always be guaranteed, even if the improvement direction $\nabla_{\boldsymbol{\theta}} J$ is not explicitly followed. However, when the scalar projection λ_t is negative, the largest safe step size is negative. This corresponds to the case in which maximizing the surrogate objective *reduces* the original one. For positive values of C_t (required improvement), there may be no way to safely pursue the surrogate objective. In this case, a negative step size is prescribed to follow the direction of $\nabla_{\boldsymbol{\theta}} J$ instead⁵. We can use the step size $\bar{\eta}_t$ from Theorem 4.3 to safely replace $\nabla_{\omega} J$ with the meta gradient $\nabla_{\omega} \mathcal{L}$ from (13) in the variance update:

$$\mathbf{v}_{t+1} \leftarrow \mathbf{v}_t + \bar{\alpha}_t \nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega_t), \quad (17)$$

$$\omega_{t+2} \leftarrow \omega_{t+1} + \bar{\eta}_{t+1} \nabla_{\omega} \mathcal{L}(\mathbf{v}_{t+1}, \omega_{t+1}), \quad (18)$$

where $\omega_{t+1} \equiv \omega_t$ and $\mathbf{v}_{t+2} \equiv \mathbf{v}_{t+1}$, as the two set of parameters cannot be safely updated together.

4.2 Approximate framework

In practice, exact gradients are not available and must be estimated from data. In this section, we show how to adapt the safe step sizes from Section 4 to take gradient estimation errors into account. Let $\hat{\nabla}_{\mathbf{v}}^N J$, $\hat{\nabla}_{\omega}^N J$ and $\hat{\nabla}_{\omega}^N \mathcal{L}$ be unbiased estimators of $\nabla_{\mathbf{v}} J$, $\nabla_{\omega} J$ and $\nabla_{\omega} \mathcal{L}$, respectively, each using a batch of N trajectories. As for MEPG, the first two can be GPOMDP estimators (2) and meta-gradient estimation is discussed in Appendix B. We make the following assumption on the gradient estimators⁶:

Assumption 4.4. *For every $\delta \in (0, 1)$ there exists a non-negative constant ϵ_{δ} such that, with probability at least $1 - \delta$:*

$$\begin{aligned} \|\nabla_{\mathbf{v}} J(\mathbf{v}, \omega) - \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}, \omega)\| &\leq \frac{\epsilon_{\delta}}{\sqrt{N}}, \\ \|\nabla_{\omega} J(\mathbf{v}, \omega) - \hat{\nabla}_{\omega}^N J(\mathbf{v}, \omega)\| &\leq \frac{\epsilon_{\delta}}{\sqrt{N}}, \end{aligned}$$

for every $\mathbf{v} \in \Upsilon$, $\omega \in \Omega$ and $N \geq 1$. \square

Here ϵ_{δ} represents an upper bound on the gradient estimation error. This can be characterized using various statistical inequalities (Papini et al., 2017). A possible one, based on ellipsoidal confidence regions, is described in Appendix C. Under Assumption 4.4, for a sufficiently large batch size, the safe step size for the mean update can be adjusted as follows:

$$\tilde{\alpha}_t = \frac{\sigma_{\omega_t}^2 \left(\left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega_t) \right\| - \frac{\epsilon_{\delta}}{\sqrt{N}} \right)}{F \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega_t) \right\|} \left(1 + \sqrt{1 - \frac{C_t}{C_t^*}} \right),$$

where $C_t^* = \frac{\sigma_{\omega_t}^2 \left(\left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega_t) \right\| - \frac{\epsilon_{\delta}}{\sqrt{N}} \right)^2}{2F}$ and F is from Lemma 2.1. Similarly, the safe step size for the vari-

⁵ The special case when the two gradients are orthogonal ($\lambda_t = 0$) is discussed in Appendix A.

⁶We do not need a similar assumption on the meta-gradient estimator $\hat{\nabla}_{\omega} \mathcal{L}$, since our improvement requirements are always on the performance J (see Appendix A.2).

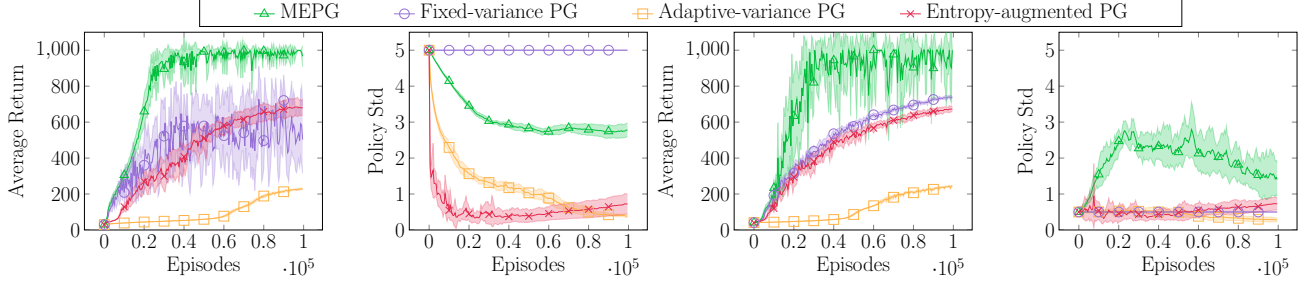


Figure 1: Average return (undiscounted) and policy standard deviation per episode of MEPG, fixed-variance PG, adaptive-variance PG and entropy-augmented PG on the continuous Cart-Pole task, starting from $\sigma = 5$ (left) and $\sigma = 0.5$ (right); averaged over 10 independent runs with 95% Student’s t-confidence intervals.

ance update can be adjusted as follows:

$$\tilde{\eta}_t = \frac{|\hat{\lambda}_t| - \frac{\epsilon_\delta}{\sqrt{N}}}{G \|\hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}_t, \omega_t)\|} \left(\text{sign}(\hat{\lambda}_t) + \sqrt{1 - \frac{C_t}{C_t^*}} \right),$$

where $\hat{\lambda}_t$ is the scalar projection of $\hat{\nabla}_\omega^N J(\mathbf{v}_t, \omega_t)$ onto $\hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}_t, \omega_t)$, $C_t^* = \frac{(|\hat{\lambda}_t| - \frac{\epsilon_\delta}{\sqrt{N}})^2}{2G}$ and G is from Theorem 4.2. We call the variant of MEPG that alternates mean and variance updates using these step sizes Stably Exploring Policy Gradient (SEPG), detailed in Algorithm 2. The SEPG algorithm satisfies our bounded-worsening constraint (6) with high probability:

Theorem 4.5. *Under Assumption 4.4, provided $C_t \leq 0$ and $N > \epsilon_\delta^2 / \hat{\lambda}_t^2$, Algorithm 2 guarantees $J(\theta_{t+1}) - J(\theta_t) \geq C_t$ with probability at least $1 - \delta$ for any $t \geq 1$. \square*

In Appendix F we prove a stronger version of this theorem that also allows strictly positive improvements. The requirement on the batch size ensures that estimation errors are smaller than the estimates themselves. If this requirement is not satisfied, we can either collect more samples or terminate. This typically happens close to stationary points anyway. The step sizes proposed by SEPG may be excessively conservative. This is similar to what happens in supervised learning, where the convergence-guaranteeing step sizes are rarely used in practice, typically replaced by heuristics (Kingma and Ba, 2015). However, since policy updates in online RL can have concrete consequences, we will use the prescribed step sizes in our experiments, leaving the possibility of explicitly relaxing the safety requirements.

5 Experiments

In this section, we test the proposed methods on simulated continuous control tasks. More details on the simulation environments are provided in Appendix E.

MEPG We test MEPG on a continuous-action version of the Cart-Pole balancing task (Barto et al., 1983). Figure 1 shows the performance (1000 is the maximum) and the policy standard deviation of MEPG and three versions of PG (with the GPOMDP gradient estimator). In fixed-variance PG, the policy variance parameter is kept constant. In adaptive-variance PG, it is learned via gradient ascent as any other parameter. Entropy-augmented PG is the same, but with entropy regularization (5). For each algorithm, the best hyper-parameters (step sizes and entropy coefficient τ) have been selected by grid search (see Appendix F). Two very different initializations are considered for the standard deviation: $\sigma_{\omega_0} = 5$ (on the left of Figure 1) and $\sigma_{\omega_0} = 0.5$ (on the right). As shown by the behavior of fixed-variance GPOMDP, the former constant value is too large to achieve optimal performance at convergence, while the latter is too small to properly explore the environment. As expected, adaptive-variance GPOMDP is too greedy and ends up always reducing the standard deviation. Besides preventing exploration, divergence issues force us to use a smaller step size ($\alpha = 0.01$ instead of 0.1), resulting in slower learning. This problem is fixed by the entropy bonus, which prevents the policy from becoming deterministic and allows to use the larger learning rate ($\alpha = 0.1$). However, entropy-augmented PG does not perform significantly better than its fixed-variance counterpart on this task, as the amount of exploration needed to find the global optimum is not maintained (or is pursued too late). Instead, MEPG is able to settle on an intermediate value with both variance initializations. This allows both to learn faster and to achieve optimal performance, although non-negligible oscillations can be observed. These oscillations are partly due to the variance of the meta-gradient estimator, and can be mitigated by the conservative step sizes prescribed by SEPG.

