

Global Convergence in Reinforcement Learning

Federico Betti, Saber Salehkaleybar, Patrick Thiran
CSE Semester Project II - EPFL Lausanne, Switzerland

January 6, 2023

Abstract

The validity of a gradient dominance property is often assumed for the convergence study of many reinforcement learning algorithms Liu et al. [2020]. A usual assumption which guarantees the latter condition in weak form is the non-degeneracy of the Fisher information matrix, together with a boundedness requirement on the transferred compatible function approximation error Masiha et al. [2022]. Under a soft-max policy, the Fisher information matrix becomes degenerate when the policy gets arbitrarily close to a greedy one; as a consequence, only a non-uniform version of the Polyak-Łojasiewicz has been established in the literature Agarwal et al. [2020]. Mei et al. [2020] has shown that the PL constant can be bounded along the trajectories of policy gradient algorithms in the deterministic case, but these algorithms require access to the full gradient of the objective which is often unfeasible or too expensive to compute; on the other hand, the study in this direction for much more practical stochastic first order and second order methods is still somehow limited. In this work, we show empirically that also in the stochastic setting a gradient dominance property holds along the trajectories of stochastic policy gradient methods and stochastic second order methods. From the obtained results, we gain further insight on the faster convergence of the stochastic cubic regularized Newton method (SCRN) over first order methods. To go further in understanding the performance of second order methods, starting from the work in Furnstion et al. [2016], we make a first attempt at analyzing the natural policy gradient method Kakade [2001] and the differences with the cubic regularized Newton’s method by comparing the second order information enclosed in the Fisher information and in the Hessian of the expected return. Finally, we draw conclusions about the work, and we present possible developments of the study presented here. Implementation details are available at <https://github.com/federicobetti99/Global-Convergence-in-RL>.

1 Introduction

The validity of a gradient dominance property is often assumed for the convergence study of many reinforcement learning algorithms Liu et al. [2020]. The latter condition is usually satisfied under a boundedness assumption on the transferred compatible function approximation error and a Fisher non-degeneracy assumption, which guarantees the agent is able to explore sufficiently well the environment under the given policy. This two assumptions are satisfied by a large class of policies, such as Gaussian policies Masiha et al. [2022]. On the other hand, under a soft-max policy, the Fisher information matrix becomes degenerate when the policy gets arbitrarily close to a greedy one; as a consequence, only a suboptimal non-uniform Łojasiewicz inequality has been shown to hold Agarwal et al. [2020], and the boundedness of the PL constant was established only in the deterministic setting Mei et al. [2020], which is poorly practical as computing the full gradient of the expected return may be unfeasible. In the much more practical stochastic setting, the empirical and theoretical study of the gradient dominance is still limited in the literature Ding et al. [2021].

In this work, we show empirically the validity of a gradient dominance property for both the non-regularized and entropy regularized expected return in a reinforcement learning setting, by testing the latter condition in different representative environments and with different stochastic algorithms. The outline of the work reads as follows: Section 2 presents in deep detail the theoretical background behind the problem at hand and the motivation for the used and tested algorithms, of both first order and second order. Section 3 presents the numerical results giving an affirmative answer for the validity of a gradient dominance property along the trajectories of the tested algorithms, where we always restrict to discrete state and action spaces. In Section 4, we present in detail the natural policy gradient algorithm Kakade [2001] and we compare the second order information enclosed in the Fisher information matrix and in the Hessian matrix. Moreover, the design of an approximate Newton solver is proposed in a similar fashion to the work of Furnstion et al. [2016]. Finally, we draw conclusions about the work in Section 5, and we discuss possible developments of the latter.

2 Theoretical Background

Reinforcement Learning Setup

Consider a discrete Markov decision process (MDP) $M = (\mathcal{S}, \mathcal{A}, \mathbb{P}, R, \rho, \gamma)$, where \mathcal{S} is the state space and \mathcal{A} is the action space. $\mathbb{P}(s'|s, a)$ denotes the probability of state transition from s to s' after taking action a and $R(\cdot, \cdot) : \mathcal{S} \times \mathcal{A} \rightarrow [-R_{\max}, R_{\max}]$ is a bounded reward function. $\rho : \mathcal{S} \rightarrow [0, 1]$ represents the initial distribution on the state space \mathcal{S} and $\gamma \in (0, 1)$ is the discount factor for future rewards. In this work, we restrict ourselves to discrete reinforcement learning environments, namely $|\mathcal{S}| < +\infty$ and $|\mathcal{A}| < +\infty$.

At each time t , the agent executes an action $a_t \in \mathcal{A}$ given the current state $s_t \in \mathcal{S}$; action selection is driven by a stochastic policy $\pi(\cdot|s) : \mathcal{A} \rightarrow \Delta(\mathcal{A})$, i.e., $a_t \sim \pi(\cdot|s_t)$. Then, given the state-action pair (s_t, a_t) , the agent observes a reward $R_t = R(s_t, a_t)$. A fundamental quantity in reinforcement learning is the state-action value function $Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ defined as

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

Similarly, one can define the value function $V : \mathcal{S} \rightarrow \mathbb{R}$ defined as

$$V(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} [Q^\pi(s, a)] \quad (2)$$

and the advantage function $A^\pi(s, a) = Q^\pi(s, a) - V(s)$.

In this work, we assume that the policy $\pi_\theta(\cdot|s)$ is parametrized for each state s by a vector $\theta_s \in \mathbb{R}^{|\mathcal{A}|}$. In particular, we assume a soft-max policy for action selection, namely actions in state s are sampled according to

$$\pi_\theta(a_i|s) = \frac{\exp\{\theta_s^T \mathbb{1}_{a_i}\}}{\sum_j \exp\{\theta_s^T \mathbb{1}_{a_j}\}}, \quad (3)$$

where $\mathbb{1}_{a_i}$ denotes a one-hot encoding vector of the i -th available action in state s , namely $\mathbb{1}_{a_i} = e_i$ where e_i is the i -th canonical vector. Hence, this ultimately reduces to

$$\pi_\theta(a_i|s) = \frac{\exp\{\theta_{s, a_i}\}}{\sum_j \exp\{\theta_{s, a_j}\}}.$$

Suppose that the initial state s_0 is drawn from some distribution ρ . Then, the reinforcement learning task is to find the optimal policy that maximizes the expected discounted return. Namely, if $\tau = \{s_t, a_t\}_{t \geq 0} \sim p(\tau|\pi_\theta)$ is a trajectory generated by the policy π_θ , with

$$p(\tau|\pi_\theta) := \rho(s_0) \prod_{t=0}^{\infty} \pi_\theta(a_t|s_t) P(s_{t+1}|s_t, a_t), \quad (4)$$

the goal is to maximize the expected return under π_θ , given by

$$J(\pi_\theta) = \mathbb{E}_{\tau \sim p(\cdot|\pi_\theta)} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]. \quad (5)$$

In the sequel, we make the assumption that π_θ is differentiable with respect to θ . For notational convenience, we denote $J(\pi_\theta)$ by $J(\theta)$ for the rest of the work.

Statement of the problem

The goal of policy-based RL is to find $\theta^* = \arg\max_\theta J(\theta)$. However, $J(\theta)$ may be non-concave and one settles instead for an ε -FOSP $\hat{\theta}$ with

$$\|\nabla_\theta J(\hat{\theta})\|_2 \leq \varepsilon. \quad (6)$$

Note that because such a $\hat{\theta}$ may not lead to a large $J(\theta)$. Hence, in general, a better goal is to have

$$J(\theta^*) - J(\theta) \leq \varepsilon + \mathcal{O}(\sqrt{\varepsilon_{\text{bias}}}), \quad (7)$$

where $\varepsilon_{\text{bias}}$ may for example reflect the inherent error related to the possibly limited expressive power of the policy parametrization π_θ . As we will soon discuss, for a soft-max policy we can take $\varepsilon_{\text{bias}} = 0$.

Standard policy gradient methods update the parameter θ of the policy by following the direction of the gradient. The gradient of the objective (5) is given by

$$\nabla J(\theta) = \mathbb{E}_{\tau \sim p(\cdot|\pi_\theta)} \left[\sum_{h=0}^{\infty} \left(\sum_{t=h}^{\infty} \gamma^t R(s_t, a_t) \right) \nabla \log \pi_\theta(a_h|s_h) \right] \quad (8)$$

with $\nabla \log \pi_\theta(a_h|s_h) = \mathbb{1}_{a_h} - \pi_\theta(\cdot|s_h)$ is the so-called *score function* under the soft-max policy (3). Note that this update rule recovers the well known paradigm of “activity minus expected activity” which often pops up in reinforcement learning optimization algorithms. Assuming that one can explicitly compute (8), this gives rise to the policy gradient update function

$$\theta \leftarrow \theta + \alpha \nabla J(\theta), \quad (9)$$

where $\alpha > 0$ is a learning rate. In practice, most of the times, the gradient computation is possibly truncated to some finite length horizon and estimated by sample batches along the trajectories taken by the agent. The combination of the two operations aforementioned gives rise to the stochastic policy gradient method (SPG): we describe the latter in more detail in the next subsection.

Before moving onto the presentation of the work, we still need to introduce some important quantities and related results. Let $d_\rho^{\pi_\theta}(\cdot)$ be the state visitation measure induced by policy π_θ and initial distribution ρ , which is defined as

$$d_\rho^{\pi_\theta}(s) = (1 - \gamma) \mathbb{E}_{s_0 \sim \rho} \left[\sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s | s_0, \pi_\theta) \right], \quad (10)$$

and let $\nu_\rho^{\pi_\theta}(\cdot, \cdot)$ the state-action visitation measure induced by the policy π_θ and initial distribution ρ , which is defined as

$$\nu_\rho^{\pi_\theta}(s, a) = (1 - \gamma) \mathbb{E}_{s_0 \sim \rho} \left[\sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s, a_t = a | s_0, \pi_\theta) \right], \quad (11)$$

Then, the transferred compatible function approximation error is defined by

$$L_{\nu_\rho^{\pi_\theta}(s,a)}(w; \theta) = \mathbb{E}_{(s,a) \sim \nu_\rho^{\pi_\theta}(s,a)} \left[(A^\pi(s, a) - (1 - \gamma)w^T \nabla \log \pi_\theta(s, a))^2 \right]. \quad (12)$$

Definition 2.1. (*Fisher information matrix*). The Fisher information matrix induced by a stochastic policy π_θ and an initial state distribution ρ is defined as

$$F_\rho(\theta) = \mathbb{E}_{(s,a) \sim \nu_\rho^{\pi_\theta}(s,a)} \left[\nabla \log \pi_\theta(a|s) \nabla \log \pi_\theta(a|s)^T \right]. \quad (13)$$

Definition 2.2. (*Fisher non-degeneracy*). A Fisher non-degeneracy condition is said to hold true if for all $\theta \in \mathbb{R}^d$, there exists some constant $\mu_J > 0$ such that the Fisher information matrix $F_\rho(\theta)$ defined in (13) induced by the policy π_θ and the initial state distribution ρ satisfies

$$F_\rho(\theta) \succcurlyeq \mu_J Id. \quad (14)$$

Remark 2.1. Condition (14) implies that all the singular values of $F_\rho(\theta)$ are at least μ_J . Equivalently, we have $\|F_\rho(\theta)^{-1}\|_2 \leq \frac{1}{\mu_J}$.

