

Policy gradient in Lipschitz Markov Decision Processes

Matteo Pirotta · Marcello Restelli · Luca Bascetta

Received: 28 September 2014 / Accepted: 6 February 2015 / Published online: 3 March 2015
© The Author(s) 2015

Abstract This paper is about the exploitation of Lipschitz continuity properties for Markov Decision Processes to safely speed up policy-gradient algorithms. Starting from assumptions about the Lipschitz continuity of the state-transition model, the reward function, and the policies considered in the learning process, we show that both the expected return of a policy and its gradient are Lipschitz continuous w.r.t. policy parameters. By leveraging such properties, we define policy-parameter updates that guarantee a performance improvement at each iteration. The proposed methods are empirically evaluated and compared to other related approaches using different configurations of three popular control scenarios: the linear quadratic regulator, the mass-spring-damper system and the ship-steering control.

Keywords Reinforcement learning · Markov Decision Process · Lipschitz continuity · Policy gradient algorithm

1 Introduction

In the last years, policy-gradient methods have emerged among the most effective Reinforcement-Learning (RL) techniques for complex real-world control problems with continuous, high-dimensional, and partially-observable properties, such as robotic control systems (Peters and Schaal 2006). Given a parameterized policy space, usually designed to incorporate domain knowledge, policy-gradient algorithms update policy parameters along

Editors: Concha Bielza, João Gama, Alípio M. Jorge, and Indrè Žliobaitė.

M. Pirotta (✉) · M. Restelli · L. Bascetta
Department of Electronics, Information, and Bioengineering,
Politecnico di Milano, Piazza Leonardo Da Vinci, 32, I-20133 Milan, Italy
e-mail: matteo.pirotta@polimi.it

M. Restelli
e-mail: marcello.restelli@polimi.it

L. Bascetta
e-mail: luca.bascetta@polimi.it

an estimated ascent direction of the expected return. Under some mild assumptions on the *step size* used to update the parameters (Moré and Thuente 1994), policy-gradient methods are guaranteed to converge at least to a locally optimal solution.

The research in policy gradient has mainly focused on defining convenient ascent directions and low-variance, model-free estimators of the policy gradient. The oldest policy-gradient approaches are finite-difference methods (Spall 1992), that estimate gradient direction by resolving a regression problem based on the performance evaluation of policies associated to different small perturbations of the current parametrization. Finite-difference methods have some advantages: they are easy to implement, do not need assumptions on the differentiability of the policy w.r.t. the policy parameters, and are efficient in deterministic settings. On the other hand, when used on real systems, the choice of parameter perturbations may be difficult and critical for system safeness. Furthermore, the presence of uncertainties may significantly slow down the convergence rate. Such drawbacks have been overcome by likelihood ratio methods (Williams 1992; Baxter and Bartlett 2001; Sutton et al. 1999), since they do not need to generate policy parameter variations and quickly converge even in highly stochastic systems. Several studies have addressed the problem to find minimum variance estimators by the computation of optimal baselines (Peters and Schaal 2008b). To further improve the efficiency of policy-gradient methods, natural-gradient approaches (where the steepest ascent is computed w.r.t. the Fisher information metric) have been considered (Kakade 2001; Peters and Schaal 2008a). Natural gradients still converge to locally optimal policies, are independent from the policy parametrization, need less data to attain good gradient estimates, and are less affected by plateaus. For recent and comprehensive surveys on policy search and policy gradient methods we refer the reader to Grondman et al. (2012) and Deisenroth et al. (2013).

Unfortunately, a good estimate of the policy gradient is not enough to guarantee effective learning. In fact, even when the exact policy gradient is known, the choice of the step size strongly influences the number of iterations needed to attain a (local) maximum or, even worse, can make convergence unfeasible (Wagner 2011). In general unconstrained programming, the value of the step size is determined through line-search algorithms (Moré and Thuente 1994), that require to evaluate the function to be optimized at points generated along the gradient direction by a sequence of candidate values for the step size. In the policy-gradient framework, being policy evaluations quite expensive, line search is impractical and step-size parameters are usually kept fixed or decreased over time according to some annealing schedule, requiring significant amounts of hand tuning. Convergence issues can be solved by making the step-size parameter decrease according to the Robbins-Monro conditions (Robbins and Monro 1951), but it usually turns out to show very slow convergence.

In spite of the strong impact of the step size over the performance of policy-gradient methods, so far little research has addressed such issue, with a few notable exceptions. Kober and Peters (2008) and Vlassis et al. (2009) studied policy-search methods based on expectation-maximization. Under some assumptions on the reward and policy models, expectation-maximization algorithms have properties similar to the ones of policy gradients, but without the need of specifying any step size. In Pirodda et al. (2013) we have directly addressed the problem of computing a step size that guarantees a policy improvement at any iteration. The idea is to use the data collected using the current policy to lower bound the expected return of any policy. The step size is then chosen to maximize such lower bound along the policy-gradient direction.

The main limitation of previous approaches is the looseness of the lower bounds to the expected return, that usually leads to conservative policy updates. In order to mitigate this drawback, we focus on Lipschitz-continuous MDPs, that represent a relevant subclass of

MDPs. In fact, many real-world problems are characterized by continuous state and action spaces (e.g., robotics, automatic control problems, natural resource management, etc.), where it is reasonable that when similar actions are executed in similar states their effects will be similar. That is what the Lipschitz assumptions want to capture. In this paper, we show that, under Lipschitz continuity assumptions on the Markov Decision Process (MDP) and the policy model (Sect. 2), the expected return of each policy and its policy-gradient components are Lipschitz w.r.t. policy parameters (Sects. 3 and 4). As shown by Armijo (1966), the Lipschitz continuity of the gradient can be used to select the value of the step size so as to guarantee a performance improvement at each iteration. In particular, the smaller are the Lipschitz constants, the larger are the step sizes and the expected improvements. We introduce how to compute the Lipschitz constant related to each component of the gradient and we show how such constants can be used to guarantee a performance improvement either by automatically identifying a proper step size along the gradient direction, or by defining new ascent directions with better guarantees (Sect. 5). Besides the theoretical contributions, we will also provide an empirical analysis to highlight advantages and limitations of the proposed approach (Sect. 6).

2 Preliminaries

In this section, we introduce notation and basic concepts about MDPs, Lipschitz MDPs, and policy gradients.

2.1 Markov Decision Process

A discrete-time continuous MDP is defined as a 6-tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \mu \rangle$, where \mathcal{S} is the continuous state space, \mathcal{A} is the continuous action space, \mathcal{P} is a Markovian transition model where $\mathcal{P}(s'|s, a)$ defines the transition density between state s and s' under action a , $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow [-R, R]$ is the reward function, such that $\mathcal{R}(s, a)$ is the expected immediate reward for the state-action pair (s, a) and R is the maximum absolute reward value, $\gamma \in [0, 1)$ is the discount factor for future rewards, and μ is the initial state distribution. We assume state and action spaces to be complete, separable metric (Polish) spaces $(\mathcal{S}, d_{\mathcal{S}})$ and $(\mathcal{A}, d_{\mathcal{A}})$, equipped with their σ -algebras $\sigma_{\mathcal{S}}, \sigma_{\mathcal{A}}$ of Borel sets, respectively. We assume—as done in Hinderer (2005)—that joint state-action space is endowed with the following taxicab norm: $d_{\mathcal{S}\mathcal{A}}((s, a), (\hat{s}, \hat{a})) = d_{\mathcal{S}}(s, \hat{s}) + d_{\mathcal{A}}(a, \hat{a})$. A stationary policy $\pi(\cdot|s)$ specifies for each state s the density function over the Borel action space $(\mathcal{A}, d_{\mathcal{A}}, \sigma_{\mathcal{A}})$.

We consider infinite-horizon problems where the future rewards are exponentially discounted with γ . For each state s , we define the utility of following a stationary policy π as:

$$V^\pi(s) = \mathbb{E}_{\substack{a_t \sim \pi \\ s_t \sim \mathcal{P}}} \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(s_t, a_t) \mid s_0 = s \right].$$

It is known that, under mild assumptions (Bertsekas and Shreve 1978), V^π solves the following recursive (Bellman) equation:

$$V^\pi(s) = \int_{\mathcal{A}} \left(\mathcal{R}(s, a) + \gamma \int_{\mathcal{S}} V^\pi(s') P(ds'|s, a) \right) \pi(da|s).$$

For model-free control purposes, the value function V is usually replaced by the action-value function Q , where action value $Q^\pi(s, a)$ is the expected return of taking action a in state s and following a policy π thereafter:

$$Q^\pi(s, a) = \mathcal{R}(s, a) + \gamma \iint_{\mathcal{S}\mathcal{A}} Q^\pi(s', a') \pi(da'|s') \mathcal{P}(ds'|s, a).$$

Policies can be ranked by their expected discounted reward starting from the state distribution μ :

$$\begin{aligned} J_\mu^\pi &= \int_{\mathcal{S}} V^\pi(s) \mu(ds) = \frac{1}{1-\gamma} \iint_{\mathcal{S}\mathcal{A}} \mathcal{R}(s, a) \pi(da|s) \delta_\mu^\pi(ds) \\ &= \frac{1}{1-\gamma} \mathbb{E}_{s \sim \delta_\mu^\pi(\cdot)} \mathbb{E}_{a \sim \pi(\cdot|s)} [\mathcal{R}(s, a)], \end{aligned}$$

where $\delta_\mu^\pi(s) = (1-\gamma) \sum_{t=0}^\infty \gamma^t Pr(s_t = s | \pi, \mu)$ is the γ -discounted future state distribution for a starting state distribution μ (Sutton et al. 1999). It is possible to rewrite previous equation in terms of the joint distribution $\zeta(\delta_\mu^\pi, \pi)$ between the future state distribution δ_μ^π and the stationary policy π , that can be written as a function of only π . Let $\zeta(\delta_\mu^\pi, \pi) = \zeta_\mu^\pi$ be a probability distribution over $\mathcal{S} \times \mathcal{A}$, such that:

$$J_\mu^\pi = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \zeta_\mu^\pi} [\mathcal{R}(s, a)].$$

Solving an MDP means finding a policy π^* that maximizes the expected long-term reward: $\pi^* \in \text{arg max}_{\pi \in \Pi} J_\mu^\pi$. For any MDP there exists at least one deterministic stationary optimal policy that simultaneously maximizes $V^\pi(s)$, $\forall s \in \mathcal{S}$ (Puterman 1994).