SEPG Figure 2 shows the performance and the policy standard deviation of SEPG on the one-

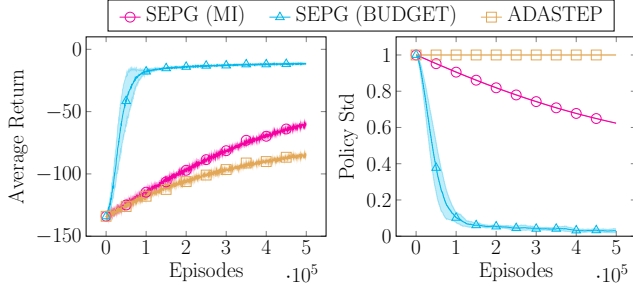


Figure 2: Average return (undiscounted) and standard deviation per episode of SEPG and ADASTEP on the LQG task, averaged over 10 independent runs with 95% Student’s t-confidence intervals.

dimensional LQG (Linear-Quadratic Gaussian regulator) task (Peters and Schaal, 2008). SEPG with a monotonic improvement constraint ($C_t \equiv 0$) is compared with the adaptive-step-size algorithm (ADASTEP in the figure) by Pirodda et al. (2013). Starting from $\sigma_{\omega_0} = 1$, SEPG achieves higher returns by safely lowering it, while ADASTEP has no way to safely update this parameter. Both algorithms use $\delta = 0.2$ and a large batch size ($N = 500$). We also consider a looser constraint, already discussed in Section 2.2 (BUDGET in the figure): that of never doing worse than the initial performance ($C_t = J(\theta_0) - J(\theta_t)$). As expected, this allows faster learning, leading to optimal performance within a reasonable time.

On the Cart-Pole task, MEPG showed an oscillatory behavior. Motivated by this fact, we run SEPG with a fixed, *negative* improvement threshold C_t . Recall that the meaning of such a constraint is to limit per-update performance worsening. Figure 3 shows the results for different values of the threshold, starting from $\sigma_{\omega_0} = 5$ and neglecting the gradient estimation error (i.e., by setting $\delta = 1$). Even under this simplifying assumption, only a very large value of C_t allows to reach optimal performance within a reasonable time. This is due to the over-conservativeness of the step sizes proposed by SEPG, as already discussed. Note how oscillations are reduced w.r.t. MEPG, and how policy standard deviation is first reduced and then *safely increased* again.

6 Discussion and Future Work

We have highlighted the special role of stochasticity in PG with Gaussian policies, complementing the empirical observations from (Ahmed et al., 2019) with theoretical insights from the Safe PG literature (Papini et al., 2019). We have proposed a variant of GPOMDP for this setting, called Meta-Exploring Policy Gradient (MEPG), which is able to adapt the vari-

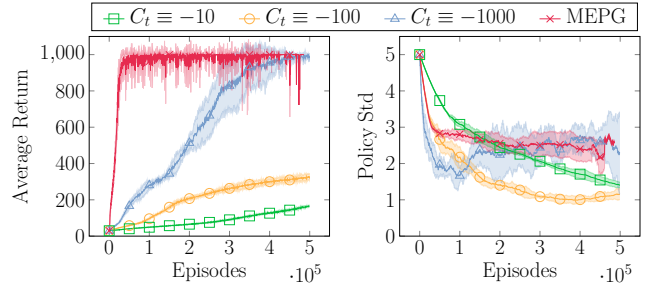


Figure 3: Average return (undiscounted) and standard deviation per episode of SEPG on the Cart-Pole task for different values of C_t , averaged over 5 runs with 95% confidence intervals. The learning curve of MEPG is reported as a reference.

ance parameter in a more far-sighted way than vanilla or entropy-augmented gradient ascent. We have empirically shown the effectiveness of this approach on the Cart-Pole balancing task, where the entropy bonus is able to prevent divergence but not to escape local optima. This should be intended as a proof of concept: entropy augmentation is still a natural choice for most applications due to its simplicity. Future work on MEPG algorithms should study its applicability to larger control problems and overcome its potential bottlenecks, such as the second-order term in (12).

Furthermore, we have generalized the existing performance-improvement bounds for Gaussian policies to the adaptive-variance case and proposed SEPG, a variant of MEPG with guarantees of stable improvement. Experiments confirmed several intuitions provided by the theory. Unfortunately, learning speed is heavily degraded for meaningful values of the improvement requirement C_t , due to over-conservative step sizes. The desired balance between speed and stability can still be achieved by hand-tuning it as a hyper-parameter. Replacing the theoretically sound upper bounds with estimates (Allen-Zhu, 2018) could bridge the gap between theory and practice. Another possible improvement is to employ an adaptive batch size as proposed in (Papini et al., 2017, 2019). Future work should also study in more depth the issue of local optima. Recent work (Agarwal et al., 2019) shows how (natural) PG can achieve global optimality in some cases, and highlights the importance of exploration on this matter. Moreover, adaptively changing the policy variance as in MEPG can be seen as an online version of graduated optimization (Hazan et al., 2016), a popular technique for finding the global optimum of a nonconvex function. Finally, further aspects of *safe exploration* (Hans et al., 2008; Turchetta et al., 2016; Cohen et al., 2018) should be considered for the online selection of policy stochasticity in real-world control tasks.

Acknowledgements

The study was partially funded by Lombardy Region (Announcement PORFESR 2014-2020).

References

- Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. (2019). Optimality and approximation with policy gradient methods in markov decision processes. *CoRR*, abs/1908.00261.
- Ahmed, Z., Roux, N. L., Norouzi, M., and Schuurmans, D. (2019). Understanding the impact of entropy on policy optimization. In *ICML*, volume 97 of *Proceedings of Machine Learning Research*, pages 151–160. PMLR.
- Allen-Zhu, Z. (2018). Natasha 2: Faster non-convex optimization than SGD. In *NeurIPS*, pages 2680–2691.
- Amari, S. (1998). Natural gradient works efficiently in learning. *Neural Computation*, 10(2):251–276.
- Amodi, D., Olah, C., Steinhardt, J., Christiano, P. F., Schulman, J., and Mané, D. (2016). Concrete problems in AI safety. *CoRR*, abs/1606.06565.
- Barto, A. G., Sutton, R. S., and Anderson, C. W. (1983). Neuronlike adaptive elements that can solve difficult learning control problems. *IEEE Trans. Systems, Man, and Cybernetics*, 13(5):834–846.
- Baxter, J. and Bartlett, P. L. (2001). Infinite-horizon policy-gradient estimation. *J. Artif. Intell. Res.*, 15:319–350.
- Bellemare, M. G., Srinivasan, S., Ostrovski, G., Schaul, T., Saxton, D., and Munos, R. (2016). Unifying count-based exploration and intrinsic motivation. In *NeurIPS*, pages 1471–1479.
- Brafman, R. I. and Tennenholtz, M. (2002). R-MAX - A general polynomial time algorithm for near-optimal reinforcement learning. *J. Mach. Learn. Res.*, 3:213–231.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. (2016). Openai gym. *CoRR*, abs/1606.01540.
- Bubeck, S. and Cesa-Bianchi, N. (2012). Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122.
- Cohen, A., Yu, L., and Wright, R. (2018). Diverse exploration for fast and safe policy improvement. In *AAAI*, pages 2876–2883. AAAI Press.
- Dann, C. and Brunskill, E. (2015). Sample complexity of episodic fixed-horizon reinforcement learning. In *NeurIPS*, pages 2818–2826.
- Dann, C., Lattimore, T., and Brunskill, E. (2017). Unifying PAC and regret: Uniform PAC bounds for episodic reinforcement learning. In *NeurIPS*, pages 5713–5723.
- Deisenroth, M. P., Neumann, G., and Peters, J. (2013). A survey on policy search for robotics. *Foundations and Trends in Robotics*, 2(1-2):1–142.
- Duan, Y., Chen, X., Houthoofd, R., Schulman, J., and Abbeel, P. (2016). Benchmarking deep reinforcement learning for continuous control. In *ICML*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 1329–1338. JMLR.org.
- García, J. and Fernández, F. (2015). A comprehensive survey on safe reinforcement learning. *J. Mach. Learn. Res.*, 16:1437–1480.
- Haarnoja, T., Tang, H., Abbeel, P., and Levine, S. (2017). Reinforcement learning with deep energy-based policies. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 1352–1361. PMLR.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. (2018). Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *ICML*, volume 80 of *Proceedings of Machine Learning Research*, pages 1856–1865. PMLR.
- Hans, A., Schneegaß, D., Schäfer, A. M., and Udfluft, S. (2008). Safe exploration for reinforcement learning. In *ESANN*, pages 143–148.
- Härdle, W. and Simar, L. (2012). *Applied multivariate statistical analysis*, volume 22007. Springer.
- Hazan, E., Levy, K. Y., and Shalev-Shwartz, S. (2016). On graduated optimization for stochastic non-convex problems. In *ICML*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 1833–1841. JMLR.org.
- Jaksch, T., Ortner, R., and Auer, P. (2010). Near-optimal regret bounds for reinforcement learning. *J. Mach. Learn. Res.*, 11:1563–1600.
- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. (2018). Is q-learning provably efficient? In *NeurIPS*, pages 4868–4878.
- Kakade, S. M. (2001). A natural policy gradient. In *NeurIPS*, pages 1531–1538. MIT Press.
- Kakade, S. M. and Langford, J. (2002). Approximately optimal approximate reinforcement learning. In *ICML*, pages 267–274. Morgan Kaufmann.
- Kearns, M. J. and Singh, S. P. (2002). Near-optimal reinforcement learning in polynomial time. *Machine Learning*, 49(2-3):209–232.
- Kingma, D. P. and Ba, J. (2015). Adam: A method for stochastic optimization. In *ICLR (Poster)*.