Definition 2.3. (*Gradient dominance property*). Let $J(\theta) : \mathbb{R}^d \rightarrow \mathbb{R}$ a differentiable function in θ . Then J is said to satisfy a gradient dominance property if there exists $\alpha \geq 1$ such that for any $\theta \in \mathbb{R}^d$ there exists a constant τ_J and $\varepsilon' > 0$ such that

$$J(\theta^*) - J(\theta) \leq \tau_J \|\nabla J(\theta)\|_2^\alpha + \varepsilon'. \quad (15)$$

Definition 2.4. (*Weak/strong gradient dominance property*). Let $J(\theta) : \mathbb{R}^d \rightarrow \mathbb{R}$ a differentiable function in θ . Then J is said to satisfy a weak gradient dominance property if (15) holds with $\alpha = 1$, namely

$$J(\theta^*) - J(\theta) \leq \tau_J \|\nabla J(\theta)\|_2 + \varepsilon', \quad (16)$$

and J is said to satisfy a strong gradient dominance property if (15) holds with $\alpha = 2$, namely

$$J(\theta^*) - J(\theta) \leq \tau_J \|\nabla J(\theta)\|_2^2 + \varepsilon'. \quad (17)$$

Assumption 2.1. (*Transferred compatible function approximation error*). Let $u_* = (F_\rho(\theta))^\dagger \nabla J(\theta)$, where A^\dagger denotes the Moore-Penrose inverse Barata and Hussein [2012] of a matrix A . For all $\theta \in \mathbb{R}^d$, there exists $\varepsilon_{\text{bias}} > 0$ such that the transferred compatible function approximation error with $(s, a) \sim p(\cdot|\pi_\theta)$ satisfies

$$\mathbb{E}_{(s,a) \sim \nu_\rho^{\pi_\theta}(s,a)} \left[(A^{\pi_\theta}(s, a) - (1 - \gamma)u_*^T \nabla \log \pi_\theta(a|s))^2 \right] \leq \varepsilon_{\text{bias}}. \quad (18)$$

Lemma 2.1. *For a softmax tabular policy, Assumption 2.1 holds true with $\varepsilon_{\text{bias}} = 0$.*

Proof. It is straightforward to verify that for a soft-max policy parametrization

$$A^{\pi_\theta}(s, a) - (1 - \gamma)u_*^T \nabla \log \pi_\theta(a|s) = 0.$$

□

Theorem 2.1. *Let $\|\nabla \log \pi_\theta(a|s)\|_2 \leq G_1$ for some constant $G_1 > 0$. If Fisher non-degeneracy (14) holds true with lower bound $\mu_J > 0$ and Assumption 2.1 holds true with bias term $\varepsilon_{\text{bias}}$, a weak gradient dominance property holds true. In particular,*

$$J(\theta^*) - J(\theta) \leq \frac{G_1}{\mu_J} \|\nabla J(\theta)\| + \frac{1}{1 - \gamma} \sqrt{\varepsilon_{\text{bias}}}.$$

Proof. By the performance difference lemma Kakade and Langford [2002] we have

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

Using Assumption 2.1 and Jensen's inequality, we have

$$\begin{aligned} \varepsilon_{\text{bias}} &\geq \mathbb{E}_{(s,a) \sim \nu_{\rho}^{\pi_\theta}(s,a)} [(A_{\pi_\theta}(s, a) - (1 - \gamma)u_*^T \nabla \log \pi_\theta(a|s))^2] \\ &\geq \left(\mathbb{E}_{(s,a) \sim \nu_{\rho}^{\pi_\theta}(s,a)} [A_{\pi_\theta}(s, a) - (1 - \gamma)u_*^T \nabla \log \pi_\theta(a|s)] \right)^2, \end{aligned}$$

where $u_* = (F_\rho(\theta))^{-1} \nabla J(\theta)$. Thus

$$\begin{aligned} \sqrt{\varepsilon_{\text{bias}}} &\geq \mathbb{E}_{(s,a) \sim \nu_{\rho}^{\pi_\theta}(s,a)} [A_{\pi_\theta}(s, a) - (1 - \gamma)u_*^T \nabla \log \pi_\theta(a|s)] \\ &= (1 - \gamma)(J(\theta^*) - J(\theta)) - (1 - \gamma)\mathbb{E}_{(s,a) \sim \nu_{\rho}^{\pi_\theta}(s,a)} [u_*^T \nabla \log \pi_\theta(a|s)]. \end{aligned}$$

Rearranging terms gives

$$\begin{aligned} J(\theta^*) - J(\theta) &\leq \mathbb{E}_{(s,a) \sim \nu_{\rho}^{\pi_\theta}(s,a)} [u_*^T \nabla \log \pi_\theta(a|s)] + \frac{1}{1 - \gamma} \sqrt{\varepsilon_{\text{bias}}} \\ &\leq \mathbb{E}_{(s,a) \sim \nu_{\rho}^{\pi_\theta}(s,a)} [\|u_*^T\| \|\nabla \log \pi_\theta(a|s)\|] + \frac{1}{1 - \gamma} \sqrt{\varepsilon_{\text{bias}}} \\ &\leq G_1 \|u_*\| + \frac{1}{1 - \gamma} \sqrt{\varepsilon_{\text{bias}}} \leq \frac{G_1}{\mu_J} \|\nabla J(\theta)\| + \frac{1}{1 - \gamma} \sqrt{\varepsilon_{\text{bias}}}. \end{aligned}$$

Thus, the claim follows. □

Remark 2.2. *For a soft-max policy parametrization, we have $\|\nabla \log \pi_\theta(a|s)\| \leq 2$ for every $(s, a) \in \mathcal{S} \times \mathcal{A}$. Furthermore, we already discussed in Lemma 2.1 that $\varepsilon_{\text{bias}} = 0$. This implies that we can concentrate on the validity of an inequality of the form $J(\theta^*) - J(\theta) \leq \tau_J \|\nabla J(\theta)\|_2^\alpha$ for some $\alpha \geq 1$.*

The Fisher non-degeneracy assumption (14) is commonly used in the literature. The Fisher-non-degenerate setting implicitly guarantees that the agent is able to explore the state-action space under the considered policy class. Fisher non-degeneracy can be satisfied by certain Gaussian policies, where $\pi_\theta(\cdot|s) = \mathcal{N}(\mu_\theta(s), \Sigma)$ with the parametrized mean function $\mu_\theta(s)$ and the fixed covariance matrix $\Sigma \succcurlyeq 0$, provided that the Jacobian of $\mu_\theta(s)$ is full-row rank for all $\theta \in \mathbb{R}^d$. In addition, such an assumption holds more generally for every full-rank exponential family parametrization with their mean parameterized by $\mu_\theta(s)$ if $\mu_\theta(s)$ is full-row rank for all $\theta \in \mathbb{R}^d$. On the other hand, we know that for the softmax policy, the Fisher information matrix becomes degenerate as soon as the policy becomes deterministic, which in turn means that close to the optimum the constant τ_J may explode or not be well defined. Without the Fisher non-degeneracy assumption, establishing a gradient dominance property or any convergence property of the algorithms becomes difficult.

Still, previous work in Mei et al. [2020] has shown that a non-uniform version of the Łojasiewicz inequality (16) holds true under a soft-max tabular policy for action selection.

Theorem 2.2. *Assume that the initial state distribution satisfies $\min_s \rho(s) > 0$. Furthermore, denote by $a^*(s)$ the optimal action in state s under a deterministic policy π_{θ^*} . Then it holds that*

$$\|\nabla J(\theta_t)\|_2 \geq \frac{\min_s \pi_{\theta_t}(a^*(s)|s)}{\sqrt{|\mathcal{S}|} \|d_{\rho}^{\pi_{\theta^*}}/d_{\rho}^{\pi_{\theta_t}}\|_\infty} (J(\theta^*) - J(\theta_t)). \quad (19)$$

Assuming $\min_s \rho(s) > 0$, we guarantee that $\|d_{\rho}^{\pi_{\theta^*}}/d_{\rho}^{\pi_{\theta_t}}\|_{\infty} > 0$, and hence the denominator of the PL constant is well defined. However, this gradient dominance inequality is cumbersome and suboptimal as it is not possible, in general, to bound from below the optimal action probabilities given by $\min_s \pi_{\theta}(a^*(s)|s)$. The difficulties in doing so, especially in the stochastic setting, were already highlighted in Agarwal et al. [2020]. For the deterministic case, when the exact PG is available, Mei et al. [2020] have shown that along the trajectories of the policy gradient algorithm for a softmax policy, $\min_s \pi_{\theta}(a^*(s)|s) > 0$ and a weak gradient dominance property holds true.

Lemma 2.2. *Using the update rule (9), it holds that $c = \inf_{s \in \mathcal{S}, t \geq 1} \pi_{\theta_t}(a^*(s)|s) > 0$.*

Similar results can be shown for the entropy regularized policy gradient methods, which we now briefly introduce. To penalize deterministic policies, entropy regularized reinforcement learning considers a surrogate objective

$$J^{\lambda}(\theta) = J(\theta) + \lambda \mathcal{H}(p, \pi_{\theta}), \quad (20)$$

for a regularization parameter $\lambda > 0$ which determines the strength of the penalty and $J(\theta)$ the usual objective defined in (5). In (20), $\mathcal{H}(p, \pi)$ is the discounted entropy defined as

$$\mathcal{H}(p, \pi) = \mathbb{E}_{\tau \sim p(\cdot|\pi_{\theta})} \left[\sum_{t=0}^{\infty} -\gamma^t \log \pi_{\theta}(a_t|s_t) \right].$$

For an entropy regularized objective Grandvalet and Bengio [2006], a better version of the non-uniform Lojasiewicz-inequality (15) holds true.

Lemma 2.3. *Let $J^{\lambda}(\theta)$ denote the entropy regularized objective (20), and assume that $\rho(s) > 0$ for all $s \in \mathcal{S}$. Then it holds that*

$$\|\nabla J^{\lambda}(\theta_t)\|_2^2 \geq \frac{2\lambda \min_s \rho(s) \min_{s,a} (\pi_{\theta_t}(a|s))^2}{|\mathcal{S}| \|d_{\rho}^{\pi_{\theta_t}}/d_{\rho}^{\pi_{\theta^*}}\|_{\infty}} (J^{\lambda}(\theta^*) - J^{\lambda}(\theta_t)), \quad (21)$$

where the index λ for π_{θ} denotes the fact that $\pi_{\theta_{\lambda}^*}$ is the policy which optimizes the entropy regularized objective.

Remark 2.3. *Due to the presence of regularization, the optimal solution will be biased with the bias disappearing as $\lambda \rightarrow 0$. More precisely, the optimal policy $\pi_{\theta_{\lambda}^*}$ of the entropy-regularized problem could also be nearly optimal in terms of the unregularized objective function, as long as the regularization parameter λ is chosen to be small. In particular, concerning the difference between the two objectives (5) and (20), one can show that for any initial state distribution ρ it holds that*

$$J(\theta_{\lambda}^*) \leq J(\theta^*) \leq J(\theta_{\lambda}^*) + \lambda \frac{\log |\mathcal{A}|}{1 - \gamma},$$

where θ_{λ}^* is the maximizer of the regularized objective.

Lemma 2.4. *Assume that the initial state distribution satisfies $\min_s \rho(s) > 0$. Then along the trajectories $\{\theta_t\}$ generated by the entropy regularized policy gradient method it holds that $c := \inf_{t \geq 1} \min_{s,a} \pi_{\theta_t}(a|s) > 0$.*

Remark 2.4. *The main difference to the previous versions of the non-uniform Lojasiewicz inequality is that the sub-optimality gap appears under the square root, or equivalently the gradient norm appears squared. For small sub-optimality gaps this means that the gradient must be larger, i.e. a stronger “signal”.*

Having established under which conditions a gradient dominance property holds true with $\alpha = 1$ for the objective (5) and with $\alpha = 2$ for the entropy regularized objective (20), we now introduce some technical definitions to claim the convergence rates proved in Mei et al. [2020] upon the validity of (15). A commonly used assumption which is satisfied by the expected return (20) for $\lambda \neq 0$ and $\lambda = 0$ is a smoothness assumptions, which guarantees boundedness of the Hessian norm. We show that this condition holds true for the expected return under a soft-max policy. Here, for space reasons, we only show basic results on the log-policy function (note that this is sufficient) and refer the interested reader to Agarwal et al. [2021].