2.2 Lipschitz MDP

In this section, we introduce the basic concepts of Lipschitz continuity. Given two metric sets (X, d_X) and (Y, d_Y) , where d_X and d_Y denote the corresponding metric functions, a function $f : X \rightarrow Y$ is called L_f -Lipschitz continuous (L_f -LC) if

$$\forall (x_1, x_2) \in X^2, \quad d_Y(f(x_1), f(x_2)) \leq L_f d_X(x_1, x_2). \tag{1}$$

The smallest constant L_f for which (1) holds is called the Lipschitz constant of f . Define $\|f\|_L = \sup_{x_1 \neq x_2} \left\{ \frac{d_Y(f(x_1), f(x_2))}{d_X(x_1, x_2)} : x_1, x_2 \in X \right\}$ to be the Lipschitz semi-norm over the function space $\mathcal{F}(X, Y)$. Furthermore, we call f pointwise Lipschitz continuous¹ (PLC) in state x if there exists a constant $L_f(x)$ such that:

$$\forall x' \in X, \quad d_Y(f(x), f(x')) \leq L_f(x) d_X(x, x') \quad \text{where } \forall x \in X, L_f(x) \leq L_f.$$

For real-valued functions (e.g., the reward function), we will use the Euclidean distance as metric for the codomain. On the other hand, for the state-transition model and the policies we need to introduce a distance between probability distributions. Following Hinderer (2005) and Rachelson and Lagoudakis (2010), we will consider the Kantorovich or L^1 -Wasserstein metric on probability measures p and q :

$$\mathcal{K}(p, q) = \sup_f \left\{ \left| \int_X f d(p - q) \right| : \|f\|_L \leq 1 \right\}. \tag{2}$$

We decided to use this metric, instead of other more common and easier metrics, like the Total Variation (TV) one, because it is “less demanding”, that is, MDPs that are Lipschitz according

¹ Notice that our definition of pointwise Lipschitz function differs from the traditional one.

to TV are always Lipschitz also w.r.t. the L^1 -Wasserstein metric, while the vice versa is not true. For instance, MDPs with deterministic transitions are never Lipschitz according to TV, while they can be Lipschitz using L^1 -Wasserstein metric. Finally, the choice of the L^1 -Wasserstein metric rather than other sophisticated distribution distances is motivated by the fact that it has been frequently used for MDPs (Rachelson and Lagoudakis 2010; Hinderer 2005; Ferns et al. 2005).

The analysis proposed in this paper is based on the assumption that the MDP is Lipschitz continuous. A Lipschitz MDP is a standard MDP enhanced by the information that the transition model is $L_{\mathcal{P}}$ -LC, and the reward model is $L_{\mathcal{R}}$ -LC.

Assumption 1 (Lipschitz MDP) A Lipschitz MDP must satisfy the following conditions:

$$\begin{aligned} \forall (s, \widehat{s}, a, \widehat{a}) \in \mathcal{S}^2 \times \mathcal{A}^2, \quad \mathcal{K}(\mathcal{P}(\cdot|s, a), \mathcal{P}(\cdot|\widehat{s}, \widehat{a})) &\leq L_{\mathcal{P}} d_{\mathcal{S}\mathcal{A}}((s, a), (\widehat{s}, \widehat{a})), \\ \forall (s, \widehat{s}, a, \widehat{a}) \in \mathcal{S}^2 \times \mathcal{A}^2, \quad |\mathcal{R}(s, a) - \mathcal{R}(\widehat{s}, \widehat{a})| &\leq L_{\mathcal{R}} d_{\mathcal{S}\mathcal{A}}((s, a), (\widehat{s}, \widehat{a})). \end{aligned}$$

If π is an L_{π} -LC policy— $\forall (s, \widehat{s}) \in \mathcal{S}^2, \mathcal{K}(\pi(\cdot|s), \pi(\cdot|\widehat{s})) \leq L_{\pi} d_{\mathcal{S}}(s, \widehat{s})$ —, under Assumption 1, it is possible to prove the LC of the corresponding value functions.

Lemma 1 (Rachelson and Lagoudakis 2010, Lemma 1, Theorem 1) *Given an $(L_{\mathcal{P}}, L_{\mathcal{R}})$ -LC MDP and a L_{π} -LC stationary policy π , if $\gamma L_{\mathcal{P}}(1 + L_{\pi}) < 1$, then the Q -function Q^{π} is $L_{Q^{\pi}}$ -LC and the V -function is $L_{V^{\pi}}$ -LC w.r.t. the joint state-action space:*

$$L_{Q^{\pi}} = \frac{L_{\mathcal{R}}}{1 - \gamma L_{\mathcal{P}}(1 + L_{\pi})}; \quad L_{V^{\pi}} = L_{Q^{\pi}}(1 + L_{\pi}).$$

All these conditions are related to state and action variables². In particular, the Lipschitz continuity of the V - and Q -functions means that: $\forall (s, \widehat{s}, a, \widehat{a}) \in \mathcal{S}^2 \times \mathcal{A}^2$

$$\begin{aligned} |Q^{\pi}(s, a) - Q^{\pi}(\widehat{s}, \widehat{a})| &\leq L_{Q^{\pi}} d_{\mathcal{S}\mathcal{A}}((s, a), (\widehat{s}, \widehat{a})); \\ |V^{\pi}(s) - V^{\pi}(\widehat{s})| &\leq L_{V^{\pi}} d_{\mathcal{S}}(s, \widehat{s}). \end{aligned}$$

In the following, we will consider the Lipschitz continuity related to the policy parametrization.

2.3 Policy space

We consider the problem of finding a policy that maximizes the expected discounted reward over a class of parameterized policies $\Pi_{\Theta} = \{\pi^{\theta} : \theta \in \Theta \subset \mathbb{R}^d\}$, where π^{θ} is a compact representation of $\pi^{\theta}(a|s)$. Moreover, we assume that (Θ, d_{Θ}) is a metric space. For ease of notation, in the following we will use θ to denote the dependence on π^{θ} where possible.

The exact gradient of the expected discounted reward J_{μ}^{θ} w.r.t. the policy parameters is (Sutton et al. 1999):

$$\begin{aligned} \nabla_{\theta} J_{\mu}^{\theta} &= \frac{1}{1-\gamma} \iint_{\mathcal{S}\mathcal{A}} \nabla_{\theta} \log \pi^{\theta}(a|s) Q^{\theta}(s, a) \pi^{\theta}(da|s) d_{\mu}^{\theta}(ds) \\ &= \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \xi_{\mu}^{\theta}} [\nabla_{\theta} \log \pi^{\theta}(a|s) Q^{\theta}(s, a)]. \end{aligned} \tag{3}$$

² It can be noticed that the results in Lemma 1 are quite similar to the ones obtained by Hinderer (2005, Theorem 4.1a). The main difference is due to the fact that Hinderer focuses on the Lipschitz continuity of the optimal value function under state-dependent action spaces. In this case, the role of the Lipschitz constant L_{π} is taken by the Lipschitz constant due to state-dependent action spaces. Since we consider policy-based value functions, the eventual differences between action spaces is implicitly coded in the policy.

Several studies have focused on computing the value of this gradient from sample trajectories, trying to produce estimators with low variance (Peters and Schaal 2008b). The policy parameters can be updated by following the direction of the gradient of the expected discounted reward: $\theta' = \theta + \alpha \nabla_{\theta} J_{\mu}^{\theta}$, where α is a parameter used to control the step size.

For proving the results in the next section, we need to introduce the following assumptions on the parameterized policy model.

Assumption 2 (Lipschitz policies) The policy model must satisfy the following conditions:

- 1) state-action LC: $\forall (s, \hat{s}) \in \mathcal{S}^2, \mathcal{K}(\pi^{\theta}(\cdot|s), \pi^{\theta}(\cdot|\hat{s})) \leq L_{\pi^{\theta}} d_{\mathcal{S}}(s, \hat{s})$
- 2) parametric PLC: $\forall s \in \mathcal{S}, \forall (\theta, \hat{\theta}) \in \Theta^2, \mathcal{K}(\pi^{\theta}(\cdot|s), \pi^{\hat{\theta}}(\cdot|s)) \leq L_{\pi}(\theta) d_{\Theta}(\theta, \hat{\theta})$

Assumption 3 (Lipschitz gradient of policy logarithm) The gradient of the policy logarithm must satisfy the following conditions of:

- 1) uniformly bounded gradient: $\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \forall \theta \in \Theta, \forall i = 1, \dots, d$

$$\left| \nabla_{\theta_i} \log \pi^{\theta}(a|s) \right| \leq M_{\theta}^i$$

- 2) state-action LC: $\forall (s, \hat{s}, a, \hat{a}) \in \mathcal{S}^2 \times \mathcal{A}^2, \forall \theta \in \Theta, \forall i = 1, \dots, d$

$$\left| \nabla_{\theta_i} \log \pi^{\theta}(a|s) - \nabla_{\theta_i} \log \pi^{\theta}(\hat{a}|\hat{s}) \right| \leq L_{\nabla \log \pi^{\theta}}^i d_{\mathcal{S}\mathcal{A}}((s, a), (\hat{s}, \hat{a}))$$

- 3) parametric PLC: $\forall (\theta, \hat{\theta}) \in \Theta^2, \forall (s, a) \in \mathcal{S} \times \mathcal{A}, \forall i = 1, \dots, d$

$$\left| \nabla_{\theta_i} \log \pi^{\theta}(a|s) - \nabla_{\theta_i} \log \pi^{\hat{\theta}}(a|s) \right| \leq L_{\nabla \log \pi}^i(\theta) d_{\Theta}(\theta, \hat{\theta})$$

Notice that some previously defined Lipschitz constants become θ -dependent when a parametric policy model is used, i.e., $L_{\pi^{\theta}}, L_{Q^{\theta}}$ and $L_{V^{\theta}}$, whereas $L_{\mathcal{P}}$ and $L_{\mathcal{R}}$ are θ -independent since they are not affected by the policy.

3 Lipschitz continuity of the expected return

In this section, we will show that under Assumptions 1, 2, and 3 the expected return J_{μ}^{θ} is a Lipschitz function w.r.t. policy parameters θ . Besides being an interesting objective by itself, it allows us to introduce some preliminary results that will be reused in the next section to make the proof of the Lipschitz property of the policy gradient easier.

The performance distance between two policies corresponding to parameters θ and $\hat{\theta}$ is measured by the absolute difference of their expected returns:

$$\left| J_{\mu}^{\theta} - J_{\mu}^{\hat{\theta}} \right| = \frac{1}{1 - \gamma} \left| \mathbb{E}_{(s,a) \sim \zeta_{\mu}^{\theta}} [\mathcal{R}(s, a)] - \mathbb{E}_{(s,a) \sim \zeta_{\mu}^{\hat{\theta}}} [\mathcal{R}(s, a)] \right|.$$

If the Lipschitz constant of the reward function \mathcal{R} were less or equal to 1 (i.e., $\|\mathcal{R}\|_L \leq 1$), it follows from (2) that the performance distance between policies π^{θ} and $\pi^{\hat{\theta}}$ would be upper bounded by the Kantorovich distance between the distributions ζ_{μ}^{θ} and $\zeta_{\mu}^{\hat{\theta}}$. On the other hand, it can be easily shown that if \mathcal{R} is $L_{\mathcal{R}}$ -LC then $\left\| \frac{\mathcal{R}}{L_{\mathcal{R}}} \right\|_L \leq 1$.