- Lakshmanan, K., Ortner, R., and Ryabko, D. (2015). Improved regret bounds for undiscounted continuous reinforcement learning. In *ICML*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 524–532. JMLR.org.
- Lanczos, C. (1950). *An iteration method for the solution of the eigenvalue problem of linear differential and integral operators*. United States Governm. Press Office Los Angeles, CA.
- Lattimore, T. and Hutter, M. (2014). Near-optimal PAC bounds for discounted mdps. *Theor. Comput. Sci.*, 558:125–143.
- Lattimore, T. and Szepesvári, C. (2019). *Bandit Algorithms*. Cambridge University Press (preprint).
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. (2016). Continuous control with deep reinforcement learning. In *ICLR (Poster)*.
- Miyamae, A., Nagata, Y., Ono, I., and Kobayashi, S. (2010). Natural policy gradient methods with parameter-based exploration for control tasks. In *NeurIPS*, pages 1660–1668. Curran Associates, Inc.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M. A., Fidjeland, A., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533.
- Nachum, O., Norouzi, M., Tucker, G., and Schuurmans, D. (2018). Smoothed action value functions for learning gaussian policies. In *ICML*, volume 80 of *Proceedings of Machine Learning Research*, pages 3689–3697. PMLR.
- Ok, J., Proutière, A., and Tranos, D. (2018). Exploration in structured reinforcement learning. In *NeurIPS*, pages 8888–8896.
- OpenAI (2018). Openai five. <https://blog.openai.com/openai-five/>.
- Ortner, R. and Ryabko, D. (2012). Online regret bounds for undiscounted continuous reinforcement learning. In *NeurIPS*, pages 1772–1780.
- Papini, M., Pirotta, M., and Restelli, M. (2017). Adaptive batch size for safe policy gradients. In *NeurIPS*, pages 3591–3600.
- Papini, M., Pirotta, M., and Restelli, M. (2019). Smoothing policies and safe policy gradients. *CoRR*, abs/1905.03231.
- Peters, J. and Schaal, S. (2008). Reinforcement learning of motor skills with policy gradients. *Neural Networks*, 21(4):682–697.
- Peters, J., Vijayakumar, S., and Schaal, S. (2005). Natural actor-critic. In *ECML*, volume 3720 of *Lecture Notes in Computer Science*, pages 280–291. Springer.
- Pirotta, M., Restelli, M., and Bascetta, L. (2013). Adaptive step-size for policy gradient methods. In *NeurIPS*, pages 1394–1402.
- Pirotta, M., Restelli, M., and Bascetta, L. (2015). Policy gradient in lipschitz markov decision processes. *Machine Learning*, 100(2-3):255–283.
- Puterman, M. L. (1994). *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley Series in Probability and Statistics. Wiley.
- Rajeswaran, A., Lowrey, K., Todorov, E., and Kakade, S. M. (2017). Towards generalization and simplicity in continuous control. In *NeurIPS*, pages 6550–6561.
- Schraudolph, N. N. (1999). Local gain adaptation in stochastic gradient descent.
- Schulman, J., Abbeel, P., and Chen, X. (2017). Equivalence between policy gradients and soft q-learning. *CoRR*, abs/1704.06440.
- Sehnke, F., Osendorfer, C., Rückstieß, T., Graves, A., Peters, J., and Schmidhuber, J. (2008). Policy gradients with parameter-based exploration for control. In *ICANN (1)*, volume 5163 of *Lecture Notes in Computer Science*, pages 387–396. Springer.
- Shani, L., Efroni, Y., and Mannor, S. (2018). Revisiting exploration-conscious reinforcement learning. *CoRR*, abs/1812.05551.
- Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L., Kumaran, D., Graepel, T., Lillicrap, T. P., Simonyan, K., and Hassabis, D. (2017). Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *CoRR*, abs/1712.01815.
- Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., and Riedmiller, M. A. (2014). Deterministic policy gradient algorithms. In *ICML*, volume 32 of *JMLR Workshop and Conference Proceedings*, pages 387–395. JMLR.org.
- Strehl, A. L., Li, L., and Littman, M. L. (2009). Reinforcement learning in finite mdps: PAC analysis. *J. Mach. Learn. Res.*, 10:2413–2444.
- Sutton, R. S. (1992). Adapting bias by gradient descent: An incremental version of delta-bar-delta. In *AAAI*, pages 171–176. AAAI Press / The MIT Press.
- Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.
- Sutton, R. S., McAllester, D. A., Singh, S. P., and Mansour, Y. (1999). Policy gradient methods for

- reinforcement learning with function approximation. In *NeurIPS*, pages 1057–1063. The MIT Press.
- Thomas, P. S., Theodorou, G., and Ghavamzadeh, M. (2015). High confidence policy improvement. In *ICML*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 2380–2388. JMLR.org.
- Tucker, G., Bhupatiraju, S., Gu, S., Turner, R. E., Ghahramani, Z., and Levine, S. (2018). The mirage of action-dependent baselines in reinforcement learning. In *ICML*, volume 80 of *Proceedings of Machine Learning Research*, pages 5022–5031. PMLR.
- Turchetta, M., Berkenkamp, F., and Krause, A. (2016). Safe exploration in finite markov decision processes with gaussian processes. In *NeurIPS*, pages 4305–4313.
- Veeriah, V., Zhang, S., and Sutton, R. S. (2017). Crossprop: Learning representations by stochastic meta-gradient descent in neural networks. In *ECML/PKDD (1)*, volume 10534 of *Lecture Notes in Computer Science*, pages 445–459. Springer.
- Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8:229–256.
- Xu, Z., van Hasselt, H., and Silver, D. (2018). Meta-gradient reinforcement learning. In *NeurIPS*, pages 2402–2413.
- Zhao, T., Hachiya, H., Niu, G., and Sugiyama, M. (2011). Analysis and improvement of policy gradient estimation. In *NeurIPS*, pages 262–270.

Table of Supplementary Contents

- Appendix A: Proofs
- Appendix B: Meta Gradient Estimation
- Appendix C: Characterizing the Estimation Error
- Appendix D: Extensions
- Appendix E: Task Specifications
- Appendix F: Experimental Setting

A Proofs

In this section we provide proofs for all the formal statements made in the paper. Recall that R_{\max} is the maximum absolute value reward, φ the Euclidean-norm bound on state features, γ the discount factor, and π with no subscript denotes the mathematical constant.

A.1 Exact framework

Lemma 2.1. *Let Π_Υ be the class of Gaussian policies parametrized as in (4), but with fixed variance parameter ω . Let $\mathbf{v}_t \in \mathbb{R}^m$ and $\mathbf{v}_{t+1} = \mathbf{v}_t + \alpha_t \nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega)$. For any $C_t \leq C_t^*$, the largest step size guaranteeing $J(\mathbf{v}_{t+1}, \omega) - J(\mathbf{v}_t, \omega) \geq C_t$ is:*

$$\bar{\alpha}_t := \frac{\sigma_\omega^2}{F} \left(1 + \sqrt{1 - \frac{C_t}{C_t^*}} \right), \quad (10)$$

where $F = \frac{2\varphi^2 R_{\max}}{(1-\gamma)^2} \left(1 + \frac{2\gamma}{\pi(1-\gamma)} \right)$ and $C_t^* = \frac{\sigma_\omega^2 \|\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega)\|^2}{2F}$. \square

Proof. This is just a slight adaptation of existing results from (Papini et al., 2019). Since the fixed-variance Gaussian policy is smoothing (Lemma 15 from Papini et al., 2019), we plug its smoothing constants (8) into (9) (Theorem 9 from Papini et al., 2019) to obtain, for all $\alpha \in \mathbb{R}$:

$$J(\mathbf{v}_{t+1}, \omega) - J(\mathbf{v}_t, \omega) \geq \alpha \|\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega)\|^2 - \alpha^2 \frac{F}{2\sigma_\omega^2} \|\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega)\|^2 := f(\alpha). \quad (19)$$

Thus, imposing $f(\alpha) \geq C_t$ is enough to ensure $J(\mathbf{v}_{t+1}, \omega) - J(\mathbf{v}_t, \omega) \geq C_t$. This yields a second-order inequality in α , whose solution is:

$$\frac{\sigma_\omega^2}{F} \left(1 - \sqrt{1 - \frac{C_t}{C_t^*}} \right) \leq \alpha \leq \frac{\sigma_\omega^2}{F} \left(1 + \sqrt{1 - \frac{C_t}{C_t^*}} \right), \quad (20)$$

provided $\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega) \neq 0$, which is a reasonable assumption (if $\nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega) = 0$, the algorithm has already converged and any update is void), and $C_t \leq C_t^*$, which is true by hypothesis. To conclude the proof, we just select the largest step size satisfying (20). \square

Lemma 4.1. *Let Π_Ω be the class of Gaussian policies parametrized as in (4), but with fixed mean parameter \mathbf{v} . Π_Ω is $\left(\frac{4}{\sqrt{2\pi e}}, 2, 2 \right)$ -smoothing.* \square

Proof. The definition of smoothing policies is from (Definition 4 in Papini et al., 2019) and reported in (7). Since the mean parameter \mathbf{v} is fixed, the policy parameter space can be restricted to Ω . Recall that $\sigma_\omega = e^\omega$. We need the following derivatives:

$$\begin{aligned} \nabla_\omega \log \pi_{\mathbf{v}, \omega}(a|s) &= \nabla_\omega \left(-\omega - \frac{1}{2} e^{-2\omega} (a - \mu_{\mathbf{v}}(s))^2 \right) \\ &= e^{-2\omega} (a - \mu_{\mathbf{v}}(s))^2 - 1, \end{aligned} \quad (21)$$

$$\nabla_\omega \nabla_\omega^T \log \pi_{\mathbf{v}, \omega}(a|s) = -2e^{-2\omega} (a - \mu_{\mathbf{v}}(s))^2. \quad (22)$$

Let $x := e^{-\omega}(a - \mu_{\mathbf{v}}(s))$ in the following, and note that $\nabla_{\omega}x = e^{-\omega}$. First we compute ψ :

$$\begin{aligned}
 & \sup_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_{\mathbf{v}, \omega}} [\|\nabla_{\omega} \log \pi_{\mathbf{v}, \omega}(a|s)\|] \\
 &= \sup_{s \in \mathcal{S}} \frac{e^{-\omega}}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-\frac{1}{2}e^{-2\omega}(a - \mu_{\mathbf{v}}(s))^2} |e^{-2\omega}(a - \mu_{\mathbf{v}}(s))^2 - 1| da \\
 &= \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-x^2/2} |x^2 - 1| dx \\
 &= \frac{4}{\sqrt{2\pi}e} := \psi.
 \end{aligned} \tag{23}$$