Definition 2.5. *(Smooth function). A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is said to be smooth if there exists a constant $L > 0$ such that for any $x, y \in \mathbb{R}^d$ it holds that*

$$f(x) \leq f(y) + \nabla f(y)^T (x - y) + \frac{L}{2} \|x - y\|_2^2. \quad (22)$$

For the expected return (5) without any entropy regularization, it is easy to show that

$$\frac{\partial \log \pi_\theta(\cdot|s)}{\partial \theta_s} = \mathbb{1}_a - \pi_\theta(\cdot|s), \quad \frac{\partial^2 \log \pi_\theta(\cdot|s)}{\partial \theta_s^2} = \pi_\theta(\cdot|s)\pi_\theta(\cdot|s)^T - \text{Diag}(\pi_\theta(\cdot|s)). \quad (23)$$

This immediately gives the following results.

Theorem 2.3. *The log-policy function is 2-smooth, namely for any $\theta, \theta' \in \mathbb{R}^d$ we have*

$$\|\nabla \log \pi_\theta(\cdot|s) - \nabla \log \pi_{\theta'}(\cdot|s)\| \leq 2\|\theta - \theta'\|_2. \quad (24)$$

Proof. We compute

$$\|\nabla \log \pi_\theta(\cdot|s) - \nabla \log \pi_{\theta'}(\cdot|s)\| = \|\mathbb{1}_a - \pi_\theta(\cdot|s) - \mathbb{1}_a + \pi_{\theta'}(\cdot|s)\| = \|\pi_\theta(\cdot|s) - \pi_{\theta'}(\cdot|s)\| \stackrel{(a)}{\leq} 2\|\theta - \theta'\|_2 \quad (25)$$

where at (a) we used the fact that $\pi_\theta(\cdot|s)$ is 2-Lipschitz. \square

Lemma 2.5. *The objective (5) has smoothness constant $L = \frac{8}{(1-\gamma)^3}$, namely*

$$\|\nabla J(\theta_1) - \nabla J(\theta_2)\|_2 \leq \frac{8}{(1-\gamma)^3} \|\theta_1 - \theta_2\|_2$$

Remark 2.5. *Similar results can be proven for the entropy regularized objective which we introduced in (20). In particular, the smoothness constant in the latter case is $\frac{8+4\lambda(1+2\log|\mathcal{A}|)}{(1-\gamma)^3}$, which recovers the previous result for a temperature parameter $\lambda = 0$.*

Combining Theorem 2.2 and Lemma 2.2, we conclude that a weak gradient dominance holds along the trajectories followed by policy gradient methods. Upon this condition and L -smoothness of the objective, Mei et al. [2020] show that policy gradient converges to the optimal policy at a rate $\mathcal{O}(\frac{1}{t})$. Similarly, a strong gradient dominance property holds (namely with $\alpha = 2$) for entropy regularized policy gradient methods from Theorem 2.3 and Lemma 2.4. Together with smoothness of the objective function, taking a learning rate $\alpha_t \leq \frac{1}{L}$ yields a convergence rate $\mathcal{O}(e^{-ct})$ to the optimum.

In practice, the knowledge of the full gradient of the objective is often not available, so one uses instead sample batches from the environment to build unbiased estimates of the gradient to be used in the update of the policy parameters θ . However, to our knowledge, for stochastic methods the validity of a gradient dominance property has not been established nor studied exhaustively under a soft-max policy for action selection. As a consequence, in this work we are interested in studying empirically and theoretically whether, along the trajectories of the much more practical stochastic first order and second order methods, the PL constant $\tau_J(\theta)$ can be bounded from above. Namely, we study whether for a sequence $\{\theta_t\}$ generated by a stochastic optimization algorithm it holds that for all $t \geq 1$

$$J(\theta^*) - J(\theta_t) \leq \tau_J(\theta_t) \|\nabla J(\theta_t)\|,$$

for $\lambda = 0$ in (20), or whether

$$J^\lambda(\theta^*) - J^\lambda(\theta_t) \leq \tau_{J^\lambda}(\theta_t) \|\nabla J^\lambda(\theta_t)\|^2$$

for the entropy regularized objective (20), for a positive constant $\tau_J(\theta_t) \leq C$ bounded from above by some $C > 0$.

Remark 2.6. *For a soft-max policy, we already established previously that $\varepsilon' = 0$ in (16). Hence, we can study directly only the boundedness of the PL constant in front of $\|\nabla J(\theta)\|^\alpha$ for some coefficient $\alpha \geq 1$.*

In particular, in this work we focus on the comparison and the validity of (15) in a weak or a strong form for three stochastic algorithms, plus a variance reduced variant: vanilla stochastic policy gradient (SPG), entropy regularized vanilla stochastic policy gradient (ESPG), two stages entropy regularized stochastic policy gradient Ding et al. [2021] and stochastic cubic regularized Newton (SCRN) Tripura-
neni et al. [2018]. For an exhaustive presentation of the stochastic policy gradient method, we refer the interested reader to Yuan et al. [2022], and similarly for the entropy regularized objective function Grandvalet and Bengio [2006]. The two stages entropy regularized SPG algorithm chooses a large enough batch size in the very first episodes to escape flat regions and to mitigate for the non-coercive landscape of the objective function: indeed, the search direction may be dominated by the gradient estimation noise at the region where the landscape is highly flat. On the other hand, after a certain number of iterations, which is a constant with respect to the optimality gap, the iteration will reach a region where

the landscape has enough curvature information: as a consequence, a smaller batch size at the end of training is sufficient to guarantee a fast convergence to the optimal policy. For this algorithm, a strong gradient dominance property holds with high probability. For space reasons and to focus on the main results of the work, we omit the derivation of such a result, referring the interested reader to Ding et al. [2021]. What is important for us is that, if we have a theoretical guarantee that the action probabilities are bounded away from zero under a softmax policy with high probability, comparing the validity of a gradient dominance property for the other stochastic algorithms we consider in this work with respect to how the PL constant evolves for the latter algorithm can give nice insight on the effective validity of the property along their trajectories. While we do not present a detailed presentation for the two stages algorithm, in the next subsection, starting from the use of second order methods in machine learning tasks, we introduce the SCR algorithm and its motivation.

Stochastic Cubic Regularized Newton

In machine learning tasks, we can often recast the training goal as the minimization of an empirical risk, namely we want to solve the optimization problem

$$\min_x f(x), \quad x \in \mathbb{R}^d. \quad (26)$$

For simplicity, we assume the problem to be unconstrained, but similar reasoning apply to constrained optimization tasks. To solve (26), second order methods have been shown to outperform first order methods in a variety of machine learning tasks. However, their main bottlenecks are usually the expensive access to the second order information and the higher requirements on the objective function.

Definition 2.6. (*Strongly convex function*). A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is said to be strongly convex if there exists a constant $\mu > 0$ such that for any $x, y \in \mathbb{R}^d$ it holds

$$f(x) \geq f(y) + \nabla f(y)^T(x - y) + \frac{\mu}{2}\|x - y\|_2^2. \quad (27)$$

Assumption 2.2. (*Lipschitz Hessians*). There exists a constant $\rho > 0$ such that

$$\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq \rho\|x - y\|_2, \quad (28)$$

for all $x, y \in \mathbb{R}^d$.

Remark 2.7. If a twice differentiable function f is μ -strongly convex, this implies in particular that for every $x \in \mathbb{R}^d$ $\|\nabla^2 f(x)^{-1}\| \leq \frac{1}{\mu}$. Hence, the Hessian is always non-degenerate and the Newton step

$$x_{t+1} = x_t - (\nabla^2 f(x_t))^{-1} \nabla f(x_t) \quad (29)$$

is always well defined.

For a strongly convex objective with ρ -Lipschitz Hessians $f : \mathbb{R}^d \rightarrow \mathbb{R}$, Nesterov and Polyak have shown that Newton's method converges to the global minimum x^* at a quadratic rate, namely

$$f(x_{t+1}) - f(x^*) \leq \frac{\rho}{2}\|x_t - x^*\|_2^2, \quad (30)$$

where $\{x_t\}$ is the sequence generated by the method. Even when requiring less regularity on the objective, second order information boosts training. Many works, see e.g. Battiti [1992], have been trying to understand deeply the reasons behind this. Empirically, they seem to avoid saddle points more efficiently than gradient-based methods because of the curvature information which is extrapolated from the landscape at each iteration.

A cubic regularized version of the Newton's method (29) was introduced by Nesterov and Polyak who showed a local convergence rate on functions satisfying a strong gradient dominance property Nesterov and Polyak [2006], namely condition (17). We now introduce the main idea behind the algorithm before moving to its application in a reinforcement learning setting. At every step, gradient descent finds the minimizer of a local second-order Taylor expansion at each step, i.e.

$$x_{t+1}^{GD} = \operatorname{argmin}_x \left[f(x_t) + \nabla f(x_t)^T(x - x_t) + \frac{L}{2}\|x - x_t\|_2^2 \right]. \quad (31)$$

Note that in (31), the second order term of the objective minimized at every step is only a very loose upper bound for the curvature information embedded in the landscape of f . Differentiating the above in x gives rise to the usual update rule

$$x_{t+1}^{GD} = x_t - \frac{1}{L} \nabla f(x_t).$$

Newton's Method replaces the term $\frac{L}{2} \|x - x_t\|^2$ in (31) with the tighter second order approximation given by $(x - x_t)^T \nabla^2 f(x_t)(x - x_t)$. On the other hand, the cubic regularized Newton method finds the minimizer of a local third-order Taylor expansion given by

$$x_{t+1}^{CRN} = \operatorname{argmin}_x \left[f(x_t) + \nabla f(x_t)^T (x - x_t) + \frac{1}{2} (x - x_t)^T \nabla^2 f(x_t)(x - x_t) + \frac{\rho}{6} \|x - x_t\|^3 \right], \quad (32)$$

where ρ is the Lipschitz constant for the Hessian of f . Unlike the previous cases, this problem does not admit a closed form solution, but the optimal x is usually approached by gradient descent on the objective (32) instead. We discuss this more in detail in the stochastic setting below. For the moment, we show that Assumption 2.2 is satisfied by the log-policy function.

Theorem 2.4. *The log-policy function has 6-Lipschitz Hessians, i.e. for any $\theta, \theta' \in \mathbb{R}^d$ we have*

$$\|\nabla^2 \log \pi_\theta(\cdot|s) - \nabla^2 \log \pi_{\theta'}(\cdot|s)\| \leq 6\|\theta - \theta'\|_2. \quad (33)$$

Proof. We compute

$$\begin{aligned} \|\nabla^2 \log \pi_\theta(\cdot|s) - \nabla^2 \log \pi_{\theta'}(\cdot|s)\| &= \|\pi_\theta(\cdot|s)\pi_\theta(\cdot|s)^T - \operatorname{Diag}(\pi_\theta(\cdot|s)) - \pi_{\theta'}(\cdot|s)\pi_{\theta'}(\cdot|s)^T + \operatorname{Diag}(\pi_{\theta'}(\cdot|s))\| \\ &\leq \|\pi_\theta(\cdot|s)\pi_\theta(\cdot|s)^T - \pi_{\theta'}(\cdot|s)\pi_{\theta'}(\cdot|s)^T\| + \|\operatorname{Diag}(\pi_\theta(\cdot|s)) - \operatorname{Diag}(\pi_{\theta'}(\cdot|s))\| \\ &\leq \|\pi_\theta(\cdot|s)\| \|\pi_\theta(\cdot|s) - \pi_{\theta'}(\cdot|s)\| + \|\pi_{\theta'}(\cdot|s)\| \|\pi_\theta(\cdot|s) - \pi_{\theta'}(\cdot|s)\| \\ &\leq 4\|\theta - \theta'\|_2 + 2\|\theta - \theta'\|_2 = 6\|\theta - \theta'\|_2 \end{aligned}$$

where to bound the second term we used Theorem 2.3. \square

Remark 2.8. *Theorem 2.4 can be used to verify that the objective $J(\theta)$ has Lipschitz Hessians with constant $\rho(\gamma) = \mathcal{O}\left(\frac{1}{(1-\gamma)^3}\right)$. For space reasons and to illustrate only the main parts of the works, we refer the interested reader to Masiha et al. [2022] for the full derivation.*

Given the smoothness assumptions verified for the objective at hand in a reinforcement learning setting in Theorem 2.3, we now pose our attention to the presentation of the cubic regularized Newton algorithm. Most previous work on the cubic-regularized Newton method has focused on the non-stochastic or partially-stochastic settings. In this work, we are interested in its stochastic version introduced in Tripuraneni et al. [2018]. In this setting, we can approach the solution to (32) by solving alternatively the *stochastic cubic submodel* given by

$$x_{t+1}^{SCRN} = \operatorname{argmin}_x m(x, x_t) \quad (34)$$

where for a given unbiased estimate g_t of $\nabla f(\theta)$ and B_t of $\nabla^2 f(\theta)$ we define

$$m(x, x_t) = \operatorname{argmin}_x \left[f(x_t) + g_t^T (x - x_t) + \frac{1}{2} (x - x_t)^T B_t (x - x_t) + \frac{\rho}{6} \|x - x_t\|^3 \right]. \quad (35)$$

In general, as there is no closed form solution for the solution to (32), the same holds for (34). Thus, one can employ at each iteration gradient descent to approximate the optimal x . For an exhaustive description of the algorithm and convergence results in the resolution of (34), we refer the interested reader to the pioneering work Tripuraneni et al. [2018]. Here, for space reasons, we restrict to presenting the pseudocode of the full algorithm in Algorithm 1. Similarly, the pseudocode for the Cubic Subsolver to approximate the solution to (34) is given in Algorithm 2. In both pseudocodes, the notation is the same used in the above, in particular L is always a Lipschitz constant for the gradient of the objective function and ρ is the Lipschitz constant for the Hessian of the objective function.