The following proposition gives an upper bound to absolute difference in performance between policies. Proof can be founded in ‘‘Proof of Proposition 1’’ section of Appendix.

Proposition 1 *Given an $L_{\mathcal{R}}$ -LC MDP, for any pair of stationary policies corresponding to parameters θ and $\widehat{\theta}$, the absolute difference between the performance of policy π^θ and policy $\pi^{\widehat{\theta}}$ can be bounded as follows:*

$$\left| J_\mu^\theta - J_\mu^{\widehat{\theta}} \right| \leq \frac{L_{\mathcal{R}}}{1 - \gamma} \mathcal{K} \left(\zeta_\mu^\theta, \zeta_\mu^{\widehat{\theta}} \right).$$

As a consequence, to prove the Lipschitz continuity of the expected return w.r.t. θ it suffices to show the Lipschitz continuity of the distribution ζ_μ^θ w.r.t. θ . It is worth recalling that the distribution ζ_μ^θ is defined over the joint state-action space and the probability of drawing a state-action pair (s, a) is $\delta_\mu^\theta(s) \cdot \pi^\theta(a|s)$. As it can be noticed, the probability distribution over actions and the one over states are not independent. This means that the Kantorovich distance of the joint distribution cannot be simply upper bounded by the sum of the Kantorovich distances of the γ -discounted future state distribution δ_μ and the policy π . The following Lemma gives an upper bound to $\mathcal{K} \left(\zeta_\mu^\theta, \zeta_\mu^{\widehat{\theta}} \right)$.

Lemma 2 *Given an L_{π^θ} -LC and $L_\pi(\theta)$ -PLC stationary policy π^θ , the Kantorovich distance between a pair of joint distributions ζ_μ^θ and $\zeta_\mu^{\widehat{\theta}}$ is bounded by:*

$$\mathcal{K} \left(\zeta_\mu^\theta, \zeta_\mu^{\widehat{\theta}} \right) \leq L_\pi(\theta) d_\Theta \left(\theta, \widehat{\theta} \right) + \left(1 + L_{\pi^\theta} \right) \mathcal{K} \left(\delta_\mu^\theta, \delta_\mu^{\widehat{\theta}} \right).$$

Proof The proof is divided into two parts. The first part is devoted to the analysis of the Lipschitz continuity of a term involved in the definition of the L^1 -Wasserstein metric between the joint distributions. The second part exploits this result to prove the lemma.

Define $b_f^\theta(s) = \mathbb{E}_{a \sim \pi^\theta} f(s, a)$, where the function f is 1-LC w.r.t. the joint state-action space. Given an L_{π^θ} -LC policy model, $b(s)$ is Lipschitz continuous:

$$\begin{aligned} \left| b_f^\theta(s) - b_f^\theta(\widehat{s}) \right| &= \left| \int_{\mathcal{A}} \pi^\theta(a|s) f(s, a) - \pi^\theta(a|\widehat{s}) f(\widehat{s}, a) da \right| \\ &= \left| \int_{\mathcal{A}} \left(\pi^\theta(a|s) - \pi^\theta(a|\widehat{s}) \right) f(s, a) + (f(s, a) - f(\widehat{s}, a)) \pi(a|\widehat{s}) da \right| \\ &\leq \left| \int_{\mathcal{A}} \left(\pi^\theta(a|s) - \pi^\theta(a|\widehat{s}) \right) f(s, a) da \right| + \left| \int_{\mathcal{A}} (1 \cdot d_S(s, \widehat{s})) \pi(a|\widehat{s}) da \right| \\ &\leq \mathcal{K} \left(\pi^\theta(\cdot|s), \pi^\theta(\cdot|\widehat{s}) \right) + d_S(s, \widehat{s}) \\ &\leq (L_{\pi^\theta} + 1) d_S(s, \widehat{s}). \end{aligned} \tag{4}$$

Recall that the function f in the definition of the L^1 -Wasserstein metric for the joint distributions is 1-LC w.r.t. every pair (s, a) , but as a consequence it is, at most, 1-LC for the single variables s and a . The proof follows from the previous result and some algebraic manipulations:

$$\begin{aligned} \mathcal{K} \left(\zeta_\mu^\theta, \zeta_\mu^{\widehat{\theta}} \right) &= \sup_f \left\{ \left| \int_S \delta_\mu^\theta(s) \int_{\mathcal{A}} \pi^\theta(a|s) f(s, a) dads \right. \right. \\ &\quad \left. \left. - \int_S \delta_\mu^{\widehat{\theta}}(s) \int_{\mathcal{A}} \pi^{\widehat{\theta}}(a|s) f(s, a) dads \right| : \|f\|_L \leq 1 \right\} \\ &= \sup_f \left\{ \left| \int_S \left(\delta_\mu^\theta(s) - \delta_\mu^{\widehat{\theta}}(s) \right) \int_{\mathcal{A}} \pi^\theta(a|s) f(s, a) dads \right. \right. \end{aligned} \tag{5}$$

$$\begin{aligned}
 & + \int_{\mathcal{S}} \delta_{\mu}^{\widehat{\theta}}(s) \int_{\mathcal{A}} \left(\pi^{\theta}(a|s) - \pi^{\widehat{\theta}}(a|s) \right) f(s, a) \text{d}a \text{d}s \Big| : \|f\|_L \leq 1 \Big\} \tag{6} \\
 & \leq \sup_f \left\{ \left| \int_{\mathcal{S}} \left(\delta_{\mu}^{\theta}(s) - \delta_{\mu}^{\widehat{\theta}}(s) \right) b_f^{\theta}(s) \text{d}s \right| : \|f\|_L \leq 1 \right\} \\
 & \quad + \sup_f \left\{ \left| \int_{\mathcal{S}} \delta_{\mu}^{\widehat{\theta}}(s) \int_{\mathcal{A}} \left(\pi^{\theta}(a|s) - \pi^{\widehat{\theta}}(a|s) \right) f(s, a) \text{d}a \text{d}s \right| : \|f\|_L \leq 1 \right\} \tag{7} \\
 & = (L_{\pi^{\theta}} + 1) \sup_f \left\{ \left| \int_{\mathcal{S}} \left(\delta_{\mu}^{\theta}(s) - \delta_{\mu}^{\widehat{\theta}}(s) \right) \frac{b_f^{\theta}(s)}{(L_{\pi^{\theta}} + 1)} \text{d}s \right| : \|f\|_L \leq 1 \right\} \\
 & \quad + \sup_f \left\{ \left| \int_{\mathcal{S}} \delta_{\mu}^{\widehat{\theta}}(s) \int_{\mathcal{A}} \left(\pi^{\theta}(a|s) - \pi^{\widehat{\theta}}(a|s) \right) f(s, a) \text{d}a \text{d}s \right| : \|f\|_L \leq 1 \right\} \tag{8} \\
 & \leq (L_{\pi^{\theta}} + 1) \sup_g \left\{ \left| \int_{\mathcal{S}} \left(\delta_{\mu}^{\theta}(s) - \delta_{\mu}^{\widehat{\theta}}(s) \right) g(s) \text{d}s \right| : \|g\|_L \leq 1 \right\} \\
 & \quad + \int_{\mathcal{S}} \delta_{\mu}^{\widehat{\theta}}(s) \sup_f \left\{ \left| \int_{\mathcal{A}} \left(\pi^{\theta}(a|s) - \pi^{\widehat{\theta}}(a|s) \right) f(s, a) \text{d}a \right| : \|f\|_L \leq 1 \right\} \text{d}s \tag{9} \\
 & \leq (L_{\pi^{\theta}} + 1) \mathcal{K} \left(\delta_{\mu}^{\theta}, \delta_{\mu}^{\widehat{\theta}} \right) + \sup_s \mathcal{K} \left(\pi^{\theta}(\cdot|s), \pi^{\widehat{\theta}}(\cdot|s) \right) \\
 & \leq (L_{\pi^{\theta}} + 1) \mathcal{K} \left(\delta_{\mu}^{\theta}, \delta_{\mu}^{\widehat{\theta}} \right) + L_{\pi^{\theta}}(\theta) d_{\Theta}(\theta, \widehat{\theta}). \tag{10}
 \end{aligned}$$

Equality (6) is obtained by manipulation of (5) after insertion of the quantity $\pm \int_{\mathcal{S}} \delta_{\mu}^{\widehat{\theta}}(s) \int_{\mathcal{A}} \pi^{\theta}(a|s) f(s, a) \text{d}a \text{d}s$. Eq. (8) is obtained by exploiting definition of $b_f^{\theta}(s)$ and adding the (identity) factor $\frac{L_{\pi^{\theta}} + 1}{L_{\pi^{\theta}} + 1}$. In (9) we rename $\frac{b_f^{\theta}(s)}{L_{\pi^{\theta}} + 1}$ to $g(s)$ and we note that, according to Eq. (4), $\|g(s)\|_L = \left\| \frac{b_f^{\theta}(s)}{L_{\pi^{\theta}} + 1} \right\|_L \leq 1$. By noting that $\delta_{\mu}^{\widehat{\theta}}$ is always positive, a valid upper bound to the second term in (8) is obtained by pushing the supreme over function space into the state integral. Finally, definition of the Kantorovich distance is used to obtain inequality (10) together with a maximization over the state space. The proof follows from Assumption 2. \square

The first term of the upper bound derives from the bound on the Kantorovich distance between policies w.r.t. parameters θ . The second term involves the Kantorovich distance between γ -discounted future state distributions w.r.t. parameters θ and the factor $(1 + L_{\pi^{\theta}})$ accounts for the dependence between the distribution over the actions and the one over the states: the larger is the $L_{\pi^{\theta}}$ constant the stronger is the dependence between π^{θ} and δ_{μ}^{θ} . As expected, when the policy does not depend on the state (i.e., $L_{\pi^{\theta}} = 0$), the bound reduces to the sum of the two Kantorovich distances. The following lemma shows that under Assumptions 1 and 2 also $\mathcal{K} \left(\delta_{\mu}^{\theta}, \delta_{\mu}^{\widehat{\theta}} \right)$ is Lipschitz w.r.t. θ .