Next, we compute κ :

$$\begin{aligned}
 & \sup_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_{\mathbf{v}, \omega}} [\|\nabla_{\omega} \log \pi_{\mathbf{v}, \omega}(a|s)\|^2] \\
 &= \sup_{s \in \mathcal{S}} \frac{e^{-\omega}}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-\frac{1}{2}e^{-2\omega}(a - \mu_{\mathbf{v}}(s))^2} (e^{-2\omega}(a - \mu_{\mathbf{v}}(s))^2 - 1)^2 da \\
 &= \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-x^2/2} (x^2 - 1)^2 dx = 2 := \kappa.
 \end{aligned} \tag{24}$$

Finally, we compute ξ :

$$\begin{aligned}
 & \sup_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_{\mathbf{v}, \omega}} [\|\nabla_{\omega} \nabla_{\omega}^T \log \pi_{\mathbf{v}, \omega}(a|s)\|] \\
 &= \sup_{s \in \mathcal{S}} \frac{e^{-\omega}}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-\frac{1}{2}e^{-2\omega}(a - \mu_{\mathbf{v}}(s))^2} |-2e^{-2\omega}(a - \mu_{\mathbf{v}}(s))^2| da \\
 &= \frac{2}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-x^2/2} x^2 dx = 2 := \xi.
 \end{aligned} \tag{25}$$

Indeed, the computed constants are independent from the value of ω . \square

Theorem 4.2. *Let Π_{Ω} be the class of policies defined in Lemma 4.1. Let $\omega_t \in \Omega$ and $\omega_{t+1} \leftarrow \omega_t + \beta_t \nabla_{\omega} J(\mathbf{v}, \omega_t)$. For any $C_t \leq C_t^*$, the largest step-size satisfying (6) is:*

$$\bar{\beta}_t = \frac{1}{G} \left(1 + \sqrt{1 - \frac{C_t}{C_t^*}} \right), \tag{15}$$

where $C_t^* = \frac{\|\nabla_{\omega} J(\mathbf{v}, \omega_t)\|^2}{2G}$ and $G = \frac{4R_{\max}}{(1-\gamma)^2} \left(1 + \frac{4\gamma}{\pi e(1-\gamma)} \right)$. \square

Proof. Similarly to the proof of Lemma 2.1, we plug the smoothing constants from Lemma 4.1 into (9) to obtain:

$$J(\mathbf{v}, \omega_{t+1}) - J(\mathbf{v}, \omega_t) \geq \beta \|\nabla_{\omega} J(\mathbf{v}, \omega_t)\|^2 - \beta^2 \frac{G}{2} \|\nabla_{\omega} J(\mathbf{v}, \omega_t)\|^2 := f(\beta). \tag{26}$$

Solving $f(\beta) \geq C_t$ yields:

$$\frac{1}{G} \left(1 - \sqrt{1 - \frac{C_t}{C_t^*}} \right) \leq \beta \leq \frac{1}{G} \left(1 + \sqrt{1 - \frac{C_t}{C_t^*}} \right), \tag{27}$$

where, again, we assume the policy gradient is non-zero. The proof is concluded by selecting the largest step size satisfying (27). \square

Theorem 4.3. *Let Π_{Θ} be a (ψ, κ, ξ) -smoothing policy class, $\boldsymbol{\theta}_t \in \Theta$, and $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \eta_t \mathbf{x}_t$, where $\mathbf{x}_t \in \mathbb{R}^m$ and $\eta_t \in \mathbb{R}$ is a (possibly negative) step size. Let $\lambda_t := \frac{\langle \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_t), \mathbf{x}_t \rangle}{\|\mathbf{x}_t\|}$ be the scalar projection of $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_t)$ onto \mathbf{x}_t . For any $C_t \leq C_t^*$, provided $\lambda_t \neq 0$, the largest step size guaranteeing $J(\boldsymbol{\theta}_{t+1}) - J(\boldsymbol{\theta}_t) \geq C_t$ is:*

$$\bar{\eta}_t = \frac{|\lambda_t|}{L \|\mathbf{x}_t\|} \left(\text{sign}(\lambda_t) + \sqrt{1 - \frac{C_t}{C_t^*}} \right), \tag{16}$$

where $C_t^* = \frac{\lambda_t^2}{2L}$ and $L = \frac{R_{\max}}{(1-\gamma)^2} \left(\frac{2\gamma\psi^2}{1-\gamma} + \kappa + \xi \right)$. \square

Proof. Since we are no longer dealing with a policy gradient update, we need a generalization of (9). By Theorem 8 from (Papini et al., 2019):

$$J(\boldsymbol{\theta}_t + \Delta\boldsymbol{\theta}) - J(\boldsymbol{\theta}_t) \geq \langle \Delta\boldsymbol{\theta}, \nabla J(\boldsymbol{\theta}_t) \rangle - \frac{L}{2} \|\Delta\boldsymbol{\theta}\|^2, \quad (28)$$

for any parameter update $\Delta\boldsymbol{\theta} \in \mathbb{R}^m$. In our case:

$$J(\boldsymbol{\theta}_{t+1}) - J(\boldsymbol{\theta}_t) \geq \eta_t \langle \mathbf{x}_t, \nabla J(\boldsymbol{\theta}_t) \rangle - \eta_t^2 \frac{L}{2} \|\mathbf{x}_t\|^2 := f(\eta_t), \quad (29)$$

where $\langle \cdot, \cdot \rangle$ denotes the dot product. Intuitively, the more \mathbf{x}_t agrees with the improvement direction $\nabla J(\boldsymbol{\theta}_t)$, the more improvement can be guaranteed. We first assume $\langle \mathbf{x}_t, \nabla J(\boldsymbol{\theta}_t) \rangle \neq 0$. Solving $f(\eta_t) \geq C_t$ yields:

$$\frac{|\lambda_t|}{L \|\mathbf{x}_t\|} \left(\text{sign}(\lambda_t) - \sqrt{1 - \frac{C_t}{C_t^*}} \right) \leq \eta_t \leq \frac{|\lambda_t|}{L \|\mathbf{x}_t\|} \left(\text{sign}(\lambda_t) + \sqrt{1 - \frac{C_t}{C_t^*}} \right), \quad (30)$$

from which we select the largest safe step size. For $C_t \geq 0$, depending on the sign of λ_t (i.e., whether \mathbf{x}_t agrees with the gradient or not) and on the value of C_t , the step size may be non-positive. Intuitively, if a positive improvement is required but \mathbf{x}_t is pejorative, a negative step size is used to invert it. Instead, if $C_t < 0$ (bounded worsening), a small-enough step in the direction of \mathbf{x}_t is always acceptable.

We now consider the special case $\langle \mathbf{x}_t, \nabla J(\boldsymbol{\theta}_t) \rangle = 0$, i.e., $\lambda = 0$. In this case, only non-positive values of C_t are allowed. Intuitively, \mathbf{x}_t is orthogonal to the improvement direction, so no positive improvement can be guaranteed. Under the restriction $C_t \leq 0$, the following range of step sizes is safe:

$$-\frac{1}{\|\mathbf{x}_t\|} \sqrt{-\frac{2C_t}{L}} \leq \eta_t \leq \frac{1}{\|\mathbf{x}_t\|} \sqrt{-\frac{2C_t}{L}}, \quad (31)$$

from which we select $\bar{\eta}_t = \frac{1}{\|\mathbf{x}_t\|} \sqrt{-\frac{2C_t}{L}}$. □

A.2 Approximate Framework

In the following, let $\hat{\nabla}_{\mathbf{v}}^N J$, $\hat{\nabla}_{\omega}^N J$ and $\hat{\nabla}_{\omega}^N \mathcal{L}$ be unbiased estimators of $\nabla_{\mathbf{v}} J$, $\nabla_{\omega} J$ and $\nabla_{\omega} \mathcal{L}$, respectively, each using a batch of N trajectories.

Corollary A.1. *Let Π_{Υ} be the class of Gaussian policies parametrized as in (4), but with fixed variance parameter ω . Let $\mathbf{v}_t \in \mathbb{R}^m$ and $\mathbf{v}_{t+1} = \mathbf{v}_t + \alpha_t \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega)$. Under Assumption 4.4, provided $N > \epsilon_{\delta}^2 / \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega_t) \right\|^2$, for any $C_t \leq C_t^*$, the largest step size satisfying (6) with probability at least $1 - \delta$ is:*

$$\tilde{\alpha}_t = \frac{\sigma_{\omega}^2 \left(\left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega_t) \right\| - \frac{\epsilon_{\delta}}{\sqrt{N}} \right)}{F \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega_t) \right\|} \left(1 + \sqrt{1 - \frac{C_t}{C_t^*}} \right),$$

where $C_t^* = \frac{\sigma_{\omega}^2 \left(\left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega_t) \right\| - \frac{\epsilon_{\delta}}{\sqrt{N}} \right)^2}{2F}$ and F is from Lemma 2.1. □

Proof. From (28) with $\Theta = \Upsilon$ and $\Delta\theta = \alpha_t \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega)$ we have:

$$\begin{aligned}
 J(\mathbf{v}_{t+1}, \omega) - J(\mathbf{v}_t, \omega) &\geq \alpha_t \left\langle \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega), \nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega) \right\rangle - \alpha_t^2 \frac{L}{2} \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\|^2 \\
 &= \alpha_t \left\langle \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega), \nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega) \pm \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\rangle \\
 &\quad - \alpha_t^2 \frac{L}{2} \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\|^2 \\
 &= \alpha_t \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\|^2 + \alpha_t \left\langle \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega), \nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega) - \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\rangle \\
 &\quad - \alpha_t^2 \frac{L}{2} \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\|^2 \\
 &\geq \alpha_t \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\|^2 - \alpha_t \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\| \left\| \nabla_{\mathbf{v}} J(\mathbf{v}_t, \omega) - \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\| \\
 &\quad - \alpha_t^2 \frac{L}{2} \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\|^2 \tag{32}
 \end{aligned}$$

$$\begin{aligned}
 &\geq \alpha_t \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\|^2 - \alpha_t \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\| \frac{\epsilon_\delta}{\sqrt{N}} \\
 &\quad - \alpha_t^2 \frac{L}{2} \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\|^2 \tag{33}
 \end{aligned}$$

$$\begin{aligned}
 &\geq \alpha_t \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\| \left(\left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\| - \frac{\epsilon_\delta}{\sqrt{N}} \right) \\
 &\quad - \alpha_t^2 \frac{L}{2} \left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\|^2, \tag{34}
 \end{aligned}$$

where (32) is from Cauchy-Schwartz inequality and (33) is from Assumption 4.4. The hypothesis on the batch size makes the $\left(\left\| \hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega) \right\| - \frac{\epsilon_\delta}{\sqrt{N}} \right)$ term positive. We then proceed as in the proof of Lemma 2.1. \square