Let us now discuss the details in building unbiased estimates of the gradient and of the Hessian of the expected return (5) which we consider for the task at hand. In this work, we deal with the problem of stochastic optimization of a non-concave objective

$$\max_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p(\cdot|\pi_\theta)} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]. \quad (36)$$

Algorithm 1: Stochastic Cubic Regularized Newton with stopping criterion

Data: Batch sizes n_1, n_2 , initial guess x_0 , accuracy ε , cubic penalty parameter ρ , maximum number of iterations T

```
1  $t \leftarrow 1$ 
2  $\|\Delta_0\| = \infty$ 
3 while  $\|\Delta_{t-1}\| \geq \sqrt{\varepsilon}$  or  $t \leq T$  do
4    $g_t \leftarrow \frac{1}{n_1} \sum_{i=1}^{n_1} \nabla f(x_t, \xi_i)$   $\triangleright$  Sample gradients
5    $B_t \leftarrow \frac{1}{n_2} \sum_{i=1}^{n_2} \nabla^2 f(x_t, \xi_i)$   $\triangleright$  Sample Hessians
6    $\Delta_t = \operatorname{argmin}_{\Delta \in \mathbb{R}^d} \langle g_t, \Delta \rangle + \frac{1}{2} \langle \Delta, B_t \Delta \rangle + \frac{\rho}{6} \|\Delta\|^3$   $\triangleright$  Cubic Subsolver( $g_t, B_t[\cdot], \varepsilon$ )
7    $x_{t+1} = x_t + \Delta_t$   $\triangleright$  Update  $x$ 
8    $t \leftarrow t + 1$ 
9 end
```

Algorithm 2: Cubic Subsolver

Data: $g, B[\cdot], \varepsilon$

```
1 if  $\|g\| \geq \frac{L^2}{\rho}$  then
2    $R_c \leftarrow -\frac{g^T B g}{\rho \|g\|^2} + \sqrt{\frac{g^T B g}{\rho \|g\|^2} + \frac{2\|g\|}{\rho}}$ 
3    $\Delta \leftarrow -R_c \frac{g}{\|g\|}$ 
4 else
5    $\Delta \leftarrow 0, \sigma \leftarrow c' \frac{\sqrt{\varepsilon \rho}}{L}, \eta \leftarrow \frac{1}{20L}$ 
6    $\tilde{g} \leftarrow g + \sigma \xi, \xi \sim \mathcal{U}(\mathbb{S}^{d-1})$ 
7   for  $t = 1, \dots, T(\varepsilon)$  do
8      $\Delta \leftarrow \Delta - \eta(\tilde{g} + B\Delta + \frac{\rho}{2}\|\Delta\|\Delta)$ 
9   end
10 end
```

Because, in most applications, one cannot compute (8), the latter quantity is commonly truncated up to a length horizon H , namely one considers instead

$$\nabla J_H(\theta) = \mathbb{E} \left[\sum_{h=0}^{H-1} \left(\sum_{t=h}^{H-1} \gamma^t R(s_t, a_t) \right) \nabla \log \pi_\theta(a_h | s_h) \right]. \quad (37)$$

The truncation bias introduced in this step can be fully characterized as long as the log-policy has bounded gradients, i.e. there exists a constant $G_1 > 0$ such that $\|\nabla \log \pi_\theta\|_2 \leq G_1$. This was already verified for the soft-max policy in Theorem 2.3 with $G_1 = 2$. Upon this assumptions, Masiha et al. [2022] show that

$$\|\nabla J(\theta) - \nabla J_H(\theta)\|_2 \leq \frac{G_1}{1-\gamma} \left(\frac{1}{1-\gamma} + H \right)^{1/2} \gamma^H.$$

To construct an estimate of $\nabla J_H(\theta)$, we sample m trajectories $\tau^i = \{s_t^i, a_t^i\}_{t \geq 0}$, $1 \leq i \leq m$. Upon this, we then define an unbiased estimator of $\nabla J_H(\theta)$ as

$$g := \hat{\nabla}_m J_H(\theta) = \frac{1}{m} \sum_{i=1}^m \sum_{h=0}^{H-1} \left(\sum_{t=h}^{H-1} \gamma^t R(s_t^i, a_t^i) \right) \nabla \log \pi_\theta(a_h^i | s_h^i). \quad (38)$$

Stochastic policy gradient now simply performs the update $\theta \leftarrow \theta + \alpha g$. On the other hand, stochastic cubic regularized Newton needs an unbiased estimator of the Hessian matrix $\nabla^2 J(\theta)$ to approximate the solution of the cubic submodel. It can be shown Shen et al. [2019] that the expression for the H -truncated Hessian expression $\nabla^2 J_H(\theta)$ reads as

$$\nabla^2 J_H(\theta) = \mathbb{E} [\nabla \phi(\theta; \tau) \nabla \log p(\tau | \pi_\theta)^T + \nabla^2 \phi(\theta; \tau)], \quad (39)$$

where $\phi(\theta; \tau) = \sum_{h=0}^{H-1} \left(\sum_{t=h}^{H-1} \gamma^t R(s_t, a_t) \right) \log \pi_\theta(a_h | s_h)$. Under the assumption of bounded Hessians, i.e. there exists a constant G_2 such that $\|\nabla^2 \log \pi_\theta\|_2 \leq G_2$, we have that the truncation bias satisfies

$$\|\nabla^2 J(\theta) - \nabla^2 J_H(\theta)\|_2 \leq \frac{G_2 + G_1^2}{1-\gamma} \left(\frac{1}{1-\gamma} + H \right) \gamma^H.$$

Again, for a softmax policy we have $G_2 = 1$, so the previous result applies. Finally, an unbiased estimator of $\nabla^2 J_H(\theta)$ can be constructed as

$$B := \hat{\nabla}_m^2 J_H(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla \phi(\theta; \tau^i) \nabla \log p(\tau^i | \pi_\theta) + \nabla^2 \phi(\theta; \tau^i). \quad (40)$$

Furthermore, one can show that

$$\mathbb{E} [\|\nabla J_H(\theta) - \nabla_m J_H(\theta)\|_2^2] \leq \frac{HG_1}{m(1-\gamma)^2}, \quad \mathbb{E} [\|\nabla^2 J_H(\theta) - \nabla_m^2 J_H(\theta)\|_2^2] \leq \frac{4e(H^2 G_1^4 + G_2^2)}{m(1-\gamma)^4}.$$

From the triangular inequality, i.e.

$$\|\nabla J(\theta) - \nabla_m J_H(\theta)\|_2 \leq \|\nabla J_H(\theta) - \nabla_m J_H(\theta)\|_2 + \|\nabla J(\theta) - \nabla J_H(\theta)\|_2$$

and exploiting the previous results, one can derive an expression for the expected deviation of the estimators (38) and (40) from the true expressions $\nabla J(\theta)$ and $\nabla^2 J(\theta)$.

Now that we defined how to properly deal with the stochasticity of the the gradients and of the Hessians via the definition of the stochastic unbiased estimators (38) and (40), one can solve (34) using the latter quantities to approach the exact solution to (32).

Concerning the performance of SCRNN in a reinforcement learning setting, it has been shown in Masiha et al. [2022] that if a weak gradient dominance property holds, SCRNN improves the best-known sample complexity of stochastic gradient descent in achieving ε -global optimum by a factor of $\mathcal{O}(\varepsilon^{-1/2})$.

3 Numerical results

In this section, we present numerical evidence for the validity of a gradient dominance property in the stochastic setting. For all the reinforcement learning environments considered in this section, the initial state for the agent is fixed. One could enforce exploration by randomly initialize the agent among all the passable states of the environment: we postpone the discussion about these aspects to future work. To estimate $J(\theta)$ and $\nabla J(\theta)$ given a current parameter value θ , we average over 100 testing policy episodes (with a frequency of 50 episodes during training) under a fixed policy π_θ . Estimating the objective and the gradient by testing the policy π_θ for a sufficient number of episodes comes as a cheaper alternative to the computation of $J(\theta) = V^{\pi_\theta}(s_0)$ by recursively solving the Bellman equation, as doing it during training becomes unnecessarily expensive. During the latter episodes, to study whether (15) holds true with $\varepsilon' = 0$ along the trajectories followed by the stochastic methods considered, we estimate the objective function $J(\theta)$ by sampling m trajectories $\{s_t^i, a_t^i\}_{t \geq 0}$, for $1 \leq i \leq m$, followed by the agent under a given fixed policy π_θ . We then define

$$\bar{J}(\theta) := \frac{1}{m} \sum_{j=1}^m \sum_{t=0}^{\infty} \gamma^t R(s_t^j, a_t^j) \quad (41)$$

as an unbiased estimate of $J(\theta)$. By adding an entropy bonus to the reward, we define similarly an unbiased estimate of the entropy regularized objective $J^\lambda(\theta)$ defined in (20). In particular, as explained previously, we truncate the infinite sum up to a finite horizon, but this is coherent with the implementation, where we impose a maximum number of steps H per episode. Moreover, the truncation bias can be made arbitrarily small by picking a not too large discount factor γ , which implies that the terms of the form $\gamma^t R(s_t, a_t)$ for $t \geq H + 1$ are basically negligible. A similar discussion holds for the estimate of $\nabla J(\theta)$, where we use the same estimate defined at (38) with a large enough $m = 100$ as to guarantee a good approximation of the expected value under the current policy. The optimum $J(\theta^*)$ is computed assuming a greedy policy for the agent, namely we assume

$$\pi(a|s) = \delta_{a,a^*}, \quad (42)$$

where a^* is the optimal action to be taken in state s . Upon these definitions, we define the estimated gradient dominance constant $\hat{\tau}_J$ as the minimum constant $C_J \in \mathbb{R}$ for which $J(\theta^*) - \bar{J}(\theta) \leq C_J \|\hat{\nabla} J(\theta)\|$ when we consider the standard objective (5), and $J^\lambda(\theta^*) - \bar{J}^\lambda(\theta) \leq C_J \|\hat{\nabla} \bar{J}^\lambda(\theta)\|^2$ when we consider the entropy regularized objective (20). This choice is coherent with the discussion we made in Section 2, as we know that in the deterministic case a weak gradient dominance holds along the trajectories of policy gradient, while a strong gradient dominance holds along the trajectories of entropy regularized policy

gradient, given an access to the full gradient. Finally, we study during training the boundedness from below of the action probabilities $\min_s \pi_\theta(a^*(s)|s)$ and $\min_{s,a} \pi_\theta(a|s)$ stated in inequalities (2.2) and (2.3), which we know are nicely bounded away from zero in the deterministic case.