Lemma 3 *Given an $L_{\mathcal{P}}$ -LC MDP and an $(L_{\pi^{\theta}}, L_{\pi}(\theta))$ -LC stationary policy model, if $\gamma L_{\mathcal{P}} (1 + L_{\pi^{\theta}}) < 1$, then the Kantorovich distance between a pair of γ -discounted future-state distributions is PLC w.r.t. paramters $\theta: \forall(\theta, \widehat{\theta}) \in \Theta^2$,*

$$\mathcal{K} \left(\delta_{\mu}^{\theta}, \delta_{\mu}^{\widehat{\theta}} \right) \leq L_{\delta}(\theta) d_{\Theta}(\theta, \widehat{\theta}), \quad \text{where } L_{\delta}(\theta) = \frac{\gamma L_{\mathcal{P}} L_{\pi}(\theta)}{1 - \gamma L_{\mathcal{P}} (1 + L_{\pi^{\theta}})}.$$

Proof We start the proof with some preliminary results that will be used in the rest of the proof. Let the function $g_f(s, a) = \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s,a)} f(s')$ where $\|f\|_L \leq 1$. Then, $g_f(s, a)$ is LC

w.r.t. the action variable:

$$\begin{aligned}
 |g_f(s, a) - g_f(s, \widehat{a})| &= \left| \int_{\mathcal{S}} (\mathcal{P}(s'|s, a) - \mathcal{P}(s'|s, \widehat{a})) f(s') ds' \right| \\
 &\leq \mathcal{K}(\mathcal{P}(\cdot|s, a), \mathcal{P}(\cdot|s, \widehat{a})) \leq L_{\mathcal{P}} d_{\mathcal{A}}(a, \widehat{a}). \tag{11}
 \end{aligned}$$

The second result involves the expectation of $g_f(s, a)$ w.r.t. policy π^θ . Let $h_f^\theta(s) = \mathbb{E}_{a \sim \pi^\theta(\cdot|s)} g_f(s, a)$, we can prove that it is Lipschitz continuous:

$$\begin{aligned}
 |h_f^\theta(s) - h_f^\theta(\widehat{s})| &= \left| \iint_{\mathcal{A}\mathcal{S}} f(s') \left(\pi^\theta(a|s) \mathcal{P}(s'|s, a) - \pi^\theta(a|\widehat{s}) \mathcal{P}(s'|\widehat{s}, a) \right) ds' da \right| \\
 &= \left| L_{\mathcal{P}} \int_{\mathcal{A}} \left(\pi^\theta(a|s) - \pi^\theta(a|\widehat{s}) \right) \int_{\mathcal{S}} \frac{\mathcal{P}(s'|s, a)}{L_{\mathcal{P}}} f(s') ds' da \right| \tag{12}
 \end{aligned}$$

$$+ \left| \int_{\mathcal{A}} \pi^\theta(a|\widehat{s}) \int_{\mathcal{S}} \mathcal{P}(s'|s, a) f(s') - \mathcal{P}(s'|\widehat{s}, a) f(s') ds' da \right| \tag{13}$$

$$\leq L_{\mathcal{P}} \mathcal{K}(\pi^\theta(\cdot|s), \pi^\theta(\cdot|\widehat{s})) + \sup_a \mathcal{K}(\mathcal{P}(\cdot|s, a), \mathcal{P}(\cdot|\widehat{s}, a)) \tag{14}$$

$$\leq L_{\mathcal{P}} (L_{\pi^\theta} + 1) d_{\mathcal{S}}(s, \widehat{s}), \tag{15}$$

where (13) is obtained by adding and subtracting the term $\iint_{\mathcal{A}\mathcal{S}} \pi^\theta(a|\widehat{s}) \mathcal{P}(s'|s, a) ds' da$. Inequality (14) follows from Kantorovich distance and bound (11), that is $\left\| \frac{\mathbb{E}_{s' \in \mathcal{P}} f(s')}{L_{\mathcal{P}}} \right\|_L \leq 1$, given that f is 1-LC. Then

$$\begin{aligned}
 \mathcal{K}(\delta_\mu^\theta, \widehat{\delta}_\mu^\theta) &= \sup_f \left\{ \left| \int_{\mathcal{S}} (\delta_\mu^\theta(s) - \widehat{\delta}_\mu^\theta(s)) f(s) ds \right| : \|f\|_L \leq 1 \right\} \\
 &= \sup_f \left\{ \left| \int_{\mathcal{S}} \left(\mu(s) + \gamma \int_{\mathcal{S}} \int_{\mathcal{A}} \pi^\theta(a|s') \mathcal{P}(s|s', a) \delta_\mu^\theta(s') da ds' \right) f(s) \right. \right. \\
 &\quad \left. \left. - \left(\mu(s) + \gamma \int_{\mathcal{S}} \int_{\mathcal{A}} \pi^{\widehat{\theta}}(a|s') \mathcal{P}(s|s', a) \widehat{\delta}_\mu^\theta(s') da ds' \right) f(s) ds \right| : \|f\|_L \leq 1 \right\} \tag{16}
 \end{aligned}$$

$$\begin{aligned}
 &= \gamma \sup_f \left\{ \left| \int_{\mathcal{S}} f(s) \int_{\mathcal{S}} \int_{\mathcal{A}} \mathcal{P}(s|s', a) \left(\pi^\theta(a|s') \delta_\mu^\theta(s') - \pi^{\widehat{\theta}}(a|s') \widehat{\delta}_\mu^\theta(s') \right) \right. \right. \\
 &\quad \left. \left. \times da ds' ds \right| : \|f\|_L \leq 1 \right\}
 \end{aligned}$$

$$\begin{aligned}
 &= \gamma \sup_f \left\{ \left| \int_{\mathcal{S}} f(s) \int_{\mathcal{S}} \int_{\mathcal{A}} \mathcal{P}(s|s', a) \left(\left(\pi^\theta(a|s') - \pi^{\widehat{\theta}}(a|s') \right) \delta_\mu^\theta(s') \right. \right. \right. \\
 &\quad \left. \left. + \left(\delta_\mu^\theta(s') - \widehat{\delta}_\mu^\theta(s') \right) \pi^\theta(a|s') \right) da ds' ds \right| : \|f\|_L \leq 1 \right\}
 \end{aligned}$$

$$\begin{aligned}
 &\leq \gamma \sup_f \left\{ \left| \int_{\mathcal{S}} \widehat{\delta}_\mu^\theta(s') \int_{\mathcal{A}} \left(\pi^\theta(a|s') - \pi^{\widehat{\theta}}(a|s') \right) \right. \right. \\
 &\quad \left. \left. \times \int_{\mathcal{S}} \mathcal{P}(s|s', a) f(s) ds da ds' \right| : \|f\|_L \leq 1 \right\} \tag{17}
 \end{aligned}$$

$$\begin{aligned}
 &+ \gamma \sup_f \left\{ \left| \int_{\mathcal{S}} \left(\delta_\mu^\theta(s') - \widehat{\delta}_\mu^\theta(s') \right) \int_{\mathcal{A}} \pi^\theta(a|s') \right. \right. \\
 &\quad \left. \left. \times \int_{\mathcal{S}} \mathcal{P}(s|s', a) f(s) ds da ds' \right| : \|f\|_L \leq 1 \right\} \tag{18}
 \end{aligned}$$

$$\begin{aligned}
 &= \gamma L_{\mathcal{P}} \sup_f \left\{ \left| \int_{\mathcal{S}} \delta_{\mu}^{\widehat{\theta}}(s') \int_{\mathcal{A}} \left(\pi^{\theta}(a|s') - \pi^{\widehat{\theta}}(a|s') \right) \frac{g_f(s', a)}{L_{\mathcal{P}}} da ds' \right| : \|f\|_L \leq 1 \right\} \\
 &\quad + \gamma L_{\mathcal{P}} (L_{\pi^{\theta}} + 1) \sup_f \left\{ \left| \int_{\mathcal{S}} \left(\delta_{\mu}^{\theta}(s') - \delta_{\mu}^{\widehat{\theta}}(s') \right) \frac{h_f^{\theta}(s')}{L_{\mathcal{P}}(L_{\pi^{\theta}} + 1)} ds' \right| : \|f\|_L \leq 1 \right\} \\
 &\leq \gamma L_{\mathcal{P}} \int_{\mathcal{S}} \delta_{\mu}^{\widehat{\theta}}(s') \sup_{\bar{f}} \left\{ \left| \int_{\mathcal{A}} \left(\pi^{\theta}(a|s') - \pi^{\widehat{\theta}}(a|s') \right) \bar{f}(s', a) da \right| : \|\bar{f}\|_{L,a} \leq 1 \right\} \\
 &\quad + \gamma L_{\mathcal{P}} (L_{\pi^{\theta}} + 1) \sup_{\tilde{f}} \left\{ \left| \int_{\mathcal{S}} \left(\delta_{\mu}^{\theta}(s') - \delta_{\mu}^{\widehat{\theta}}(s') \right) \tilde{f}(s') ds' \right| : \|\tilde{f}\|_L \leq 1 \right\} \tag{19}
 \end{aligned}$$

$$\leq \gamma L_{\mathcal{P}} \sup_{s'} \mathcal{K} \left(\pi^{\theta}(\cdot|s'), \pi^{\widehat{\theta}}(\cdot|s') \right) + \gamma L_{\mathcal{P}} (L_{\pi^{\theta}} + 1) \mathcal{K} \left(\delta_{\mu}^{\theta}, \delta_{\mu}^{\widehat{\theta}} \right) \tag{20}$$

$$\leq \gamma L_{\mathcal{P}} (L_{\pi}(\theta) + (1 + L_{\pi^{\theta}})L_{\delta}(\theta)) d(\theta, \widehat{\theta}). \tag{21}$$

By replacing δ_{μ}^{θ} with its definition we get equality (16). By adding and subtracting the term $\int \int \int_{\mathcal{S}\mathcal{S}\mathcal{A}} f(s)\mathcal{P}(s|s', a)\pi^{\theta}(a|s')\delta_{\mu}^{\widehat{\theta}}(s')dad s' ds$ and resorting to the triangle inequality, we derive lines (17) and (18). Such terms can be simplified by noting that they contain definition of $g_f(s, a)$ and $h_f^{\theta}(s)$, respectively. After insertion of invariant scaling factors we rename $\bar{f}(s', a) = \frac{g_f(s', a)}{L_{\mathcal{P}}}$ and $\tilde{f}(s') = \frac{h_f^{\theta}(s')}{L_{\mathcal{P}}(L_{\pi^{\theta}}+1)}$ in (19). Let $\|z(x, y)\|_{L,y}$ be the Lipschitz semi-norm of function z w.r.t. only variable y (taking the supremum over x). Then, from inequality (11), it is easy to see that $\|\bar{f}\|_{L,a} \leq 1$. Similarly, from inequality (15), we derive that $\|\tilde{f}\|_L \leq 1$. As done in the proof of Lemma 2, we push the supremum into the state integral and we maximize the Kantorovich distance [inequality (20)].

Note that inequality (21) leads to the following fixed point equation:

$$L_{\delta}(\theta) = \gamma L_{\mathcal{P}} [L_{\pi}(\theta) + (1 + L_{\pi^{\theta}})L_{\delta}(\theta)]$$

that admits a unique feasible solution only if $\gamma L_{\mathcal{P}} (1 + L_{\pi^{\theta}}) < 1$. □

As expected, the smoothness of the γ -discounted future state distribution w.r.t. to θ strongly depends on the smoothness of the state transition model and the policy model. In particular, a relevant role is played by $L_{\mathcal{P}}$ that influences both the numerator and the denominator of L_{δ} (decreasing $L_{\mathcal{P}}$ decreases the value of the numerator and increases the value of the denominator). As in the case of the Lipschitz continuity of the Q - and V -functions (see Lemma 1), the Lipschitz continuity of the γ -discounted future state distribution can be guarantee only when the condition $\gamma L_{\mathcal{P}} (1 + L_{\pi^{\theta}}) < 1$ holds. This condition emerges from the recursive nature of the considered functions and enforces the discounted Markov kernel underlying policy π to be a contraction w.r.t. the Kantorovich distance.