Corollary A.2. Let Π_Ω be the class of policies defined in Lemma 4.1. Let $\omega_t \in \Omega$ and $\omega_{t+1} \leftarrow \omega_t + \eta_t \hat{\nabla}_{\omega}^N \mathcal{L}(\mathbf{v}, \omega_t)$. Under Assumption 4.4, provided $\hat{\lambda}_t \neq 0$ and $N > \epsilon_\delta^2 / \hat{\lambda}_t^2$, for any $C_t \leq C_t^*$, the largest step-size satisfying (6) with probability at least $1 - \delta$ is:

$$\tilde{\eta}_t = \frac{\left| \hat{\lambda}_t \right| - \frac{\epsilon_\delta}{\sqrt{N}}}{G \left\| \hat{\nabla}_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t) \right\|} \left(\text{sign}(\hat{\lambda}_t) + \sqrt{1 - \frac{C_t}{C_t^*}} \right),$$

where $\hat{\lambda}_t := \frac{\langle \hat{\nabla}_{\omega} J(\mathbf{v}_t, \omega_t), \hat{\nabla}_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t) \rangle}{\left\| \hat{\nabla}_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t) \right\|}$ is the scalar projection of $\hat{\nabla}_{\omega} J(\mathbf{v}_t, \omega_t)$ onto $\hat{\nabla}_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t)$, $C_t^* = \frac{\left(\left| \hat{\lambda}_t \right| - \frac{\epsilon_\delta}{\sqrt{N}} \right)^2}{2G}$, and G is from Theorem 4.2. \square

Proof. From (28) with $\Theta = \Omega$ and $\Delta\theta = \eta_t \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t)$ we have:

$$\begin{aligned}
 J(\mathbf{v}, \omega_{t+1}) - J(\mathbf{v}, \omega_t) &\geq \eta_t \left\langle \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t), \nabla_\omega J(\mathbf{v}, \omega_t) \right\rangle - \eta_t^2 \frac{L}{2} \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t) \right\|^2 \\
 &= \eta_t \left\langle \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t), \nabla_\omega J(\mathbf{v}, \omega_t) \pm \hat{\nabla}_\omega^N J(\mathbf{v}, \omega_t) \right\rangle - \eta_t^2 \frac{L}{2} \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t) \right\|^2 \\
 &= \eta_t \left\langle \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t), \hat{\nabla}_\omega^N J(\mathbf{v}, \omega_t) \right\rangle \\
 &\quad + \eta_t \left\langle \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t), \nabla_\omega J(\mathbf{v}, \omega_t) - \hat{\nabla}_\omega^N J(\mathbf{v}, \omega_t) \right\rangle \\
 &\quad - \eta_t^2 \frac{L}{2} \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t) \right\|^2 \\
 &\geq \eta_t \left\langle \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t), \hat{\nabla}_\omega^N J(\mathbf{v}, \omega_t) \right\rangle \\
 &\quad - |\eta_t| \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t) \right\| \left\| \nabla_\omega J(\mathbf{v}, \omega_t) - \hat{\nabla}_\omega^N J(\mathbf{v}, \omega_t) \right\| \\
 &\quad - \eta_t^2 \frac{L}{2} \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t) \right\|^2 \tag{35}
 \end{aligned}$$

$$\begin{aligned}
 &\geq \eta_t \left\langle \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t), \hat{\nabla}_\omega^N J(\mathbf{v}, \omega_t) \right\rangle - |\eta_t| \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t) \right\| \frac{\epsilon_\delta}{\sqrt{N}} \\
 &\quad - \eta_t^2 \frac{L}{2} \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t) \right\|^2 \tag{36}
 \end{aligned}$$

$$\begin{aligned}
 &\geq \eta_t \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t) \right\| \left(\hat{\lambda}_t - \text{sign}(\eta_t) \frac{\epsilon_\delta}{\sqrt{N}} \right) \\
 &\quad - \eta_t^2 \frac{L}{2} \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t) \right\|^2, \tag{37}
 \end{aligned}$$

where (35) is from Cauchy-Schwartz inequality and (36) is from Assumption 4.4. Note the absolute value on η_t , which may be negative. We first consider the case $\eta_t > 0$, which yields:

$$\begin{aligned}
 J(\mathbf{v}, \omega_{t+1}) - J(\mathbf{v}, \omega_t) &\geq \eta_t \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t) \right\| \left(\hat{\lambda}_t - \frac{\epsilon_\delta}{\sqrt{N}} \right) \\
 &\quad - \eta_t^2 \frac{L}{2} \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t) \right\|^2. \tag{38}
 \end{aligned}$$

Solving the safety constraint for $\tilde{\eta}_t$ yields:

$$\tilde{\eta}_t = \frac{1}{G \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}_t, \omega_t) \right\|} \left(\hat{\lambda}_t - \frac{\epsilon_\delta}{\sqrt{N}} + \left| \hat{\lambda}_t - \frac{\epsilon_\delta}{\sqrt{N}} \right| \sqrt{1 - \frac{C_t}{C_t^*}} \right). \tag{39}$$

Given the batch size condition, the step size is indeed positive if and only if $\hat{\lambda}_t > 0$. We then consider the case $\eta_t \leq 0$, which yields:

$$\begin{aligned}
 J(\mathbf{v}, \omega_{t+1}) - J(\mathbf{v}, \omega_t) &\geq \eta_t \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t) \right\| \left(\hat{\lambda}_t + \frac{\epsilon_\delta}{\sqrt{N}} \right) \\
 &\quad - \eta_t^2 \frac{L}{2} \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}, \omega_t) \right\|^2. \tag{40}
 \end{aligned}$$

Solving the safety constraint for $\tilde{\eta}_t$ yields:

$$\tilde{\eta}_t = \frac{1}{G \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}_t, \omega_t) \right\|} \left(\hat{\lambda}_t + \frac{\epsilon_\delta}{\sqrt{N}} + \left| \hat{\lambda}_t + \frac{\epsilon_\delta}{\sqrt{N}} \right| \sqrt{1 - \frac{C_t}{C_t^*}} \right). \tag{41}$$

Given the batch size condition, the step size is indeed non-positive if and only if $\hat{\lambda}_t < 0$. The two cases can be unified as:

$$\bar{\eta}_t = \begin{cases} \frac{\hat{\lambda}_t - \frac{\epsilon_\delta}{\sqrt{N}}}{G \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}_t, \omega_t) \right\|} \left(1 + \sqrt{1 - \frac{C_t}{C_t^*}} \right) & \text{if } \hat{\lambda}_t > 0, \\ \frac{\hat{\lambda}_t + \frac{\epsilon_\delta}{\sqrt{N}}}{G \left\| \hat{\nabla}_\omega^N \mathcal{L}(\mathbf{v}_t, \omega_t) \right\|} \left(1 - \sqrt{1 - \frac{C_t}{C_t^*}} \right) & \text{if } \hat{\lambda}_t < 0, \end{cases} \tag{42}$$

which can be further simplified to obtain the $\tilde{\eta}_t$ in the thesis.

As in the exact framework, we can treat the case $\hat{\lambda}_t = 0$ separately. Under the restriction $C_t \leq 0$, the following range of step sizes is safe:

$$-\frac{1}{\|\hat{\nabla}_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t)\|} \sqrt{-\frac{2C_t}{G}} \leq \eta_t \leq \frac{1}{\|\hat{\nabla}_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t)\|} \sqrt{-\frac{2C_t}{G}}, \quad (43)$$

from which we select $\bar{\eta}_t = \frac{1}{\|\hat{\nabla}_{\omega} \mathcal{L}(\mathbf{v}_t, \omega_t)\|} \sqrt{-\frac{2C_t}{G}}$. No assumption on the batch size is requested, but only non-positive improvement constraints can be satisfied. \square

Theorem 4.5. *Under Assumption 4.4, provided $C_t \leq 0$ and $N > \epsilon_{\delta}^2 / \hat{\lambda}_t^2$, Algorithm 2 guarantees $J(\boldsymbol{\theta}_{t+1}) - J(\boldsymbol{\theta}_t) \geq C_t$ with probability at least $1 - \delta$ for any $t \geq 1$.* \square

Proof. We just combine the results from Corollary A.1 and A.2. We actually obtain a stronger versions than the one reported in the main paper: for the odd steps of Algorithm 2 (mean updates), a batch size $N > \epsilon_{\delta}^2 / \|\hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega_t)\|^2$ is enough to satisfy (6) for any $C_t \leq \frac{\sigma_{\omega_t}^2 (\|\hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega_t)\| - \frac{\epsilon_{\delta}}{\sqrt{N}})^2}{2F}$. For the even steps (variance updates), a batch size $N > \epsilon_{\delta}^2 / \hat{\lambda}_t^2$ allows to satisfy (6) for any $C_t \leq \frac{(|\hat{\lambda}_t| - \frac{\epsilon_{\delta}}{\sqrt{N}})^2}{2G}$. The looser version from the paper is then obtained by observing that $\|\hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega_t)\| > |\hat{\lambda}_t|$ always, since $\hat{\lambda}_t$ is a scalar projection of $\hat{\nabla}_{\mathbf{v}}^N J(\mathbf{v}_t, \omega_t)$, and that both upper bounds on C_t are non-negative. \square

B Meta Gradient Estimation

In this section, we propose an unbiased estimator for $\nabla_{\omega} \|\nabla_{\mathbf{v}} J(\boldsymbol{\theta})\|$, which is necessary to estimate $\nabla_{\omega} \mathcal{L}(\boldsymbol{\theta})$ (recall that $\boldsymbol{\theta} = [\mathbf{v}, \omega]$ is the full vector of policy parameters). First note that:

$$\nabla_{\omega} \|\nabla_{\mathbf{v}} J(\boldsymbol{\theta})\| = \frac{\langle \nabla_{\mathbf{v}} J(\boldsymbol{\theta}), \nabla_{\omega} \nabla_{\mathbf{v}} J(\boldsymbol{\theta}) \rangle}{\|\nabla_{\mathbf{v}} J(\boldsymbol{\theta})\|}, \quad (44)$$

which is the scalar projection of $\nabla_{\omega} \nabla_{\mathbf{v}} J$ onto $\nabla_{\mathbf{v}} J$.