As training details, during all the experiments, a random initialization of the policy is considered. In particular, we always start without any loss of generality from $\theta = 0$, for which the initial policy becomes $\pi_\theta(a|s) = \frac{1}{|A|}$. We then compare the performance of SCRNN against stochastic policy gradient methods by looking at how quickly the policy becomes deterministic and how well the algorithms escape saddle points and sub-optimal policies. In all the environments, an additional exploration bonuses was attempted, but has not provided the expected results: re-scaling the soft-max by a temperature parameter β . The scaling parameter β requires little extra effort in the computation of the gradient and of the Hessian: for a policy

$$\pi_\theta^\beta(a_i|s) = \frac{\beta \exp\{\theta_s^T \mathbb{1}_{a_i}\}}{\sum_j \exp\{\beta \theta_s^T \mathbb{1}_{a_j}\}}, \quad (43)$$

and the corresponding objective $J_\beta(\theta)$, the gradients and the Hessian (and also their unbiased estimates (38) and (40) are just pre-multiplied by β). For β small enough, exploration is vastly encouraged. However, in this particular setting the learned policy was basically independent of β for all optimization methods. For SCRNN, the addition of an entropy regularization term was also considered but did not provide better results.

As an implementation detail, note that under a soft-max policy, computing the Hessian once the gradient is computed does not require much additional computational cost. To see this, recall the estimates for $\nabla J(\theta)$ and $\nabla^2 J(\theta)$ introduced in (8) and (39). Hence, the running time of the trainings for the tested algorithms is comparable.

Hyperparameter setting

For all the environments shown in what follows, the following set of hyperparameters was used for the different algorithms:

- Number of episodes equal to 10000
- Discount factor $\gamma = 0.8$
- For all the algorithms besides the two stages policy gradient, batch sizes $n_1 = 1$ and $n_2 = 1$. For the two stages algorithm, initial batch size equal to 16.
- Learning rate $\alpha = 0.01$ for first order methods, with step decreasing schedule
- Learning rate $\alpha = 1e - 4$ for second order methods, with step decreasing schedule
- For entropy regularized objective, the temperature λ was 1

Cliff Walking

We start by considering the cliff walking RL task Sutton and Barto [2018]. The agent’s aim is to reach a goal state from a start state, avoiding a region of cells called “cliff”. The rest of the rectangular state space is free for the agent to move, but to encourage the finding of the shortest path to the goal a negative reward of -0.1 is given in all transitions. The penalty for falling into the cliff is a reward of -100, while reaching the goal on the other side of the cliff gives a reward of +100.

Figure 1 shows that initially, for all algorithms the episodes have a very short length, because the agent is unaware of the best action selection and falls into the cliff most of the episodes. The first learning step is actually the avoiding of the cliff, so the agent keeps exploring the state space, still unaware of the highest rewarded state at the extremity of the cliff region. Then, for the SCRNN and the two stages entropy regularized policy gradient methods, there is a moment of sudden discovery which makes the agent capable of reaching the goal. From this moment, in the remaining episodes the agent mostly tries to shorten the path to the goal by staying as close as possible to the cliff. This happens way sooner during training for SCRNN, which is also able to find the shortest path to the goal, with respect to first order methods. On the other hand, standard policy gradient methods are in general settling for just avoiding the cliff, which corresponds to a suboptimal policy even if $J(\theta^*) - J(\theta)$ may be contained.

This suggests that while the soft-max policy suffers from the fact that the Fisher information is not strictly positive definite around the optimum $J(\theta^*)$, the preconditioning induced by the full second order information boosts training. Indeed, as pointed out in Masiha et al. [2022], Figure 1 shows that the agent trained with stochastic cubic regularized Newton finds the goal way earlier during the training episodes, hence with a lower sample complexity. Finally, it is worth noticing that SCR N is able to find the optimal path, i.e. to stay always close to the cliff before moving down, even if probably during the first updates it is driven far away from the cliff after having discovered the very negative reward.

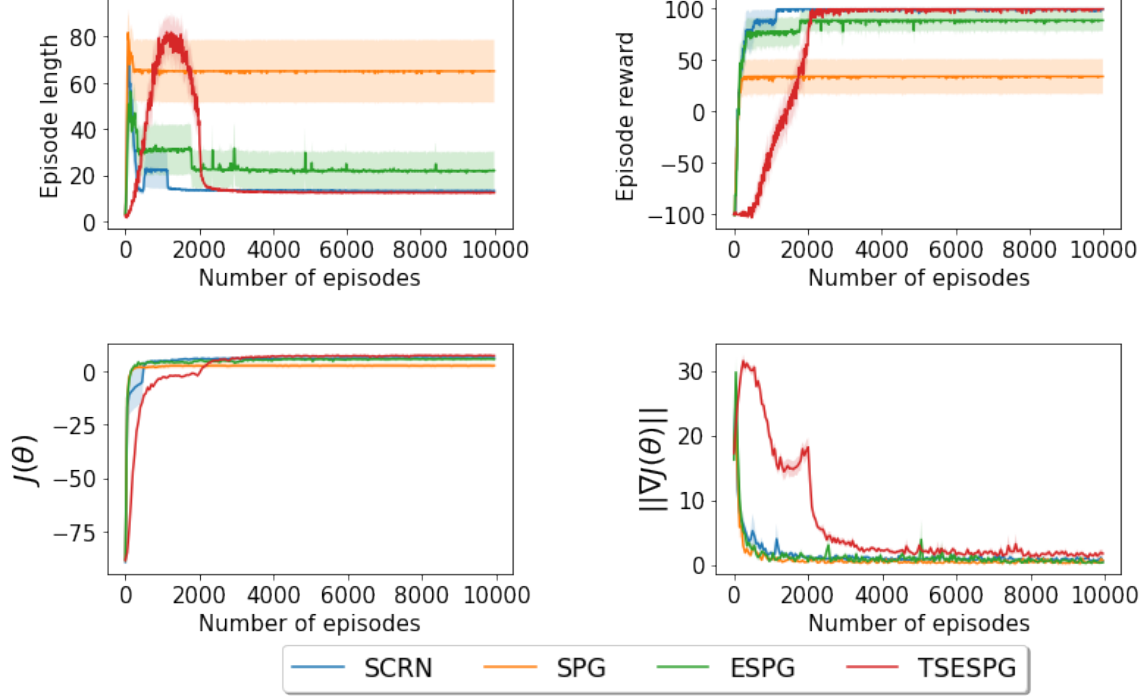


Figure 1: Training results for different RL algorithms in cliff walking optimization. Above: Episode lengths (left) and episode rewards (right) during training. Below: objective estimates (left) and gradient estimates (right) during training. Results are averaged over 10 full training runs. Shaded regions show the standard deviation.

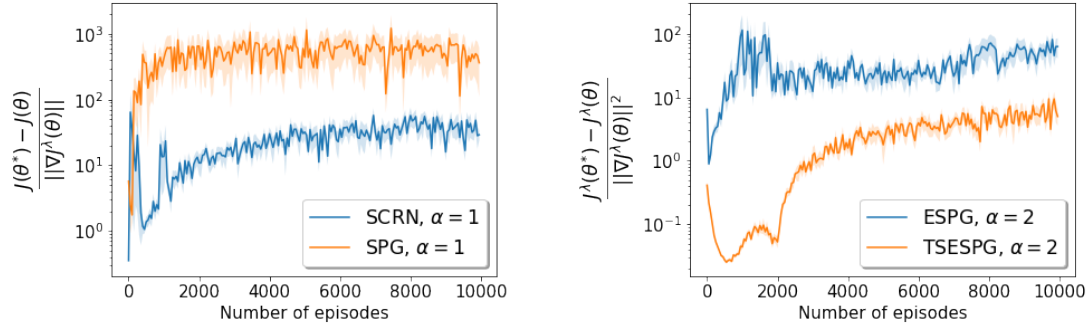


Figure 2: Estimation of τ_J along the trajectories of different RL algorithms for cliff walking optimization. Results are averaged over 10 full training runs. Shaded regions show the standard deviation.

Figure 2 shows that the empirical constant τ_J computed along the trajectories followed on average by the algorithms remains nicely bounded from below, for both the objective (5) with $\alpha = 1$ and the entropy regularized objective (20) with $\alpha = 2$. Furthermore, it is not diverging when the objective reaches an almost deterministic policy. Instead, it reaches an asymptotic value and does not show oscillations when evolving during training. The limited amount of oscillations which are observed can indeed be motivated with the noise in the estimation of $J(\theta)$ and $\nabla J(\theta)$. We can surely conclude that a gradient dominance property holds true, and the constant reaches an asymptotic value as soon as the policy converges, independently of the optimality of the latter. Figure 3 shows that the action probabilities remain equally bounded away from zero and they reach an asymptotic value as well without becoming too

small. Figure 3 (left) shows how the quantity $\min_s \pi_\theta(a^*(s)|s)$ is evolving over training when optimizing the unregularized objective (5), while Figure 3 (right) shows the decay of the quantity $\min_{s,a} \pi_\theta(a|s)$ for the entropy regularized objective (20).

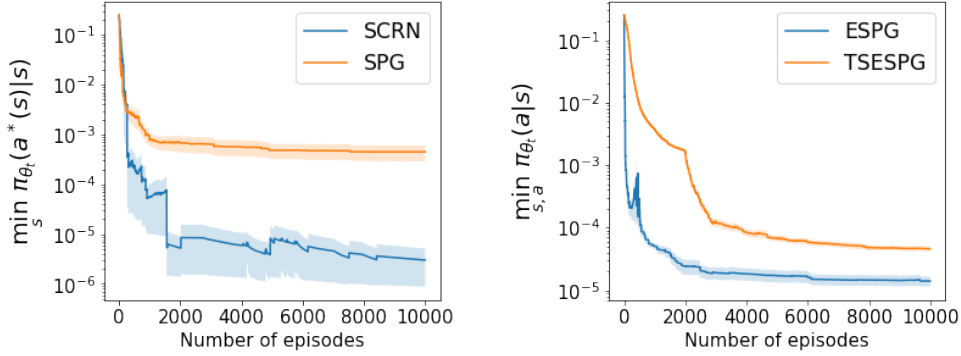


Figure 3: Action probabilities involved in the PL inequalities (19) for the objective (5) and (21) for the objective (20) in the cliff walking optimization. Results are averaged over 10 full training runs. Shaded regions show the standard deviation.

Hole discovering

We further considered the hole discovering as a RL task Frémaux et al. [2013]. This is often considered in a bigger class of reinforcement learning tasks as a navigation task, where the goal is hidden behind an obstacle and the trained agent may never become aware of the possibility of getting a very high reward by remaining in the neighbouring of the obstacle. The highest rewarded state (reward +100) is perhaps in the middle of a U -shaped wall. If the agent hits against the wall a reward -100 is given. Again, to encourage a greedy policy towards the goal, for every other action the reward is -0.1 .