Finally, combining Proposition 1 with Lemmas 2 and 3, we can derive the Lipschitz continuity of the joint distribution ζ_{μ} .

Lemma 4 *Under Assumption 1 and 2, if $\gamma L_{\mathcal{P}} (1 + L_{\pi^{\theta}}) < 1$, then the joint distribution ζ_{μ}^{θ} is $L_{\zeta}(\theta)$ -PLC w.r.t. the policy parameters θ , with:*

$$L_{\zeta}(\theta) = \frac{L_{\pi}(\theta)}{1 - \gamma L_{\mathcal{P}} (1 + L_{\pi^{\theta}})}.$$

The lemma comes directly from the application of Lemma 3 to Lemma 2. Now we have all the technicalities required to derive the main theorem.

Theorem 1 *Given an $(L_{\mathcal{P}}, L_{\mathcal{R}})$ -LC MDP and an $(L_{\pi\theta}, L_{\pi}(\theta))$ -LC stationary policy model, if $\gamma L_{\mathcal{P}}(1 + L_{\pi\theta}) < 1$, then the performance measure J_{μ}^{θ} is $L_J(\theta)$ -PLC w.r.t. the policies parameters:*

$$\left| J_{\mu}^{\theta} - J_{\mu}^{\hat{\theta}} \right| \leq L_J(\theta) d_{\Theta}(\theta, \hat{\theta}),$$

with:

$$L_J(\theta) = \frac{L_{\mathcal{R}} L_{\pi}(\theta)}{(1 - \gamma)(1 - \gamma L_{\mathcal{P}}(1 + L_{\pi\theta}))}.$$

The proof follows from the application of Lemma 4 to Proposition 1. Although L_J resembles the Lipschitz constant of the Q -function (actually $L_J = \frac{L_{\mathcal{R}}}{1-\gamma} L_{Q^{\pi}}$), it is worth underline that they define Lipschitz conditions over different spaces: the expected return is L_J -PLC w.r.t. policy parameters θ , while $L_{Q^{\pi}}$ is a Lipschitz constant w.r.t. the state-action space.

4 Lipschitz continuity of the policy gradient

Leveraging on the results presented in the previous section, here we investigate the Lipschitz continuity of the gradient of the expected discounted reward w.r.t. the policy parameters. Both the expected return of a policy π^{θ} and its gradient can be defined as expected values w.r.t. distribution ζ_{μ}^{θ} : in the former case the function to be averaged is the reward function $\mathcal{R}(s, a)$, while in the latter case is $\nabla_{\theta} \log \pi^{\theta}(s, a) Q^{\theta}(s, a)$ (see Eq. 3). For ease of notation, we define the function $\eta : \mathcal{S} \times \mathcal{A} \times \Theta \rightarrow \mathbb{R}^d$ as

$$\eta^{\theta}(s, a) = \nabla_{\theta} \log \pi^{\theta}(s, a) Q^{\theta}(s, a).$$

In particular, we consider the component-wise absolute difference between gradients corresponding to different parameterizations:

$$\left| \nabla_{\theta_i} J_{\mu}^{\theta} - \nabla_{\theta_i} J_{\mu}^{\hat{\theta}} \right| = \frac{1}{1 - \gamma} \left| \mathbb{E}_{(s,a) \sim \zeta_{\mu}^{\theta}} \left[\eta_i^{\theta}(s, a) \right] - \mathbb{E}_{(s,a) \sim \zeta_{\mu}^{\hat{\theta}}} \left[\eta_i^{\hat{\theta}}(s, a) \right] \right|.$$

It is worth nothing that, differently from the reward function in the expected-return case, functions η^{θ} do depend on the policy parameters θ . This prevents to follow immediately the same steps as done in the previous section and requires to decompose the problem by introducing the following upper bound.

Proposition 2 *For any pair of stationary policies corresponding to parameters θ and $\hat{\theta}$, the component-wise absolute difference between the gradients of the expected return can be upper bounded as follows:*

$$\begin{aligned} \left| \nabla_{\theta_i} J_{\mu}^{\theta} - \nabla_{\theta_i} J_{\mu}^{\hat{\theta}} \right| &\leq \frac{1}{1 - \gamma} \left| \mathbb{E}_{(s,a) \sim \zeta_{\mu}^{\theta}} \left[\eta_i^{\theta}(s, a) \right] - \mathbb{E}_{(s,a) \sim \zeta_{\mu}^{\hat{\theta}}} \left[\eta_i^{\theta}(s, a) \right] \right| \\ &\quad + \frac{1}{1 - \gamma} \left| \mathbb{E}_{(s,a) \sim \zeta_{\mu}^{\hat{\theta}}} \left[\eta_i^{\theta}(s, a) - \eta_i^{\hat{\theta}}(s, a) \right] \right|. \end{aligned}$$

While the second term requires a further expansion (that will be presented later in the section), the first one can be bounded following a similar argument as the one used for the expected return.

Upper bound to the first term. While the bound of the expected reward (Lemma 1) follows directly from the definition of Lipschitz MDP (Assumption 1), here, we need to prove that the function η_i^θ is Lipschitz w.r.t. the joint state-action space. Since the product of two Lipschitz functions is Lipschitz, given Assumption 2 and Lemma 1, we can show that η^θ is LC w.r.t. the state-action space (see “Proof of Lemma 5” section of Appendix for the proof).

Lemma 5 *Under Assumptions 1, 2 and 3, the i -th component of η^θ is $L_{\eta^\theta}^i$ -LC w.r.t. the state-action space, that is: $\forall (s, \widehat{s}, a, \widehat{a}) \in \mathcal{S}^2 \times \mathcal{A}^2$,*

$$\left| \eta_i^\theta(s, a) - \eta_i^\theta(\widehat{s}, \widehat{a}) \right| \leq L_{\eta^\theta}^i d_{SA}((s, a), (\widehat{s}, \widehat{a})),$$

where $L_{\eta^\theta}^i = \frac{R}{1-\gamma} L_{\nabla \log \pi^\theta}^i + M_\theta^i L_{Q^\theta}$.

Since η_i^θ is LC, it follows that

$$\left| \mathbb{E}_{(s,a) \sim \zeta_\mu^\theta} \left[\eta_i^\theta(s, a) \right] - \mathbb{E}_{(s,a) \sim \zeta_\mu^{\widehat{\theta}}} \left[\eta_i^\theta(s, a) \right] \right| \leq L_{\eta^\theta}^i \mathcal{K}(\zeta_\mu^\theta, \zeta_\mu^{\widehat{\theta}}) \leq L_{\eta^\theta}^i L_\zeta(\theta) d_\Theta(\theta, \widehat{\theta}).$$

Upper bound to the second term. To prove its Lipschitz continuity w.r.t. policy parameters, we need to introduce an upper bound to $\left| \eta_i^\theta - \eta_i^{\widehat{\theta}} \right|$. Refer to “Proof of Lemma 6” section of Appendix for the proof.

Lemma 6 *For any pair of stationary policies corresponding to θ and $\widehat{\theta}$, the absolute difference of the i -th component of functions η^θ and $\eta^{\widehat{\theta}}$ is upper bounded by: $\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \forall (\theta, \widehat{\theta}) \in \Theta^2$,*

$$\begin{aligned} \left| \eta_i^\theta(s, a) - \eta_i^{\widehat{\theta}}(s, a) \right| &\leq \frac{R}{1-\gamma} \left| \nabla_{\theta_i} \log \pi^\theta(a|s) - \nabla_{\theta_i} \log \pi^{\widehat{\theta}}(a|s) \right| \\ &\quad + M_\theta^i \left| Q^\theta(s, a) - Q^{\widehat{\theta}}(s, a) \right|. \end{aligned}$$

The first term of the bound can be upper bounded in turn by exploiting the LC assumption on $\nabla \log \pi$ w.r.t. policy parameters (see Assumption 3). For what concerns the second term, we need to show further results about the Lipschitz continuity of the Q - and V -functions w.r.t. the policy parameters. Here we extend the standard Lipschitz framework for MDPs to the case in which a parametric policy model is available.

Theorem 2 *Under Assumptions 1 and 2, the V -function and the Q -function are respectively $L_V(\theta)$ - and $L_Q(\theta)$ -PLC w.r.t. to the policy parameters, with:*

$$L_V(\theta) = \frac{L_\pi(\theta) L_{\mathcal{R}}}{(1-\gamma)(1-\gamma L_{\mathcal{P}}(1+L_{\pi^\theta}))}; \quad L_Q(\theta) = \gamma L_V(\theta).$$

Proof We need to introduce some preliminary results that will be used to prove the main theorem. First of all, we want to prove that the expected reward under two stationary policies π^θ and $\pi^{\widehat{\theta}}$ corresponding to different policy parameterizations is Lipschitz continuous for any state s :

$$\begin{aligned} \left| \mathcal{R}^\theta(s) - \mathcal{R}^{\widehat{\theta}}(s) \right| &= \int_{a \in \mathcal{A}} \left(\pi^\theta(a|s) - \pi^{\widehat{\theta}}(a|s) \right) \mathcal{R}(s, a) da \\ &\leq L_{\mathcal{R}} \mathcal{K} \left(\pi^\theta(\cdot|s), \pi^{\widehat{\theta}}(\cdot|s) \right) \\ &\leq L_{\mathcal{R}} L_\pi(\theta) d_\Theta(\theta, \widehat{\theta}). \end{aligned}$$