An estimator for $\nabla_{\mathbf{v}} J(\boldsymbol{\theta})$ is already available (2). We now show how to estimate $\nabla_{\omega} \nabla_{\mathbf{v}} J(\boldsymbol{\theta})$. First note that:

$$\begin{aligned} \nabla_{\omega} \nabla_{\mathbf{v}} \log p_{\boldsymbol{\theta}}(\tau) &= \sum_{h=0}^H \nabla_{\omega} \nabla_{\mathbf{v}} \log \pi_{\boldsymbol{\theta}}(a_h | s_h) = \sum_{h=0}^H \nabla_{\omega} \frac{a_h - \mu_{\boldsymbol{\theta}}(s_h)}{e^{2\omega}} = -2 \sum_{t=0}^H \frac{a - \mu_{\boldsymbol{\theta}}}{e^{2\omega}} \\ &= -2 \nabla_{\mathbf{v}} \log p_{\boldsymbol{\theta}}(\tau). \end{aligned} \quad (45)$$

Using the log-trick:

$$\begin{aligned} \nabla_{\omega} \nabla_{\mathbf{v}} J(\boldsymbol{\theta}) &= \nabla_{\omega} \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}} [\nabla_{\mathbf{v}} \log p_{\boldsymbol{\theta}}(\tau) R(\tau)] \\ &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}} [\nabla_{\omega} \log p_{\boldsymbol{\theta}}(\tau) \nabla_{\mathbf{v}} \log p_{\boldsymbol{\theta}}(\tau) R(\tau)] + \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}} [\nabla_{\omega} \nabla_{\mathbf{v}} \log p_{\boldsymbol{\theta}}(\tau) R(\tau)] \\ &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}} [\nabla_{\omega} \log p_{\boldsymbol{\theta}}(\tau) \nabla_{\mathbf{v}} \log p_{\boldsymbol{\theta}}(\tau) R(\tau)] - 2 \nabla_{\mathbf{v}} J(\boldsymbol{\theta}) \\ &:= \text{mix}(\boldsymbol{\theta}) - 2 \nabla_{\mathbf{v}} J(\boldsymbol{\theta}). \end{aligned} \quad (46)$$

Since we can reuse the policy gradient w.r.t. \mathbf{v} in (46), we have reduced the problem of estimating $\nabla_{\omega} \mathcal{L}(\boldsymbol{\theta})$ to that of estimating $\text{mix}(\boldsymbol{\theta}) := \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}} [\nabla_{\omega} \log p_{\boldsymbol{\theta}}(\tau) \nabla_{\mathbf{v}} \log p_{\boldsymbol{\theta}}(\tau) R(\tau)]$. The following estimator is inspired by GPOMDP (Baxter and Bartlett, 2001):

Theorem B.1. *An unbiased estimator for $\text{mix}(\boldsymbol{\theta}) := \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}} [\nabla_{\omega} \log p_{\boldsymbol{\theta}}(\tau) \nabla_{\mathbf{v}} \log p_{\boldsymbol{\theta}}(\tau) R(\tau)]$ is:*

$$\widehat{\text{mix}}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^N \sum_{h=0}^H \gamma^h r_h^n \left(\sum_{i=0}^h \nabla_{\omega} \log \pi_{\boldsymbol{\theta}}(a_i^n | s_i^n) \right) \left(\sum_{j=0}^h \nabla_{\mathbf{v}} \log \pi_{\boldsymbol{\theta}}(a_j^n | s_j^n) \right), \quad (47)$$

where subscripts denote time steps and superscripts denote trajectories. To preserve the unbiasedness of the estimator, separate trajectories must be used to compute the two inner sums. \square

Proof. Let's abbreviate action probabilities as $\pi_k = \pi_{\theta}(a_k|s_k)$ and sub-trajectory probabilities as $p_{\theta}(\tau_{h:k}) = \pi_{\theta}(a_h|s_h)\mathcal{P}(s_{h+1}|s_h, a_h) \dots \mathcal{P}(s_k|s_{k-1}, a_{k-1})$. We can split $mix(\theta)$ into the sum of four components:

$$\begin{aligned} mix(\theta) &= \int_T p_{\theta}(\tau_{0:H}) \nabla_{\omega} \log p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) R(\tau) d\tau = \\ &= \int_T p_{\theta}(\tau_{0:H}) \left(\sum_{h=0}^H \gamma^h r_h \right) \left(\sum_{i=0}^H \nabla_{\omega} \log \pi_i \right) \left(\sum_{j=0}^H \nabla_{\mathbf{v}} \log \pi_j \right) d\tau = \\ &= \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:H}) \gamma^h r_h \left(\sum_{i=0}^h \nabla_{\omega} \log \pi_i \right) \left(\sum_{j=0}^h \nabla_{\mathbf{v}} \log \pi_j \right) d\tau \end{aligned} \quad (48)$$

$$+ \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:H}) \gamma^h r_h \left(\sum_{i=0}^h \nabla_{\omega} \log \pi_i \right) \left(\sum_{j=h+1}^H \nabla_{\mathbf{v}} \log \pi_j \right) d\tau \quad (49)$$

$$+ \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:H}) \gamma^h r_h \left(\sum_{i=h+1}^H \nabla_{\omega} \log \pi_i \right) \left(\sum_{j=0}^h \nabla_{\mathbf{v}} \log \pi_j \right) d\tau \quad (50)$$

$$+ \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:H}) \gamma^h r_h \left(\sum_{i=h+1}^H \nabla_{\omega} \log \pi_i \right) \left(\sum_{j=h+1}^H \nabla_{\mathbf{v}} \log \pi_j \right) d\tau. \quad (51)$$

Next, we show that (49), (50) and (51) all evaluate to 0:

$$\begin{aligned} (49) &= \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:H}) \gamma^h r_h \left(\sum_{i=0}^h \nabla_{\omega} \log \pi_i \right) \left(\sum_{j=h+1}^H \nabla_{\mathbf{v}} \log \pi_j \right) d\tau \\ &= \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:h}) \gamma^h r_h \left(\sum_{i=0}^h \nabla_{\omega} \log \pi_i \right) d\tau \int_T p_{\theta}(\tau_{h+1:H}) \left(\sum_{j=h+1}^H \nabla_{\mathbf{v}} \log \pi_j \right) d\tau \\ &= \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:h}) \gamma^h r_h \left(\sum_{i=0}^h \nabla_{\omega} \log \pi_i \right) d\tau \int_T p_{\theta}(\tau_{h+1:H}) \nabla_{\mathbf{v}} \log p_{\theta}(\tau_{h+1:H}) d\tau \\ &= \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:h}) \gamma^h r_h \left(\sum_{i=0}^h \nabla_{\omega} \log \pi_i \right) d\tau \int_T \nabla_{\mathbf{v}} p_{\theta}(\tau_{h+1:H}) d\tau \\ &= \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:h}) \gamma^h r_h \left(\sum_{i=0}^h \nabla_{\omega} \log \pi_i \right) d\tau \nabla_{\mathbf{v}} \int_T p_{\theta}(\tau_{h+1:H}) d\tau \\ &= 0. \end{aligned}$$

Analogously, we can say that (50) = 0. Finally:

$$\begin{aligned} (51) &= \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:H}) \gamma^h r_h \left(\sum_{i=h+1}^H \nabla_{\omega} \log \pi_i \right) \left(\sum_{j=h+1}^H \nabla_{\mathbf{v}} \log \pi_j \right) d\tau \\ &= \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:h}) \gamma^h r_h d\tau \int_T p_{\theta}(\tau_{h+1:H}) \nabla_{\omega} \log p_{\theta}(\tau_{h+1:H}) \nabla_{\mathbf{v}} \log p_{\theta}(\tau_{h+1:H}) d\tau \\ &= \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:h}) \gamma^h r_h d\tau \int_T (\nabla_{\omega} \nabla_{\mathbf{v}} p_{\theta}(\tau_{h+1:H}) - \nabla_{\omega} \nabla_{\mathbf{v}} \log p_{\theta}(\tau_{h+1:H})) d\tau \\ &= \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:h}) \gamma^h r_h d\tau \int_T \nabla_{\omega} \nabla_{\mathbf{v}} p_{\theta}(\tau_{h+1:H}) d\tau \\ &\quad - \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:h}) \gamma^h r_h d\tau \int_T \nabla_{\omega} \nabla_{\mathbf{v}} \log p_{\theta}(\tau_{h+1:H}) d\tau \\ &= \sum_{h=0}^H \int_T p_{\theta}(\tau_{0:h}) \gamma^h r_h d\tau \nabla_{\omega} \nabla_{\mathbf{v}} \int_T p_{\theta}(\tau_{h+1:H}) d\tau \end{aligned}$$

$$\begin{aligned}
 & - \sum_{h=0}^H \int_T p_{\boldsymbol{\theta}}(\tau_{0:h}) \gamma^h r_h d\tau \int_T -2p_{\boldsymbol{\theta}}(\tau_{h+1:H}) \nabla_{\mathbf{v}} \log p_{\boldsymbol{\theta}}(\tau_{h+1:H}) d\tau \\
 & = 2 \sum_{h=0}^H \int_T p_{\boldsymbol{\theta}}(\tau_{0:h}) \gamma^h r_h d\tau \int_T p_{\boldsymbol{\theta}}(\tau) \nabla_{\mathbf{v}} \log p_{\boldsymbol{\theta}}(\tau_{h+1:H}) d\tau \\
 & = 2 \sum_{h=0}^H \int_T p_{\boldsymbol{\theta}}(\tau_{0:h}) \gamma^h r_h d\tau \int_T \nabla_{\mathbf{v}} p_{\boldsymbol{\theta}}(\tau_{h+1:H}) d\tau = 0.
 \end{aligned}$$