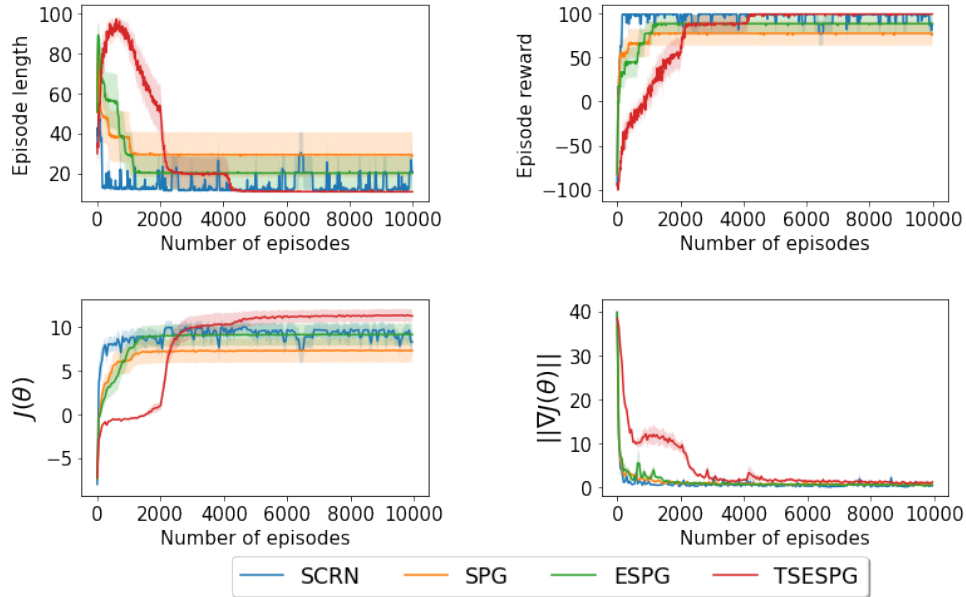


Figure 4: Training results for different RL algorithms in hole discovering optimization. Above: Episode lengths (left) and episode rewards (right) during training. Below: objective estimates (left) and gradient estimates (right) during training. Results are averaged over 10 full training runs. Shaded regions show the standard deviation.

Figure 4 shows that, contrarily to the cliff walking optimization shown in Figure 1, for all the algorithms the agent is able to capture immediately the presence of a wall in the middle of the state space, thus performing better on average. However, once again, SCRN finds the shortest path in the smaller number of episodes, while for first order methods the agent cannot find a short path to the goal. Figure 4 shows how SCRN exhibits some instability and oscillatory behaviour close to the optimum, which is

not the case instead for first order methods and in particular for the two stages entropy regularized SPG algorithm, which shows more stability after reaching a greedy policy which leads every time to the goal state.

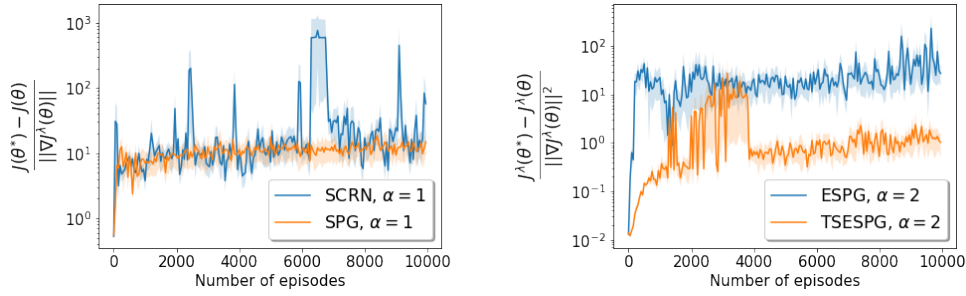


Figure 5: Estimation of τ_J along the trajectories of different RL algorithms for hole discovering optimization. Results are averaged over 10 full training runs. Shaded regions show the standard deviation.

Still, Figure 5 shows that the constant τ_J remains bounded from above and reaches an asymptotic value, without exploding, when the algorithms reach a greedy policy (optimal or not). Oscillations observed for SCRn and the two stages algorithm are due to the fact that for the former, the algorithm is inherently unstable close to the greedy policy, and for the latter is motivated by the sharp transition happening when switching from a large batch size to a small batch size, which makes it difficult to get accurate estimates for $J(\theta)$ and $\nabla J(\theta)$. Note that a cutoff effect in this transition is also observed in the quantity of interest showed in Figure 6 (right). However, Figure 6 shows again that the action probabilities remain nicely bounded away from zero. Most importantly, they reach an asymptotic value in the number of episodes which was used for training and they do not become arbitrarily small.

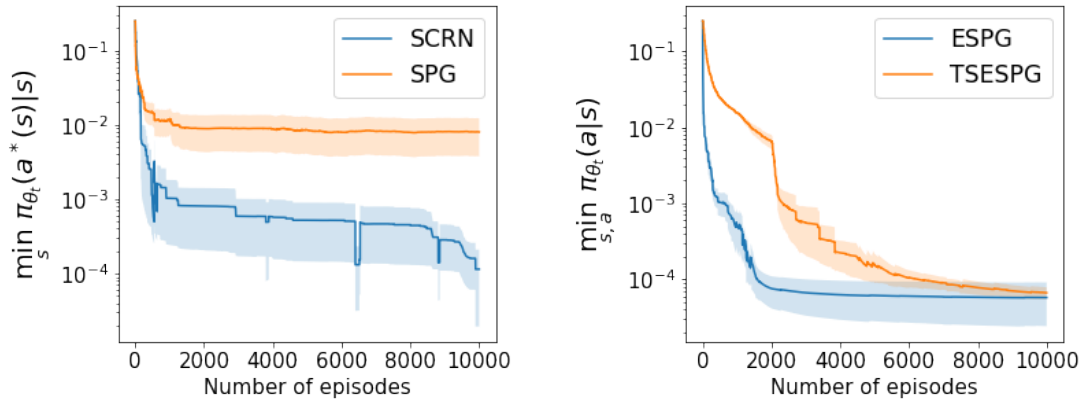


Figure 6: Action probabilities involved in the PL inequalities (19) for the objective (5) and (21) for the objective (20) in the hole discovering optimization. Results are averaged over 10 full training runs. Shaded regions show the standard deviation.

Random shape maze

We further considered the random mazes as a RL task Zuo [2018]. The goal in this case is to find the goal (reward +100) at the end of a randomly placed sequence of obstacles. If the agent hits the obstacles a reward -100 is given, for every other action the reward is -0.1. Figure 7 shows that SCRn is the only algorithm which is able to consistently find a path to the goal, while other algorithms are just settling for avoiding the randomly placed obstacles rather than finding the highest rewarded state on the other hand of the maze. The two stages algorithm is again the only first order method which is able to catch up or at least come close to the performance of SCRn. However, we notice again that the latter exhibits an unstable behaviour close to the optimum. Figure 8 shows nevertheless that the constant τ_J remains bounded from above during training and even after policy convergence, and Figure 9 shows that action probabilities are bounded from below coherently. In particular, while we know that Fisher non-degeneracy is violated as the softmax policy converges to a deterministic one, a gradient dominance property seems to hold without any diverging behaviour. Hence, empirically (15) holds even if we get arbitrarily close to the maximum of $J(\theta)$.

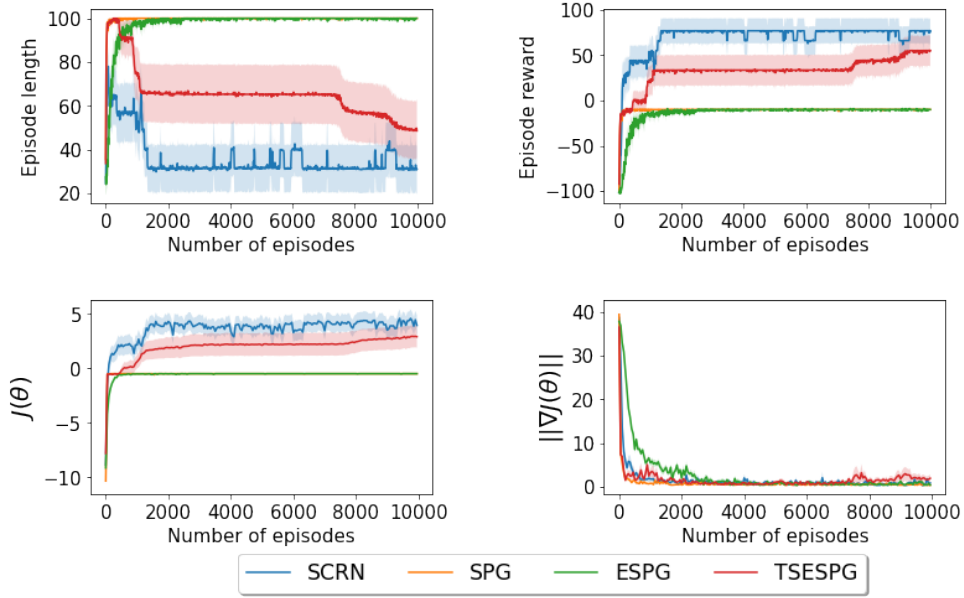


Figure 7: Training results for different RL algorithms in random shape maze optimization. Above: Episode lengths (left) and episode rewards (right) during training. Below: objective estimates (left) and gradient estimates (right) during training. Results are averaged over 10 full training runs. Shaded regions show the standard deviation.

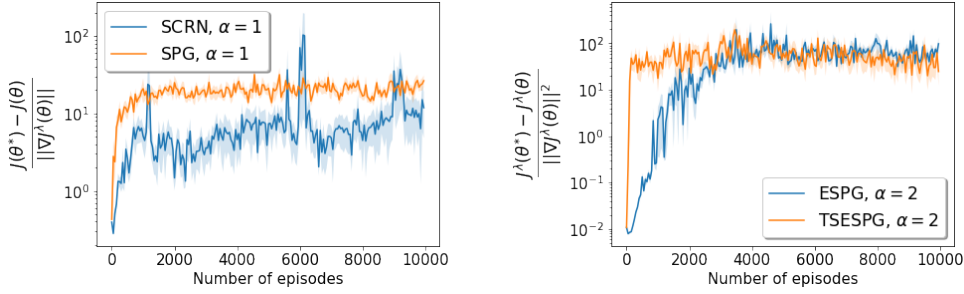


Figure 8: Estimation of τ_J along the trajectories of different RL algorithms for random mazes optimization. Results are averaged over 10 full training runs. Shaded regions show the standard deviation.

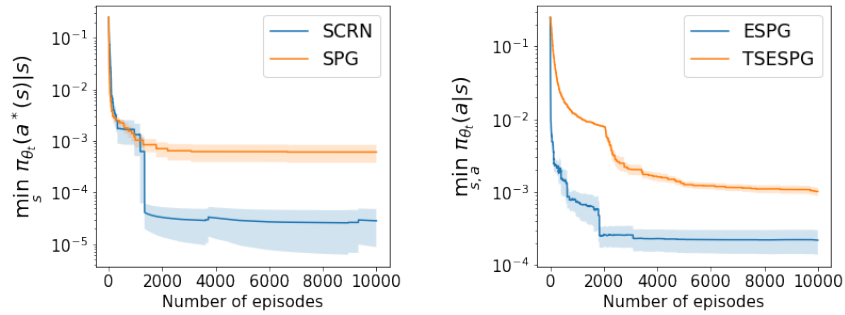


Figure 9: Action probabilities involved in the PL inequalities (19) for the objective (5) and (21) for the objective (20) in random shape maze optimization. Results are averaged over 10 full training runs. Shaded regions show the standard deviation.

U-Maze

We further considered the U shaped mazes as a RL task Zuo [2018]. The goal in this case is to find the goal (reward +100) on the other side of a U-shaped maze, where the obstacle is represented by the area which is not touched by the U. If the agent hits such a obstacle a reward -100 is given, for every other

action the reward is -0.1 to encourage the agent to find the optimal path as well. Figure 10 shows that stochastic policy gradient with and without entropy regularization are both failing to find the highest rewarded state while SCRNN finds a good path (but not optimal) to the goal, and a similar discussion holds true for the two stages entropy regularized SPG.