The following equation gives an upper bound to the maximum absolute difference of the V -functions associated to two policy parameterizations: $\forall s \in \mathcal{S}$

$$\begin{aligned}
 |V^\theta(s) - V^{\hat{\theta}}(s)| &= \left| \int_{\mathcal{A}} \left(\mathcal{R}(s, a) + \gamma \int_{\mathcal{S}} V^\theta(s') \mathcal{P}(ds'|s, a) \right) \pi^\theta(da|s) \right. \\
 &\quad \left. - \int_{\mathcal{A}} \left(\mathcal{R}(s, a) + \gamma \int_{\mathcal{S}} V^{\hat{\theta}}(s') \mathcal{P}(ds'|s, a) \right) \pi^{\hat{\theta}}(da|s) \right| \\
 &\leq \left| \int_{\mathcal{A}} \mathcal{R}(s, a) \left(\pi^\theta(a|s) - \pi^{\hat{\theta}}(a|s) \right) da \right| \\
 &\quad + \gamma \left| \iint_{\mathcal{S}\mathcal{A}} \mathcal{P}(s'|s, a) \left(\pi^\theta(a|s) V^\theta(s') - \pi^{\hat{\theta}}(a|s) V^{\hat{\theta}}(s') \right) da ds' \right| \\
 &\leq \left| \mathcal{R}^\theta(s) - \mathcal{R}^{\hat{\theta}}(s) \right| + \gamma \left| \int_{\mathcal{A}} \left(\pi^\theta(a|s) - \pi^{\hat{\theta}}(a|s) \right) \right. \\
 &\quad \left. \times \int_{\mathcal{S}} \mathcal{P}(s'|s, a) V^\theta(s') ds' da \right| \\
 &\quad + \gamma \left| \int_{\mathcal{A}} \pi^{\hat{\theta}}(a|s) \int_{\mathcal{S}} \mathcal{P}(s'|s, a) \left(V^\theta(s') - V^{\hat{\theta}}(s') \right) ds' da \right| \\
 &= \left| \mathcal{R}^\theta(s) - \mathcal{R}^{\hat{\theta}}(s) \right| + \gamma L_V L_{\mathcal{P}} \left| \int_{\mathcal{A}} \left(\pi^\theta(a|s) - \pi^{\hat{\theta}}(a|s) \right) \right. \\
 &\quad \left. \times \int_{\mathcal{S}} \frac{\mathcal{P}(s'|s, a) V^\theta(s')}{L_V L_{\mathcal{P}}} ds' da \right| \\
 &\quad + \gamma \left| \int_{\mathcal{A}} \pi^{\hat{\theta}}(a|s) \int_{\mathcal{S}} \mathcal{P}(s'|s, a) \left(V^\theta(s') - V^{\hat{\theta}}(s') \right) ds' da \right| \\
 &\leq \left| \mathcal{R}^\theta(s) - \mathcal{R}^{\hat{\theta}}(s) \right| + \gamma L_V L_{\mathcal{P}} \mathcal{K} \left(\pi^\theta(\cdot|s), \pi^{\hat{\theta}}(\cdot|s) \right) + \gamma L_V(\theta) d_\Theta(\theta, \hat{\theta}) \\
 &\leq (L_{\mathcal{R}} + \gamma L_V L_{\mathcal{P}}) L_\pi(\theta) d_\Theta(\theta, \hat{\theta}) + \gamma L_V(\theta) d_\Theta(\theta, \hat{\theta}) \\
 &\leq \frac{L_\pi(\theta)}{1 - \gamma} (L_{\mathcal{R}} + \gamma L_V L_{\mathcal{P}}) d_\Theta(\theta, \hat{\theta}).
 \end{aligned}$$

The proof follows from the substitution of the value L_V with its definition (Rachelson and Lagoudakis 2010). The following equations provides the Lipschitz continuity of the Q -function: $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned}
 |Q^\theta(s, a) - Q^{\hat{\theta}}(s, a)| &= \gamma \left| \int_{\mathcal{S}} \mathcal{P}(s'|s, a) \left(V^\theta(s') - V^{\hat{\theta}}(s') \right) ds' \right| \\
 &\leq \gamma \int_{\mathcal{S}} \mathcal{P}(s'|s, a) \left| V^\theta(s') - V^{\hat{\theta}}(s') \right| ds' \\
 &\leq \gamma \sup_{s'} \left| V^\theta(s') - V^{\hat{\theta}}(s') \right| \leq \gamma L_V(\theta) d_\Theta(\theta, \hat{\theta}).
 \end{aligned}$$

□

Similarly as done for Lemma 5, we can observe that η^θ is the product of two functions that are Lipschitz continuous w.r.t. parameters θ . Proof can be founded in “Proof of Lemma 7” section of Appendix.

Lemma 7 Under Assumptions 1, 2 and 3, the i -th component of η is $L^i_\eta(\theta)$ -PLC w.r.t. the policy parameters, that is: $\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \forall (\theta, \hat{\theta}) \in \Theta^2,$

$$\left| \eta_i^\theta(s, a) - \eta_i^{\hat{\theta}}(s, a) \right| \leq L^i_\eta(\theta) d_\Theta(\theta, \hat{\theta}),$$

where $L^i_\eta(\theta) = \frac{R}{1-\gamma} L^i_{\nabla \log \pi}(\theta) + M^i_\theta L_Q(\theta).$

Finally, combining Lemmas 5 and 7 (see “Proof of Theorem 3” section of Appendix for the proof), we are ready to state our main result about the Lipschitz continuity of the policy gradient w.r.t. policy parameters $\theta.$

Theorem 3 Under Assumptions 1, 2, and 3, the i -th component of the gradient $\nabla_\theta J$ of the expected return is $L^i_{\nabla J}(\theta)$ -PLC, that is: $\forall (\theta, \hat{\theta}) \in \Theta^2,$

$$\left| \nabla_{\theta_i} J_\mu^\theta - \nabla_{\theta_i} J_\mu^{\hat{\theta}} \right| \leq L^i_{\nabla J}(\theta) d_\Theta(\theta, \hat{\theta}),$$

where $L^i_{\nabla J}(\theta) = \frac{1}{1-\gamma} \left(L^i_{\eta^\theta} L_\zeta(\theta) + L^i_\eta(\theta) \right).$

Given the vector $\mathbf{L}_{\nabla J}(\theta) = [L^1_{\nabla J}(\theta), \dots, L^d_{\nabla J}(\theta)],$ the policy gradient $\nabla_\theta J_\mu^\theta$ is $L_{\nabla J}(\theta)$ -LC in $\Theta,$ where $L_{\nabla J}(\theta) = \|\mathbf{L}_{\nabla J}(\theta)\|_2$ when $d_\Theta(\theta, \hat{\theta}) = \|\theta - \hat{\theta}\|_2.$

5 Updating policy parameters

In this section, we will show how the Lipschitz continuity of the policy gradient discussed in the previous section can be exploited to update the policy parameters and to guarantee performance improvement. The following lemma exploits the Taylor expansion and the Lipschitz continuity of the policy gradient to derive a lower bound to the policy performance improvement.

Lemma 8 If the policy gradient is Lipschitz continuous, the policy performance improvement between policy $\hat{\theta}$ and policy θ can be lower bounded as follows:

$$J_\mu^{\hat{\theta}} - J_\mu^\theta \geq \nabla_\theta J_\mu^{\theta^T} \cdot (\hat{\theta} - \theta) - \frac{1}{2} \|\hat{\theta} - \theta\|_p \mathbf{L}_{\nabla J}(\theta)^T \cdot |\hat{\theta} - \theta|,$$

where $|v|$ denotes the component-wise absolute value when v is a vector.

Proof According to the definition of Lipschitz gradient given in Theorem 3, we have that, for some $t \in [0, 1]:$

$$\begin{aligned} \left| \nabla_{\theta_i} J_\mu^{\hat{\theta}} - \nabla_{\theta_i} J_\mu^\theta \right| &= \left\| \int_0^1 \nabla(\nabla_{\theta_i} J_\mu)(t \hat{\theta} + (1-t)\theta)^T \cdot (\hat{\theta} - \theta) dt \right\|_p \\ &\leq \|\hat{\theta} - \theta\|_p \int_0^1 \|\nabla(\nabla_{\theta_i} J_\mu)(t \hat{\theta} + (1-t)\theta)\|_p dt \\ &\leq L^i_{\nabla J}(\theta) \|\hat{\theta} - \theta\|_p. \end{aligned}$$

If we only consider the first term of the Taylor’s expansion, we can write the following upper bound:

$$\begin{aligned}
 J_{\mu}^{\widehat{\theta}} &= J_{\mu}^{\theta} + \int_0^1 \nabla_{\theta} J_{\mu}^t \widehat{\theta}^{+(1-t)\theta^T} \cdot (\widehat{\theta} - \theta) dt \pm \nabla_{\theta} J_{\mu}^{\theta^T} \cdot (\widehat{\theta} - \theta) \\
 &= J_{\mu}^{\theta} + \nabla_{\theta} J_{\mu}^{\theta^T} \cdot (\widehat{\theta} - \theta) + \int_0^1 \left(\nabla_{\theta} J_{\mu}^t \widehat{\theta}^{+(1-t)\theta} - \nabla_{\theta} J_{\mu}^{\theta} \right)^T \cdot (\widehat{\theta} - \theta) dt \\
 &= J_{\mu}^{\theta} + \nabla_{\theta} J_{\mu}^{\theta^T} \cdot (\widehat{\theta} - \theta) + \sum_i \int_0^1 \left(\nabla_{\theta_i} J_{\mu}^t \widehat{\theta}^{+(1-t)\theta} - \nabla_{\theta_i} J_{\mu}^{\theta} \right) (\widehat{\theta}_i - \theta_i) dt \\
 &\geq J_{\mu}^{\theta} + \nabla_{\theta} J_{\mu}^{\theta^T} \cdot (\widehat{\theta} - \theta) - \sum_i \int_0^1 \left| \left(\nabla_{\theta_i} J_{\mu}^t \widehat{\theta}^{+(1-t)\theta} - \nabla_{\theta_i} J_{\mu}^{\theta} \right) \right| |\widehat{\theta}_i - \theta_i| dt \\
 &\geq J_{\mu}^{\theta} + \nabla_{\theta} J_{\mu}^{\theta^T} \cdot (\widehat{\theta} - \theta) - \int_0^1 \|t \widehat{\theta} + (1-t)\theta - \theta\|_p dt \sum_i L_{\nabla J}^i |\widehat{\theta}_i - \theta_i| \\
 &\geq J_{\mu}^{\theta} + \nabla_{\theta} J_{\mu}^{\theta^T} \cdot (\widehat{\theta} - \theta) - \frac{1}{2} \|\widehat{\theta} - \theta\|_p \sum_i L_{\nabla J}^i |\widehat{\theta}_i - \theta_i|.
 \end{aligned}$$

□

In Sect. 2.3, we have seen that the policy parameters are updated as follows

$$\theta_{t+1} = \theta_t + \Delta \theta_t,$$

where, in the steepest ascent approaches, $\Delta \theta_t = \alpha_t \nabla_{\theta} J_{\mu}^{\theta_t}$ and α_t is a parameter that determines the step-size length. Recall that the steepest ascent direction of $J_{\mu}^{\theta_t}$ is defined as the vector $\Delta \theta_t$ that maximizes $J_{\mu}^{\theta_t + \Delta \theta_t}$ under the constraint that the change in the parameters ($\|\Delta \theta_t\|_p$) is sufficiently small.

In the following we describe three approaches to determine the step size exploiting the Lipschitz continuity of the policy gradient.