Hence, $\widehat{mix}(\boldsymbol{\theta})$ is equal to (48) alone, of which the proposed $\widehat{mix}(\boldsymbol{\theta})$ is a Monte Carlo estimator. \square

Similarly to what has been done for GPOMDP (Peters and Schaal, 2008), we can introduce a baseline to reduce the variance of the estimator. Let:

$$\widehat{mix}_h(\boldsymbol{\theta}) = \mathbb{E} \left[\underbrace{\left(\sum_{k=0}^h \nabla_{\mathbf{v}} \log \pi_k \right)}_{G_h} \underbrace{\left(\sum_{k=0}^h \nabla_{\omega} \log \pi_k \right)}_{H_h} \left(\underbrace{\gamma^h r_h}_{F_h} - b_h \right) \right], \quad (52)$$

where b_h is a generic baseline that is independent from actions a_k . Any baseline $b_h = \tilde{b}_h \left(\frac{G_h + H_h}{G_h H_h} \right)$ will keep the estimator unbiased for any value of \tilde{b}_h , as long as different data are used for each multiplicative term:

$$\begin{aligned}
 E[G_h H_h (F_h - b_h)] &= E[G_h H_h F_h] - E[G_h H_h b_h] \\
 &= E[G_h H_h F_h] - E \left[G_h H_h \tilde{b}_h \left(\frac{G_h + H_h}{G_h H_h} \right) \right] \\
 &= E[G_h H_h F_h] - \tilde{b}_h E[G_h] - \tilde{b}_h E[H_h] \\
 &= E[G_h H_h F_h].
 \end{aligned}$$

We choose \tilde{b}_h as to minimize the variance of \widehat{mix}_h :

$$\begin{aligned}
 Var[\widehat{mix}_h] &= E[\widehat{mix}_h^2] - E[\widehat{mix}_h]^2 \\
 &= E \left[G_h^2 H_h^2 \left(F_h^2 - 2F_h \tilde{b}_h \frac{G_h + H_h}{G_h H_h} + \tilde{b}_h^2 \frac{(G_h + H_h)^2}{G_h^2 H_h^2} \right) \right] - E[G_h H_h F_h]^2 \\
 &= E[G_h^2 H_h^2 F_h^2] - 2\tilde{b}_h E[G_h H_h F_h (G_h + H_h)] + \tilde{b}_h^2 E[(G_h + H_h)^2] \\
 &\quad - E[G_h H_h F_h]^2.
 \end{aligned}$$

Setting the gradient to zero yields:

$$b_h^* = \arg \min_{\tilde{b}_h} Var[\widehat{mix}_h] = \frac{E[G_h H_h F_h (G_h + H_h)]}{E[(G_h + H_h)^2]}.$$

Hence the estimator has minimum variance with baseline:

$$b_h = b_h^* \frac{G_h + H_h}{G_h H_h} = \frac{G_h H_h F_h (G_h + H_h)}{(G_h + H_h)^2} \frac{G_h + H_h}{G_h H_h},$$

which can be estimated from samples as in (Peters and Schaal, 2008).

Finally:

$$\begin{aligned}
 \widehat{\nabla}_\omega \|\nabla_{\mathbf{v}} J(\boldsymbol{\theta})\| &= \frac{\langle \widehat{\nabla}_{\mathbf{v}} J(\boldsymbol{\theta}), \widehat{\nabla}_\omega \nabla_{\mathbf{v}} J(\boldsymbol{\theta}) \rangle}{\|\widehat{\nabla}_{\mathbf{v}} J(\boldsymbol{\theta})\|}, \\
 &= \frac{\langle \widehat{\nabla}_{\mathbf{v}} J(\boldsymbol{\theta}), \widehat{mix}(\boldsymbol{\theta}) - 2\widehat{\nabla}_{\mathbf{v}} J(\boldsymbol{\theta}) \rangle}{\|\widehat{\nabla}_{\mathbf{v}} J(\boldsymbol{\theta})\|} \\
 &= \frac{\langle \widehat{\nabla}_{\mathbf{v}} J(\boldsymbol{\theta}), \widehat{mix}(\boldsymbol{\theta}) \rangle}{\|\widehat{\nabla}_{\mathbf{v}} J(\boldsymbol{\theta})\|} - 2.
 \end{aligned} \tag{53}$$

To preserve the unbiasedness of this estimator we need to employ three independent sets of sample trajectories: two for \widehat{mix} and a third one for the (normalized) policy gradient estimator w.r.t. \mathbf{v} . For the latter, we can re-use the same data used to compute the other additive terms of $\widehat{\nabla}_\omega \mathcal{L}$. Additional, independent data are needed for the variance-minimizing baseline. However, as often in practice, the bias introduced by using a single batch of trajectories to compute the estimator (and its baseline) is too small to justify the variance introduced by splitting the batch in order to preserve unbiasedness. Hence, we never split our batches in the experiments.

C Estimation Error Characterization

The safe step sizes for the approximate framework presented in Section 4.2 require an upper bound on the policy gradient estimator error. For simplicity and generality, in this section we will not distinguish between mean and variance parameters. Thus, we seek an $\epsilon_\delta > 0$ such that, for all $\delta \in (0, 1)$ and $N \geq 1$:

$$\|\widehat{\nabla}^N J(\boldsymbol{\theta}_t) - \nabla J(\boldsymbol{\theta}_t)\| \leq \frac{\epsilon_\delta}{\sqrt{N}}, \tag{54}$$

with probability at least $1 - \delta$, where $\widehat{\nabla}^N J(\boldsymbol{\theta}_t)$ is an *unbiased* estimator of $\nabla J(\boldsymbol{\theta}_t)$ employing N sample trajectories. A formal way to obtain such an ϵ_δ , based on Chebychev's inequality and an upper bound on the variance of GPOMDP (Zhao et al., 2011) is provided in (Papini et al., 2019). However, this solution tends to be over-conservative (Papini et al., 2017). With a small additional assumption, i.e., Gaussianity of $\widehat{\nabla} J(\boldsymbol{\theta}_t)$, we can use Gaussian confidence regions instead⁷. Since we care about the magnitude error of an m -dimensional random vector, we propose to employ ellipsoidal confidence regions. For any $\delta \in (0, 1)$, let E_δ be the following set:

$$E_\delta = \left\{ \mathbf{x} \in \mathbb{R}^m : \left(\widehat{\nabla}^N J(\boldsymbol{\theta}_t) - \mathbf{x} \right)^T S^{-1} \left(\widehat{\nabla}^N J(\boldsymbol{\theta}_t) - \mathbf{x} \right) < \frac{m}{N-m} F_{1-\delta, m, N-m} \right\}, \tag{55}$$

where S is the sample covariance of $\widehat{\nabla}^1 J(\boldsymbol{\theta}_t)$ and $F_{1-\delta, m, N-m}$ is the quantile $(1 - \delta)$ of the F-distribution with m and $n - m$ degrees of freedom. This set is centered in $\widehat{\nabla}^N J(\boldsymbol{\theta}_t)$ and is delimited by an ellipsoid. It is a standard result (Härdle and Simar, 2012) that, with probability $1 - \delta$, the true gradient is contained in this region, i.e., $\mathbb{P}(\nabla J(\boldsymbol{\theta}_t) \in E_\delta) = 1 - \delta$. Equivalently, the difference $\widehat{\nabla}^N J(\boldsymbol{\theta}_t) - \nabla J(\boldsymbol{\theta}_t)$ is contained within the following origin-centered ellipsoid:

$$\mathcal{E}_\delta = \{ \mathbf{x} \in \mathbb{R}^d : \mathbf{x}^T A_\delta \mathbf{x} = 1 \}, \tag{56}$$

where $A_\delta = \left(\frac{m F_{1-\delta, m, N-m}}{N-m} S \right)^{-1}$. Thus, the estimation error $\|\widehat{\nabla}^N J(\boldsymbol{\theta}_t) - \nabla J(\boldsymbol{\theta}_t)\|$ cannot be larger than the largest semi-axis of \mathcal{E}_δ . Simple algebraic computations yield the following:

$$\|\widehat{\nabla}^N J(\boldsymbol{\theta}_t) - \nabla J(\boldsymbol{\theta}_t)\| \leq \sqrt{\frac{m F_{1-\delta, m, N-m} \|S\|}{N-m}}, \tag{57}$$

with probability at least $1 - \delta$, where $\|S\|$ denotes the spectral norm (i.e., the largest eigenvalue) of the sample covariance. The latter can be computed efficiently with the Lanczos method (Lanczos, 1950). Finally, we define

⁷The plausibility of this assumption relies on the Central Limit Theorem, hence is only justified by sufficiently large batch sizes.

the following error bound:

$$\epsilon_\delta = \sqrt{\frac{NmF_{1-\delta,m,n-m}\|S\|}{N-m}}, \quad (58)$$

which can be directly used in Algorithm 2.

D Extensions

In this section, we consider possible generalizations of the results of the paper.