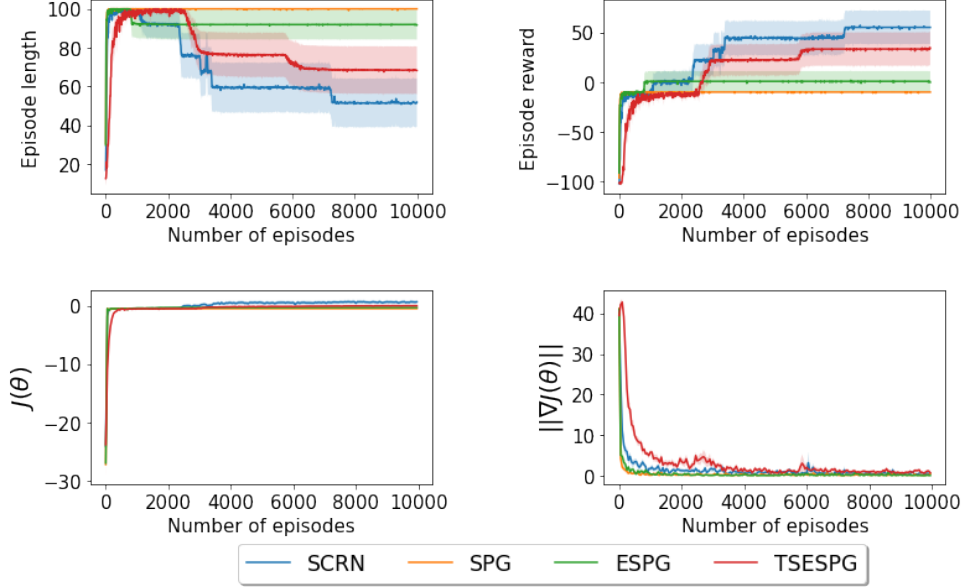


Figure 10: Training results for different RL algorithms in U-shaped maze optimization. Above: Episode lengths (left) and episode rewards (right) during training. Below: objective estimates (left) and gradient estimates (right) during training. Results are averaged over 10 full training runs. Shaded regions show the standard deviation.

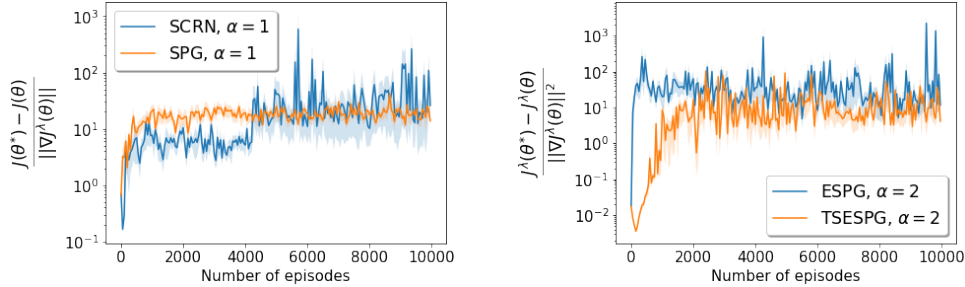


Figure 11: Estimation of τ_J along the trajectories of different RL algorithms for U-shaped maze optimization. Results are averaged over 10 full training runs. Shaded regions show the standard deviation.

In particular, while we know that Fisher non-degeneracy is violated as the softmax policy converges to a deterministic one, a gradient dominance property seems to hold without any diverging behaviour, as we can see in Figure 11. Hence, we observe once again that empirically (15) holds even if we get arbitrarily close to the maximum of $J(\theta)$. We observe a higher noise with respect to the previous environment in the estimation of the PL constant probably because the policy still needs to converge and hence the number of episodes used for testing may not be large enough. Coherently, action probabilities remain bounded as it can be observed in Figure 12, but the quantities $\min_s \pi_\theta(a^*(s)|s)$ for the standard objective $J(\theta)$ and $\min_{s,a} \pi_\theta(a|s)$ are exhibiting higher oscillations. Still, in both cases, no diverging behaviour for the PL constant or degenerate behaviour for the action probabilities seems to be observed.

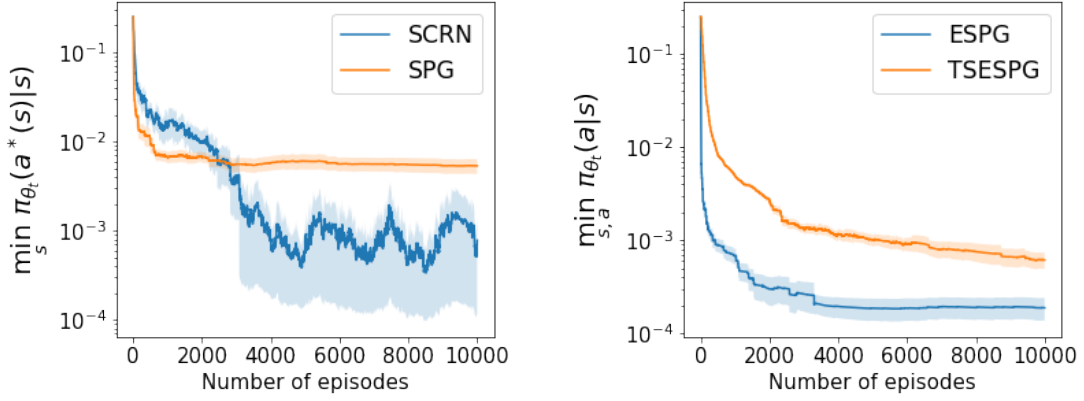


Figure 12: Action probabilities involved in the PL inequalities (19) for the objective (5) and (21) for the objective (20) in U maze walking optimization. Results are averaged over 10 full training runs. Shaded regions show the standard deviation.

4 A comparison between CRN and NPG

Natural policy gradient has been the starting point for the theoretical development of many successful algorithms in policy gradient’s based reinforcement learning such as TRPO (Trust Region Policy Optimization) Schulman et al. [2015] and PPO (Proximal Policy Optimization) Schulman et al. [2017]. Natural policy gradient belongs to a general class of update rule of the form

$$\theta_t \leftarrow \theta_t + \alpha_t (M(\theta_t))^{-1} \nabla J(\theta_t), \quad (44)$$

where $M(\theta)$ is a preconditioner for the gradient which needs to be chosen with a trade-off between computational cost and performance in revealing information about the landscape of the objective function. The standard policy gradient rule takes naively $M(\theta_t) = I$ as the identity matrix for all $t \geq 0$, and hence no preconditioning is applied. Considering the second order information given by $\nabla^2 J(\theta)$, one recovers the Newton’s Method which we briefly introduced in Section 2. Natural Policy Gradient chooses the Fisher information matrix defined in (13) as a preconditioner for the update. It is clear from the form of (44) that for such an algorithm to work well, Fisher non-degeneracy must hold true such that the inverse of $M(\theta_t)$ is always well defined. However, as we extensively clarified throughout the elaborate, for a soft-max policy, the Fisher information matrix becomes degenerate when the policy gets arbitrarily close to a greedy one, whether the latter is optimal or not. This makes it difficult to study the performance of such an algorithm under a soft-max policy. In this section, we aim at a first comparison between the second order information for preconditioning the gradient update when using the whole Hessian matrix and the Fisher information (13), or estimates of the latter quantities. We then reproduce the results provided in the pioneering work about NPG methods Kakade [2001] and we discuss further improvements of the current analysis.

Natural Policy Gradient

Let us start by doing a brief recap on the motivation of natural policy gradients. Traditional policy gradient methods are fundamentally flawed. In such methods, the gradient $\nabla J(\theta)$ only gives the direction of the weight update. It does not tell how far to step into that direction, and because derivatives are defined on an infinitesimal interval, the gradient is valid only locally and may be completely different at another part of the function (5). Hence, the update (9) does not take into account the possibility of overshooting, and in a reinforcement learning setting this may result in a poorly updated policy from which the agent is not able to recover, because future sample batches obtained by interacting with the environment may not provide much meaningful information. A naive solution to prevent overshooting is to pick a small enough learning rate α , but this slows down convergence. Another possibility would be to cap the update changes to some ε tolerance, namely to define

$$\Delta\theta = \operatorname{argmax}_{\|\Delta\theta\| \leq \varepsilon} J(\theta + \Delta\theta), \quad (45)$$

but capping the parameter space does not effectively cap the statistical (and nonlinear) manifold of random policies on which the algorithm is operating. To prevent overshooting during an update, one should

impose a constraint on the policy change instead of the parameters. This is what **natural policy gradient** methods perform. Capping the difference between the old policy and the updated one is a much more powerful way to control the optimization problem on the statistical manifold of policies. This is precisely the main advantage of second order information: by preconditioning the gradient with curvature information capturing the sensitivity of the manifold to marginal parameter changes, we ensure that the old and the new policy are not too far from each other, and this on one hand prevents overshooting and on the other hand makes learning faster in flat, suboptimal plateaus where first order methods as policy gradients may remain stuck.

To define the distance between two policies (i.e. probability distributions defined on the state-action space of a RL environment), the **Kullback-Leibner divergence** is often considered. The latter is defined as

$$\mathcal{D}_{KL}(\pi_\theta \| \pi_{\theta+\Delta\theta}) = \mathbb{E}_{\pi_\theta} \left[\log \left(\frac{\pi_\theta}{\pi_{\theta+\Delta\theta}} \right) \right]. \quad (46)$$

The locally defined Hessian of the KL divergence is actually equivalent to the Fisher information matrix (13). Indeed, the latter is often defined as the Riemannian metric describing the curvature of a statistical manifold, and it may be viewed as a correction to the distance that accounts for the curvature. If the Fisher matrix is an identity matrix, the distance over the manifold is simply the Euclidean distance. In that case, traditional and natural policy gradients are equivalent.

If we now modify (45) with a constraint on the policy change rather than the parameter change, we may define the local optimization problem as

$$\Delta\theta = \operatorname{argmax}_{\mathcal{D}_{KL}(\pi_\theta \| \pi_{\theta+\Delta\theta}) \leq \varepsilon} J(\theta + \Delta\theta). \quad (47)$$

However, computing the KL divergence for large RL environments is often not feasible. This is where the Fisher information matrix becomes a more appealing solution. We start by solving the unconstrained optimization problem

$$\Delta\theta^* = \operatorname{argmax}_{\Delta\theta} J(\theta + \Delta\theta) - \lambda(\mathcal{D}_{KL}(\pi_\theta \| \pi_{\theta+\Delta\theta}) - \varepsilon), \quad (48)$$

as a surrogate (and easier to handle) objective with respect to (47). By Taylor expansion in a neighbourhood of the current policy θ we get

$$\Delta\theta^* \sim \operatorname{argmax}_{\Delta\theta} J(\theta) + \nabla J(\theta)^T \Delta\theta - \frac{1}{2} \lambda (\Delta\theta^T \nabla^2 \mathcal{D}_{KL}(\pi_\theta \| \pi_{\theta+\Delta\theta}) \Delta\theta) + \lambda \varepsilon.$$

Because we are considering the local KL divergence, we may substitute its Hessian by the Fisher information matrix. Note that the Fisher information matrix can be represented as the outer product of policy gradients, and while the expression is locally equivalent to the Hessian matrix of the KL divergence, it is way computationally cheaper to compute. Dropping all the terms which are not depending on $\Delta\theta$ explicitly in the above now gives

$$\Delta\theta^* \sim \nabla J(\theta)^T \Delta\theta - \frac{1}{2} \lambda \Delta\theta^T F_\rho(\theta) \Delta\theta.$$

Now, we aim at maximizing the parameter change; we can then use the Lagrangian multiplier λ to impose a small change in the policy. If Fisher non-degeneracy holds, it is straightforward to verify by differentiating in $\Delta\theta$ that the objective $\nabla J(\theta)^T \Delta\theta - \frac{1}{2} \lambda \Delta\theta^T F_\rho(\theta) \Delta\theta$ is strongly concave with a unique maximizer given by $\Delta\theta^* = -\frac{1}{\lambda} F_\rho(\theta)^{-1} \nabla J(\theta)$. λ is then used to impose a dynamical learning rate satisfying the KL constraint. Finally, this gives the update

$$\theta \leftarrow \theta + \alpha F_\rho(\theta)^{-1} \nabla J(\theta), \quad (49)$$

which is the standard gradient preconditioned with the inverse Fisher matrix, accounting for the curvature of the Riemannian space.