5.1 Single step size from single Lipschitz constant (SSS–SLC)

The scenario where a single Lipschitz constant is available for the gradient has been widely studied in the optimization literature. As suggested by Armijo (1966), if policy gradient is $L_{\nabla J}$ -LC and the L_2 -norm is used as metric, fixing the step size to the reciprocal value of the Lipschitz constant³:

$$\alpha = \frac{1}{L_{\nabla J}},$$

guarantees to improve at each iteration. Furthermore, when the function is convex, it allows to converge to an ϵ -optimal solution in $O(\frac{1}{\epsilon})$ iterations. It can be easily shown that similar results hold when the Lipschitz constant is replaced with its pointwise version ($\alpha_t = L_{\nabla J}(\theta_t)^{-1}$), with the advantage of inducing larger step sizes—being $L_{\nabla J} = \sup_{\theta} L_{\nabla J}(\theta)$ —with better performance improvements. Such result is simply obtained by the maximization w.r.t. the step size of the following quadratic lower bound to the performance improvement derived

³ See (Armijo 1966) for conditions under which convergence occurs and a proof of convergence.

from Lemma 8 using the L_2 -norm as metric:

$$\begin{aligned} J_{\mu}^{\theta_{t+1}} - J_{\mu}^{\theta_t} &\geq \nabla_{\theta} J_{\mu}^{\theta_t^T} \cdot \Delta\theta_t - \frac{L_{\nabla J}(\theta_t)}{2} \|\Delta\theta_t\|_2^2 \\ &= \alpha_t \left\| \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_2^2 - \alpha_t^2 \frac{L_{\nabla J}(\theta_t)}{2} \left\| \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_2^2. \end{aligned}$$

Using the step size α_t that maximizes the above lower bound, we are guaranteed that, at each iteration, the policy improvement is larger than $\frac{\left\| \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_2^2}{2L_{\nabla J}(\theta_t)}$. As expected, the smaller are the Lipschitz constants of the MDP and the policy model, the smaller is the Lipschitz constant of the policy gradient, and the larger are the step size and the guaranteed improvement at each iteration.

5.2 Single step size from multiple Lipschitz constants (SSS–MLC)

When a Lipschitz constant for each gradient component is available, it is convenient to exploit such information. By maximizing the bound in Lemma 8:

$$\begin{aligned} J_{\mu}^{\theta_{t+1}} - J_{\mu}^{\theta_t} &\geq \nabla_{\theta} J_{\mu}^{\theta_t^T} \cdot \Delta\theta_t - \frac{1}{2} \|\Delta\theta_t\|_2 L_{\nabla J}(\theta_t)^T \cdot |\Delta\theta_t| \\ &= \alpha_t \left\| \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_2^2 - \frac{1}{2} \alpha_t^2 \left\| \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_2 L_{\nabla J}(\theta_t)^T \cdot \left| \nabla_{\theta} J_{\mu}^{\theta_t} \right|, \end{aligned}$$

the following step size is obtained:

$$\alpha_t = \frac{\left\| \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_2}{L_{\nabla J}(\theta_t)^T \cdot \left| \nabla_{\theta} J_{\mu}^{\theta_t} \right|}.$$

Notice that, differently from the single Lipschitz constant case, the step size directly exploits the information of the gradient at time t . In this case, the policy performance improves at least by $\frac{\left\| \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_2^3}{2L_{\nabla J}(\theta_t)^T \cdot \left| \nabla_{\theta} J_{\mu}^{\theta_t} \right|}$, that is never worse than the improvement that is guaranteed with the SSS–SLC approach.

5.3 Multiple step sizes from multiple Lipschitz constants (MSS–MLC)

The two previous approaches update the parameters along the steepest ascent direction. However, to mitigate the main drawbacks of the steepest ascent method, many gradient approaches follow different directions (e.g. Amari and Douglas 1998). In particular, when a Lipschitz constant for each gradient component is available, it is interesting to consider a different step size for each parameter: $\Delta\theta_t^i = \alpha_t^i \nabla_{\theta} J_{\mu}^{\theta_t}$. As a result, the new candidate solution θ_{t+1} may lie outside the policy gradient direction $\nabla_{\theta} J_{\mu}^{\theta_t}$. All the step sizes can be obtained by maximizing the following concave⁴ function w.r.t. the step sizes $[\alpha_t^1, \dots, \alpha_t^d]$:

$$J_{\mu}^{\theta_{t+1}} - J_{\mu}^{\theta_t} \geq \nabla_{\theta} J_{\mu}^{\theta_t^T} \cdot \mathbf{A}_t \cdot \nabla_{\theta} J_{\mu}^{\theta_t} - \frac{1}{2} \left\| \mathbf{A}_t \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_p L_{\nabla J}(\theta_t)^T \cdot \mathbf{A}_t \cdot \left| \nabla_{\theta} J_{\mu}^{\theta_t} \right|, \quad (22)$$

where $\mathbf{A}_t = \text{diag}(\alpha_t^1, \dots, \alpha_t^d)$ and $\alpha_t^i \geq 0, \forall i \in \{1, \dots, d\}$.

⁴ The reader may refer to “Proof of Concavity of Function” section of Appendix for the proof.

Algorithm 1 MSS-MLC Algorithm

```

function MSS-MLC( $\nabla_{\theta} J_{\mu}, \mathbf{L}_{\nabla J}$ )
  initialize:  $V = 0$ 
   $\phi = \arccos\left(\frac{\mathbf{L}_{\nabla J}^{\top} |\nabla_{\theta} J_{\mu}|}{\|\mathbf{L}_{\nabla J}\|_2 \|\nabla_{\theta} J_{\mu}\|_2}\right)$ 
  if  $\phi \leq \arctan(\sqrt{1/8})$  then
     $\psi = \arctan\left(\frac{1 - \sqrt{1 - 8 \tan(\phi)^2}}{4 \tan(\phi)}\right)$ 
     $V = \frac{\|\nabla_{\theta} J_{\mu}\|_2^2 \cos(\psi)^2}{2\|\mathbf{L}_{\nabla J}\|_2 \cos(\psi + \phi)}$ 
    Create vector  $\Delta\theta = \frac{|\nabla_{\theta} J_{\mu}| \cos(\psi)}{\|\mathbf{L}_{\nabla J}\|_2 \cos(\psi + \phi)}$ 
    Rotate  $\Delta\theta$  of angle  $\phi$  around bivector  $\mathbf{L}_{\nabla J} \wedge |\nabla_{\theta} J_{\mu}|$ 
  end if
   $i^* = \arg \min_i \frac{|\nabla_{\theta} J_{\mu}^i|}{\mathbf{L}_{\nabla J}^i}$ 
  Build  $\widetilde{\nabla_{\theta} J_{\mu}}$  and  $\widetilde{\mathbf{L}_{\nabla J}}$  by removing component  $i^*$ 
   $[\widetilde{V}, \widetilde{\alpha}] = \text{MSS-MLC}(\widetilde{\nabla_{\theta} J_{\mu}}, \widetilde{\mathbf{L}_{\nabla J}})$ 
  if  $\widetilde{V} > V$  then
    return  $\widetilde{V}$  and  $[\widetilde{\alpha}^1, \dots, \widetilde{\alpha}^{i^*-1}, 0, \widetilde{\alpha}^{i^*}, \dots, \widetilde{\alpha}^d]$ 
  else
    return  $V$  and  $\left[\frac{\Delta\theta^1}{|\nabla_{\theta} J_{\mu}^1|}, \frac{\Delta\theta^2}{|\nabla_{\theta} J_{\mu}^2|}, \dots, \frac{\Delta\theta^d}{|\nabla_{\theta} J_{\mu}^d|}\right]$ 
  end if
end function

```

Being the performance-improvement bound a concave function in the step size, we could resort to one of the standard tools in convex optimization to find the optimal learning rates. Nonetheless, we propose a more efficient algorithm (reported in Algorithm 1) that optimizes the performance improvement bound by taking the best among a set of d closed-form candidate solutions (the algorithm is presented as a recursive function for conciseness, but an iterative implementation is straightforward). The computational complexity of MSS-MLC is $\mathcal{O}(d^3)$.

To explain how the algorithm works we need to rewrite the performance improvement bound as follows:

$$\begin{aligned}
 J_{\mu}^{\theta_{t+1}} - J_{\mu}^{\theta_t} &\geq \nabla_{\theta} J_{\mu}^{\theta_t \top} \Delta\theta_t - \frac{1}{2} \|\Delta\theta_t\|_2 \mathbf{L}_{\nabla J}(\theta_t)^{\top} |\Delta\theta_t| \\
 &= \left\| \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_2 \|\Delta\theta_t\|_2 \cos(\beta_{\nabla J_{\mu}}) - \frac{1}{2} \|\mathbf{L}_{\nabla J}(\theta_t)\|_2 \|\Delta\theta_t\|_2^2 \cos(\beta_L).
 \end{aligned}$$

To compute the length of the optimal step size $\|\Delta\theta_t\|_2$ w.r.t. to the above bound, we can compute the derivative and put it equal to zero:

$$\|\Delta\theta_t\|_2 = \frac{\left\| \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_2 \cos(\beta_{\nabla J_{\mu}})}{\|\mathbf{L}_{\nabla J}(\theta_t)\|_2 \cos(\beta_L)},$$

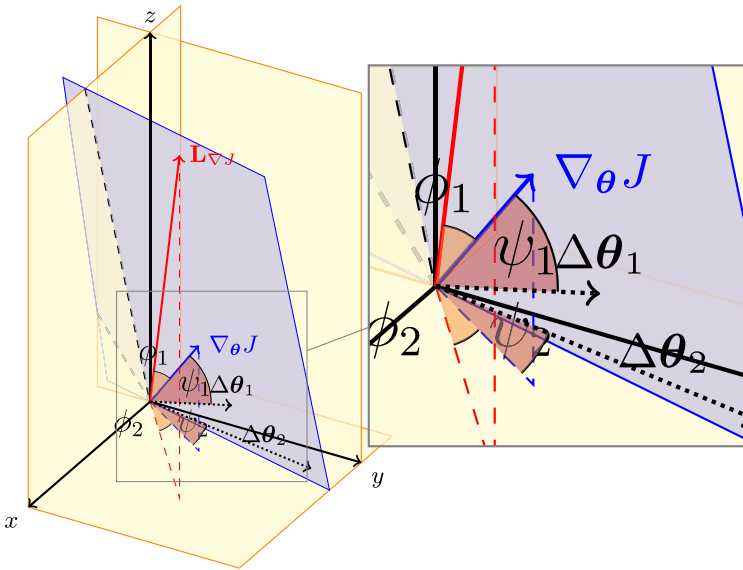


Fig. 1 Geometrical interpretation of the MSS-MLC algorithm

that leads to the following bound:

$$J_{\mu}^{\theta_{t+1}} - J_{\mu}^{\theta_t} \geq \frac{\|\nabla_{\theta} J_{\mu}^{\theta_t}\|_2^2 \cos(\beta_{\nabla J_{\mu}})^2}{2 \|\mathbf{L}_{\nabla J}(\theta_t)\|_2 \cos(\beta_L)}.$$