D.1 Multi-dimensional actions

In the main paper, we only considered scalar actions, i.e., $\mathcal{A} \subseteq \mathbb{R}$. However, many continuous RL tasks involve multi-dimensional actions (Brockman et al., 2016). The natural generalization of (3) for the case $\mathcal{A} \in \mathbb{R}^d$ is a multi-variate Gaussian policy:

$$\pi_\theta(a|s) = \frac{1}{(2\pi)^{d/2}|\Sigma_\omega|^{1/2}} \exp\left\{-\frac{1}{2}(a - \mu_v(s))^T \Sigma_\omega^{-1}(a - \mu_v(s))\right\}, \quad (59)$$

where Σ_ω is a $d \times d$ covariance matrix parametrized by ω . In this case the policy parameters are $\theta = [v^T | \omega^T]^T$. We denote with m_1 the dimensionality of v , with m_2 the dimensionality of ω and with $m = m_1 + m_2$ the dimensionality of the full parameter vector θ .

The simplest case is $\Sigma_\omega = e^{2\omega}\mathbb{I}$, where \mathbb{I} denotes the identity matrix. This corresponds to a factored Gaussian policy where the action dimensions are independent and the same variance is used for them all. The results of the paper extend directly to this case.

A more common parametrization (Duan et al., 2016) is $\Sigma_\omega = \text{diag}(\exp^{2\omega})$, where the covariance matrix is diagonal. This corresponds to a factored Gaussian policy where the actions dimensions are independent, but a different variance is employed for each of them. It can be useful when the actions have different scales, or when the environment is more sensitive to particular actions. Again, the results on scalar actions can be generalized quite easily by considering each action dimension separately.

Finally, Σ_ω can be a full matrix. This allows to capture correlations among different actions. A possible parametrization is $\Sigma_\omega = L_\omega L_\omega^T$, where L_ω is a lower triangular matrix with positive diagonal entries:

$$L_\omega = \begin{bmatrix} e^{\omega_{11}} & \omega_{12} & \dots & \omega_{1d} \\ \omega_{21} & e^{\omega_{22}} & \dots & \omega_{2d} \\ \vdots & \ddots & \ddots & \vdots \\ \omega_{d1} & \omega_{d2} & \dots & e^{\omega_{dd}} \end{bmatrix}. \quad (60)$$

Generalizing our results to this case is non-trivial. However, full covariance matrices are rarely used in practice.

D.2 Heteroskedastic exploration

Another generalization is to make the policy variance state-dependent. This allows to concentrate the stochasticity on those regions of the state space where there is more need of exploration. We can employ a linear parametrization as done for the mean:

$$\sigma_\omega = \exp\{\omega^T \phi(s)\}, \quad (61)$$

where $\phi(s)$ is a vector of bounded state features, i.e., $\sup_{s \in \mathcal{S}} \|\phi(s)\| \leq \varphi$. Our results extend quite easily to this case. It is enough to adjust the smoothing constants from Lemma 4.1 as follows:

$$\psi = \frac{4\varphi}{\sqrt{2\pi e}}, \quad \kappa = \xi = 2\varphi^2. \quad (62)$$

E Task Specifications

In this section, we provide a detailed description of the environments employed by the numerical simulations of Section 5.

E.1 Linear-Quadratic Gaussian regulator (LQG)

This is a 1D continuous control task. The state is initialized uniformly at random in $[-4, 4]$, which is also the state space. The task is deterministic otherwise. The action space is $[-4, 4]$ as well. The next state is $s_{h+1} = s_h + a_h$ (linear) and the reward is $r_h = -0.9s_h^2 - 0.9a_h^2$ (quadratic). The induced goal is to bring the state in 0 with the minimum effort. The episode is always 20 steps long. A discount factor of $\gamma = 0.9$ is used for this task.

E.2 Continuous-action Cart-Pole

This is a 2D continuous control task. The goal is to balance (keep upright) a pole situated on a cart, by applying forces to the cart in the horizontal direction. The cart has a mass of $1Kg$ and the pole has a mass of $0.1Kg$ and is $0.5m$ long. The state is four dimensional and includes the cart’s position x , the cart’s horizontal speed \dot{x} , the pole’s angle w.r.t. the upright position θ and the pole’s angular velocity $\dot{\theta}$. The action (force) that can be applied is $a \in [-10, 10]$. The agent receives a reward of 1 for each step. All state variables are initialized uniformly at random in $[-0.05, 0.05]$. The task is deterministic otherwise. The episode terminates when the pole falls ($|\theta| > 12$ degrees), when the cart goes too far from the initial position ($|x| > 2.4$), and anyway after 1000 time steps. The control frequency is 50Hz. A discount factor of $\gamma = 0.99$ is used for this task.

F Experimental Setting

In this section, we provide further details on the experiments for the sake of reproducibility. Table 1 provides a recap of the selected hyper-parameters.

Algorithm	α	η	τ
MEPG ($\sigma_{\omega_0} = 5.0$)	0.1	0.01	
MEPG ($\sigma_{\omega_0} = 0.5$)	1.0	0.1	
Fixed-variance PG ($\sigma_{\omega_0} = 5.0$)	1.0		
Fixed-variance PG ($\sigma_{\omega_0} = 0.5$)	0.1		
Adaptive-variance PG	0.01		
Entropy-augmented PG	0.1		0.1

Table 1: Hyper-parameters for the Cart-Pole experiment, optimized via grid search. When not specified, the same value was selected for the two initializations.

F.1 MEPG experiment (Figure 1)

Five independent random seeds were employed for the hyper-parameter selection. Average return (i.e., averaged both over learning iterations and across different seeds) was used as a metric. This metric is proportional to the area under the learning curve. Step sizes causing divergence issues were discarded. In Tables 2-5, we report the metric, averaged over the five seeds, with 95% confidence intervals, for all the candidate hyper-parameter settings. The final choice for each experiment is in bold. The sign \perp corresponds to settings that caused divergence issues and were hence discarded.

The best step sizes for MEPG were searched among $(\alpha \times \eta) \in \{10, 1, 0.1\} \times \{1, 0.1, 0.01\}$; $(0.1, 0.01)$ was the best choice for $\sigma_{\omega_0} = 5$, while $(1, 0.1)$ performed better in the case $\sigma_{\omega_0} = 0.5$ (Table 2).

		$\alpha = 10$	$\alpha = 1.0$	$\alpha = 0.1$
$\sigma = 5.0$	$\eta = 1.0$	\perp	\perp	258.38 ± 47.48
	$\eta = 0.1$	\perp	\perp	753.42 ± 14.34
	$\eta = 0.01$	\perp	733.68 ± 47.96	822.88 ± 11.41
$\sigma = 0.5$	$\eta = 1.0$	\perp	139.14 ± 32.01	381.23 ± 49.14
	$\eta = 0.1$	\perp	751 ± 27.06	238.30 ± 28.52
	$\eta = 0.01$	\perp	472.02 ± 11.71	78.42 ± 1.56

Table 2: Hyper-parameter selection of MEPG

The best step size for fixed-variance PG was searched among $\alpha \in \{10, 1, 0.1, 0.01\}$; $\alpha = 1$ turned out to be the best choice for $\sigma_{\omega_0} = 5$, while $\alpha = 0.1$ was better for $\sigma_{\omega_0} = 0.5$ (Table 3).

	$\alpha = 10$	$\alpha = 1$	$\alpha = 0.1$	$\alpha = 0.01$
$\sigma = 5.0$	\perp	472.35 ± 39.6	182.36 ± 0.80	39.57 ± 0.15
$\sigma = 0.5$	\perp	\perp	515.14 ± 4.67	114.52 ± 0.85

Table 3: Hyper-parameter selection of fixed-variance PG

The best step size for adaptive-variance PG was searched among $\alpha \in \{1, 0.1, 0.01, 0.001\}$; $\alpha = 0.01$ turned out to be the best choice both for $\sigma_{\omega_0} = 5$ and $\sigma_{\omega_0} = 0.5$ (Table 4).

	$\alpha = 1$	$\alpha = 0.1$	$\alpha = 0.01$	$\alpha = 0.001$
$\sigma = 5.0$	\perp	\perp	98.64 ± 1.16	35.40 ± 0.13
$\sigma = 0.5$	\perp	\perp	117.65 ± 1.00	42.34 ± 0.11

Table 4: Hyper-parameter selection of adaptive-variance PG

The hyper-parameters for entropy-augmented PG were searched among $(\alpha, \tau) \in \{1, 0.1, 0.01\} \times \{0.05, 0.1, 0.2\}$; $(0.1, 0.1)$ was the best choice both for $\sigma_{\omega_0} = 5$ and $\sigma_{\omega_0} = 0.5$ (Table 5).

		$\alpha = 1$	$\alpha = 0.1$	$\alpha = 0.01$
$\sigma = 5.0$	$\tau = 0.05$	22.86 ± 0.26	355.46 ± 20.13	84.91 ± 0.55
	$\tau = 0.1$	22.93 ± 0.31	371.43 ± 5.64	68.24 ± 1.08
	$\tau = 0.2$	22.43 ± 0.15	328.02 ± 11.75	50.68 ± 0.39
$\sigma = 0.5$	$\tau = 0.05$	\perp	318.54 ± 31.75	103.36 ± 0.78
	$\tau = 0.1$	\perp	382.80 ± 3.02	90.11 ± 0.70
	$\tau = 0.2$	\perp	348.92 ± 4.43	64.04 ± 0.44

Table 5: Hyper-parameter selection of entropy-augmented PG

The final results are averaged over 10 separate random seeds. The figure also reports 95% Student's t confidence intervals.

The batch size is $N = 500$ for all algorithms.

F.2 SEPG LQG experiment (Figure 2)

All the considered algorithms tune the step sizes automatically.

A batch size of $N = 500$ and a confidence parameter of $\delta = 0.2$ were used for all the algorithms.

The results are averaged over 10 random seeds. The figure also reports 95% Student's t confidence intervals.

F.3 SEPG Cart-Pole experiment (Figure 3)

SEPG tunes the step sizes automatically.

A batch size of $N = 500$ and a confidence parameter of $\delta = 1$ were used.

The results are averaged over 5 random seeds. The figure reports 95% Student's t confidence intervals. The learning curve of MEPG is also reported as a reference.

The code used for the experiments is included in the supplementary materials, including scripts for reproducing the experiments (see the README file).