Fisher information matrix and Hessian matrix

Let us begin this comparison by recalling the expression for the Fisher information matrix

$$F_\rho(\theta) = \mathbb{E}_{(s,a) \sim \nu_\rho^{\pi_\theta}(s,a)} \left[\nabla \log \pi_\theta(a|s) \nabla \log \pi_\theta(a|s)^T \right], \quad (50)$$

for the quantities defined previously. From Theorem 3 in Furnston et al. [2016], we know that we can write the Hessian of the objective $J(\theta)$ as

$$\begin{aligned} \nabla^2 J(\theta) = \mathbb{E}_{\tau \sim p(\cdot|\pi_\theta)} & \left[\sum_{h=0}^{H-1} \left(\sum_{t=h}^{H-1} \gamma^t R(s_t, a_t) \right) \nabla \log \pi_\theta(a_h|s_h) \nabla \log \pi_\theta(a_h|s_h)^T \right] \\ & + \left[\sum_{h=0}^{H-1} \left(\sum_{t=h}^{H-1} \gamma^t R(s_t, a_t) \right) \nabla^2 \log \pi_\theta(a_h|s_h) \right]. \end{aligned} \quad (51)$$

In the first term of (51), we recognize an expression similar to that of the Fisher information matrix, but weighted with the state-action value function $\sum_{t=h}^{H-1} \gamma^t R(s_t, a_t)$. Indeed, while in the Fisher information matrix the expectation is taken with respect to the geometrically weighted summation of state-action occupancy marginals of the trajectory distribution, in the first term of the Hessian there is an additional weighting from the state-action value function. Hence, the latter expression incorporates information about the reward structure of the objective function, and hence about the curvature of the objective function. The second term in (51) is the Hessian of the log-likelihood of the policy, i.e. the gradient of the score function scaled by its value in expectation with respect to the discounted state-occupancy measure over the state-space \mathcal{S} . Interestingly, it has been shown empirically in the hallway problem in Furnston et al. [2016] that when getting closer to the optimal policy, the first part of the Hessian in (51) tends to dominate the second term in spectral norm. This is used in Furnston et al. [2016] to design approximate Newton solvers for policy optimization, where only the (much cheaper to generate and store) first term of the Hessian is computed. Another possibility, whose development we leave to future work, would be to design an algorithm which passes only the dominant term in (51) to the cubic subsolver in Algorithm 1.

The case of a simple MDP

As we anticipated in the introduction, the natural policy gradient was introduced in the pioneering work of Kakade [2001]. In the latter work, the numerical experiments are carried out on a simple yet instructive MDP which we now briefly introduce and study. The considered environment depicted in Figure 13 has only two states, which we call i and j . The agent always starts from state i , and remaining in such a state by doing a self-loop gives a suboptimal reward of 1, whereas the same action performed in state j gives a reward of 2. The transitions $i \rightarrow j$ and vice-versa have no reward. Figure 13 (C) shows the average episode reward as a function of the training episodes, where every episode lasts a maximum of 20 steps. We observe that using only the first order information the agent is stuck for a very long time and keeps on getting a 1 reward, basically becoming unaware of the possibility of transitioning to state j looking for a higher reward. This realization only occurs after $\approx 10^7$ episodes. On the other hand, the agent trained with the Natural Policy Gradient algorithm becomes aware of the higher reward in state j after a few episodes. Figure 13 (D) shows how θ_i and θ_j , both increasing with the probabilities of doing a self-loop in states i and j , respectively, are evolving over time one with respect to the other. In particular, we see that for NPG (dashed line) θ_j becomes much larger soon, while for SPG (solid line) θ_i becomes much larger before θ_j grows. As we already discussed, this may be something from which the agent cannot recover in more complex environments with a bigger number of state-action pairs.

In our work, similarly to what is done in Kakade [2001], the Fisher information matrix is estimated from trajectories followed by the agent. In particular, by sampling m trajectories $\{s_t^i, a_t^i\}_{t \geq 0}$, for $1 \leq i \leq m$, followed by the agent under a given fixed policy π_θ , we define then

$$\bar{F}_\rho(\theta) := \frac{1}{m} \sum_{j=1}^m \sum_{t=0}^{\infty} \nabla \log \pi_\theta(a_t^j|s_t^j) \log \pi_\theta(a_t^j|s_t^j)^T \quad (52)$$

as an unbiased estimate of the Fisher information matrix (13). Because the possible number of trajectories followed by the agent is rather small on such an easy environment and 20 steps episodes, we expect these estimates to be accurate enough with a relatively small batch size. This approach was also followed in the original paper. Figure 14 shows the episode lengths for NPG and SPG: we see that NPG finds the optimal state in a very small number of episodes, while SPG is stuck performing suboptimal actions and the policy diverges from the optimal one. This confirms the superior behaviour of the natural policy gradient on this very simple environment.

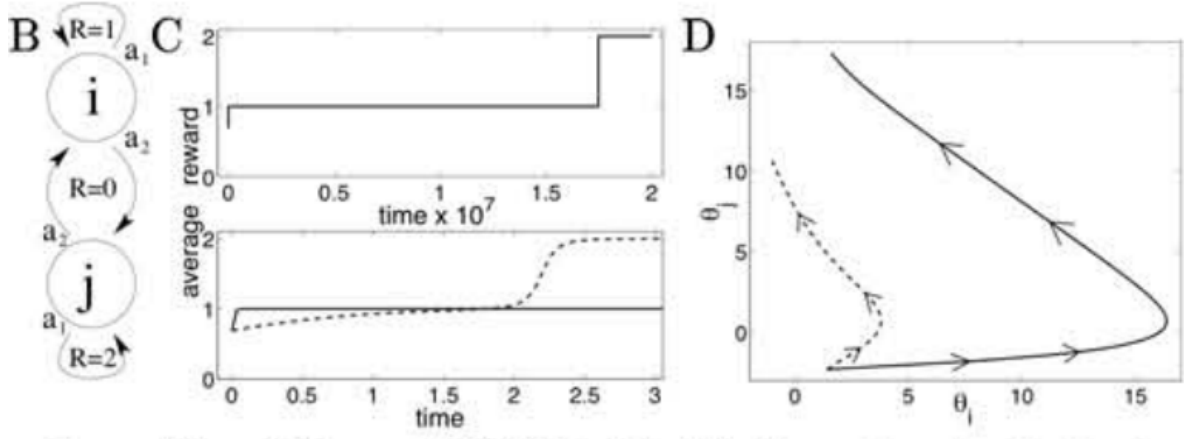


Figure 13: Simple MDP which gives insight on the superior behaviour of NPG over SPG. The agent always starts in state i and its goal is to reach state j . However, a self-loop from state i gives a reward of 1, which is a suboptimal trap in which the agent must not fall. Indeed, a self-loop from state j gives a higher reward equal to 2. Figure is taken from the original work on NPG Kakade [2001].

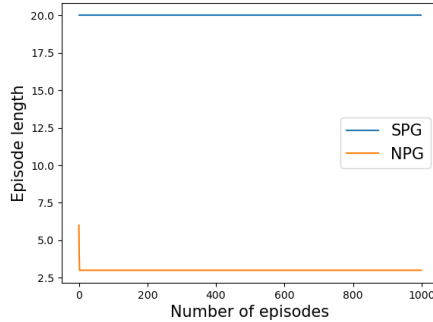


Figure 14: Results for NPG and SPG on the simple MDP shown in Figure 13. We see that SPG is stuck in state i throughout all the episodes, while NPG takes only a very small amount of episodes to go straight to the highest rewarded state.

5 Conclusions

In this work, we started from the general purpose of checking empirically the validity of a gradient dominance property along the trajectories followed by stochastic first order and second order methods in a reinforcement learning setting. The validity of such an inequality, which implies that any stationary point is a global optimum, was already established under a soft-max policy for action selection in Mei et al. [2020] for the (deterministic) policy gradient and entropy regularized policy gradient methods. Similar results for the stochastic setting are still very limited in the literature, a first successful attempt at a similar analysis being the two stages algorithm proposed in Ding et al. [2021], which we also took into account in the work. The empirical results show that the PL constant, computed with estimates both for the objective function and its gradient, remains nicely bounded from below on a variety of reinforcement learning environments and for all the stochastic first order and second order considered methods. Similarly, action probabilities which are involved in the non-uniform PL inequalities established for the expected return remain bounded from below. We observed, coherently with the condition of Fisher non-degeneracy not being necessarily satisfied, that the constant becomes larger as we get arbitrarily close to the optimum or to a greedy suboptimal policy because action probabilities may become very small. We may consider solving directly the Bellman equation to get the expected return and its gradient as a future development of the current work, as well as testing the PL inequality on bigger and more various environments. In the final part of the elaborate, we worked in the direction of comparing the second order information enclosed in the Fisher information Kakade [2001] and in the Hessian matrix of the expected return. We showed that the latter contains more information about the curvature of the function by taking into account the Q -values of the state-action pairs. On the other hand, the Hessian of the log-

likelihood in the Hessian expression has been shown to be negligible when getting close to the optimum in previous work Furnston et al. [2016]. This can be of interest for the future design of practical algorithms neglecting the second term in the Hessian, as the first term is much cheaper to compute and store.

References

- A. Agarwal, S. M. Kakade, J. D. Lee, and G. Mahajan. Optimality and approximation with policy gradient methods in markov decision processes. In *Conference on Learning Theory*, pages 64–66. PMLR, 2020.
- A. Agarwal, S. M. Kakade, J. D. Lee, and G. Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *J. Mach. Learn. Res.*, 22(98):1–76, 2021.
- J. C. A. Barata and M. S. Hussein. The Moore–Penrose pseudoinverse: A tutorial review of the theory. *Brazilian Journal of Physics*, 42(1):146–165, 2012.
- R. Battiti. First-and second-order methods for learning: between steepest descent and newton’s method. *Neural computation*, 4(2):141–166, 1992.
- Y. Ding, J. Zhang, and J. Lavaei. Beyond exact gradients: Convergence of stochastic soft-max policy gradient methods with entropy regularization. *arXiv preprint arXiv:2110.10117*, 2021.
- N. Frémaux, H. Sprekeler, and W. Gerstner. Reinforcement learning using a continuous time actor-critic framework with spiking neurons. *PLoS computational biology*, 9(4):e1003024, 2013.
- T. Furnston, G. Lever, and D. Barber. Approximate newton methods for policy search in markov decision processes. *Journal of Machine Learning Research*, 17, 2016.
- Y. Grandvalet and Y. Bengio. Entropy regularization., 2006.
- S. Kakade and J. Langford. Approximately optimal approximate reinforcement learning. In *In Proc. 19th International Conference on Machine Learning*. Citeseer, 2002.
- S. M. Kakade. A natural policy gradient. *Advances in neural information processing systems*, 14, 2001.
- Y. Liu, K. Zhang, T. Basar, and W. Yin. An improved analysis of (variance-reduced) policy gradient and natural policy gradient methods. *Advances in Neural Information Processing Systems*, 33:7624–7636, 2020.
- S. Masiha, S. Salehkaleybar, N. He, N. Kiyavash, and P. Thiran. Stochastic second-order methods provably beat sgd for gradient-dominated functions. *arXiv preprint arXiv:2205.12856*, 2022.
- J. Mei, C. Xiao, C. Szepesvari, and D. Schuurmans. On the global convergence rates of softmax policy gradient methods. In *International Conference on Machine Learning*, pages 6820–6829. PMLR, 2020.
- Y. Nesterov and B. T. Polyak. Cubic regularization of newton method and its global performance. *Mathematical Programming*, 108(1):177–205, 2006.
- J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR, 2015.
- J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Z. Shen, A. Ribeiro, H. Hassani, H. Qian, and C. Mi. Hessian aided policy gradient. In *International conference on machine learning*, pages 5729–5738. PMLR, 2019.
- R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- N. Tripuraneni, M. Stern, C. Jin, J. Regier, and M. I. Jordan. Stochastic cubic regularization for fast nonconvex optimization. *Advances in neural information processing systems*, 31, 2018.
- R. Yuan, R. M. Gower, and A. Lazaric. A general sample complexity analysis of vanilla policy gradient. In *International Conference on Artificial Intelligence and Statistics*, pages 3332–3380. PMLR, 2022.
- X. Zuo. mazelab: A customizable framework to create maze and gridworld environments, 2018.