To maximize the performance improvement we need to choose the direction of $\Delta \theta_t$ so that it is as close as possible to $\nabla_{\theta} J_{\mu}^{\theta_t}$ (to maximize the value of $\cos(\beta_{\nabla J_{\mu}})$) and as far as possible to the vector $\mathbf{L}_{\nabla J}(\theta_t)$ (so as to minimize $\cos(\beta_L)$). It can be easily shown that the combination of these two desiderata puts $\Delta \theta_t$ on the plane identified by the vectors $\nabla_{\theta} J_{\mu}^{\theta_t}$ and $\mathbf{L}_{\nabla J}(\theta_t)$ closer to vector $\nabla_{\theta} J_{\mu}^{\theta_t}$ (refer to Fig. 1 for a geometrical interpretation of the algorithm in a three-dimensional problem). It follows that, calling ϕ the angle between vector $\nabla_{\theta} J_{\mu}^{\theta_t}$ and $\mathbf{L}_{\nabla J}(\theta_t)$, and renaming $\beta_{\nabla J_{\mu}}$ as ψ , we can rewrite the bound on the performance improvement as:

$$J_{\mu}^{\theta_{t+1}} - J_{\mu}^{\theta_t} \geq \frac{\|\nabla_{\theta} J_{\mu}^{\theta_t}\|_2^2 \cos(\psi)^2}{2 \|\mathbf{L}_{\nabla J}(\theta_t)\|_2 \cos(\psi + \phi)},$$

from which we can compute the optimal value for ψ : if $\phi \leq \arctan(1/8)$ the function presents a local maximum at $\psi = \arctan\left(\frac{1 - \sqrt{1 - 8 \tan(\phi)^2}}{4 \tan(\phi)}\right)$ associated to a guaranteed performance

improvement equal to $\frac{\|\nabla_{\theta} J_{\mu}^{\theta_t}\|_2^2 \cos(\psi)^2}{2 \|\mathbf{L}_{\nabla J}(\theta_t)\|_2 \cos(\psi + \phi)}$, followed by a local minimum and then increasing with ψ . Otherwise, if $\phi > \arctan(1/8)$, the bound is monotonically increasing with ψ . However, since we have the constraint that the step sizes must be positive, the angle ψ cannot increase arbitrarily, and must stop at the intersection with the planes identified by the Cartesian axes. So we need to explore what happens on the boundaries of the optimization problem (i.e., when one of the step sizes is put to zero). To do this we have to identify which

step size is the first to reach zero as ψ increases, that is the parameter with the smallest $\frac{|\nabla_{\theta} J_{\mu}^i|}{L_{\nabla J}^i}$ ratio. Once identified such parameter, the guaranteed performance improvement is saved, the parameter is removed from the vectors $\nabla_{\theta} J_{\mu}^{\theta_t}$ and $L_{\nabla J}(\theta_t)$, and the same procedure is repeated until the vector size is reduced to one. Finally, the optimal solution is the one associated to the largest guaranteed performance improvement among the local maxima found during the execution of the algorithm. Once the optimal $\Delta\theta$ has been identified, the i -th step-size parameter is easily computed through the following ratio: $\frac{\Delta\theta^i}{|\nabla_{\theta} J_{\mu}^i|}$. Notice that, although the resulting update may not follow the steepest ascent direction, since $\alpha_i^i \geq 0$, the algorithm has the same convergence guarantees as the standard policy gradient algorithm.

6 Numerical results

In this section, we empirically evaluate the performance of the three different strategies for choosing the step size introduced in the previous section. The performance is compared against a recently proposed algorithm (Pirotta et al. 2013), that we name Adaptive Gradient Step-Size and will be denoted by *AGSS*. As far as we know, it is the most related algorithm available in RL literature that investigates the choice of the step size. Pirotta et al. (2013) determine the step size by maximizing a lower bound to the expected performance gain. In the case of Gaussian policy model, the bound is a second-order polynomial of the step-size parameter that can be easily maximized in both exact and approximate scenarios. In the exact case, the optimal step size, that is, the value that maximizes a simplified version of the lower bound, is given by:

$$\alpha_{AGSS}^* = \frac{(1 - \gamma)^3 \sqrt{2\pi} \sigma^3 \|\nabla_{\theta} J_{\mu}^{\theta}\|_2^2}{\left(\gamma \sqrt{2\pi} \sigma + 2(1 - \gamma)|\mathcal{A}|\right) RM_{\phi}^2 \|\nabla_{\theta} J_{\mu}^{\theta}\|_1^2}, \tag{23}$$

where M_{ϕ} is the maximum absolute value of the basis functions, i.e., $\forall s \in \mathcal{S}, \forall i |\phi_i(s)| \leq M_{\phi}$. This step size method ensures a monotonic policy performance improvement. Moreover, Pirotta et al. (2013) have empirically shown that *AGSS* is more robust than standard step-size selection mechanisms (e.g., constant or time varying) to changes of the MDP parameters. To make the comparison fair, in the following experiments we use a Gaussian policy model as in (Pirotta et al. 2013). This choice is compatible with the Assumptions 2 and 3.

Before showing numerical results about the learning performance of the proposed algorithms, we can state some considerations about computational times. While single step-size algorithms have a very low computational complexity, *MSS-MLC* has a complexity that is cubic in the number of policy parameters d . The computational gap between single and multi-step algorithms increases as the size of the problem increases. Nonetheless the complexity is dominated by the sampling procedure. In Table 1, we compare the per-iteration times needed to compute the step size using *SSS-SLC* and *MSS-MLC* in the *Mass-Spring-Damper* domain using different parameterizations. Tests are performed on a standard laptop using a C++ implementation. We employed a state-dependent linear policy parametrization in which we have changed the degree of the polynomial. As shown by the experiments the complexity of *MSS-MLC* is higher than the one of *SSS-SLC* but it is dominated by the data collection phase, which, in turn, is 4 order of magnitude less than the time that would be required in the real-world scenario. In summary, the time spent for the computation of the step size is negligible w.r.t. the time needed for data collection.

Table 1 Time Complexity. Given the gradient and the Lipschitz constants, the time required by SSS–SLC and MSS–MLC to compute the step size are reported in table. The number of policy parameters have been changed by considering different degrees of the polynomial: 1, 5 and 10, respectively. In addition, the time required to collect the samples is reported both in simulated and real-world scenario

Parameters	SSS–SLC	MSS–MLC	Data collection time	
			Simulated	Real
3	$1.2 \cdot 10^{-6}s$	$9.4 \cdot 10^{-5}s$	$4.4 \cdot 10^{-1}s$	$2.5 \cdot 10^3s$
21	$1.6 \cdot 10^{-6}s$	$2.2 \cdot 10^{-3}s$		
66	$1.9 \cdot 10^{-6}s$	$2.0 \cdot 10^{-2}s$		

6.1 Linear-quadratic gaussian regulator

The first scenario is a discrete-time linear-quadratic Gaussian (LQG) regulator as described by Peters and Schaal (2008b). The LQG problem is defined by the following dynamics:

$$s_{t+1} = A(s_t + a_t); \quad a_t \sim \mathcal{N}(\theta \cdot s_t, \sigma^2); \quad r_t = -B(s_t^2 + a_t^2),$$

where s_t and a_t are scalars, and B is 0.5. The range of the state space is bounded to the interval $[-2, 2]$ and the initial state is drawn uniformly at random in the same interval. The scenario is particularly instructive since it allows to easily compute all the Lipschitz constants.

The objective of this test is to show how changes in the model parameters affect the algorithm performance. In particular, we analyze the impact of changes in the transition model (i.e., A), discount factor, and standard deviation of the Gaussian policy. While the values of the discount factor and of the standard deviation of the policy are exploited by all the algorithms, the change in the transition model (that influences the Lipschitz constant $L_{\mathcal{P}}$) is exploited only by Lipschitz approaches. Notice that, since there is a single policy parameter ($d = 1$), the proposed algorithms lead to the same result. For what concerns $AGSS$, in this case its step size is constant over iterations since the L_2 -norm equals the L_1 -norm for scalar values. It is also easy to modify the Lipschitz constant $L_{\mathcal{R}}$ associated to the reward model by changing the value of B . However, any change of such value does not alter the algorithm performance since it is associated to an equal change in the gradient magnitude.

Table 2 shows how the changes in the parameters influence the number of iterations required to learn a near-optimal value of the policy parameter. The variability of the results comes from the estimation of the gradient through GPOMDP algorithm. As $L_{\mathcal{P}}$ decreases the model becomes more and more smooth, so that the proposed approach has a significant advantage w.r.t. $AGSS$. For what concerns the standard deviation of the policy, algorithm performances degrade as the policy gets more deterministic. Note that the step length are inversely related to M_{ϕ}^i that grows as the policy becomes more deterministic (see Assumption 3). Lipschitz approach is less influenced by changes in the policy than $AGSS$ (it depends quadratically on M_{ϕ}) as shown by the increased performance gap. This is a positive effect due to the choice of considering the Kantorovich distance instead of simpler metrics between distributions, like the total variation one.

Although the Lipschitz approach outperforms the $AGSS$ algorithm in most of the settings, it cannot be applied to domains with large values of γ and $L_{\mathcal{P}}$ since the MDP needs to be a contraction w.r.t. the Kantorovich distance (i.e., $\gamma L_{\mathcal{P}}(1 + L_{\pi\theta}) < 1$).

Table 2 LQR domain. Table reports the number of iterations occurred before reaching a 0.01 approximation of the optimal policy parameter. Results are averaged over 10 runs

σ	γ	$L_{\mathcal{P}} = 0.3$		$L_{\mathcal{P}} = 0.5$		$L_{\mathcal{P}} = 0.7$	
		LIP	AGSS	LIP	AGSS	LIP	AGSS
0.50	0.5	896 ± 18	1932 ± 25	1430 ± 15	2637 ± 12	1724 ± 12	2565 ± 16
0.50	0.7	3098 ± 36	7204 ± 52	4442 ± 18	8367 ± 27	5809 ± 19	6988 ± 22
0.50	0.9	25676 ± 2853	114164 ± 207	36294 ± 68	103176 ± 119	65605 ± 86	66381 ± 50
0.75	0.5	368 ± 7	500 ± 56	515 ± 4	686 ± 7	636 ± 5	677 ± 6
0.75	0.7	1188 ± 8	2131 ± 24	1435 ± 7	2068 ± 8	1770 ± 5	1633 ± 4
0.75	0.9	9643 ± 45	33195 ± 47	9281 ± 13	23089 ± 25	14167 ± 9	13086 ± 13
1.00	0.5	171 ± 4	199 ± 4	208 ± 2	219 ± 2	296 ± 2	250 ± 2
1.00	0.7	502 ± 5	761 ± 6	529 ± 2	640 ± 2	729 ± 2	578 ± 1
1.00	0.9	3546 ± 14	11273 ± 21	2788 ± 5	6409 ± 6	4652 ± 3	4216 ± 3

Bold values are used to denote the best algorithm in each experimental setting

6.2 Mass-spring-damper

The second scenario is a simple mechanical system described by a mass, a linear spring and a damper (Ammar and Taylor 2012). The objective is to reach a desired state s_g by controlling the external force applied to the system. The system dynamic is described by the following equation:

$$m\ddot{x} + c\dot{x} + kx = F,$$

where m is the mass, c is the viscous friction coefficient ($\frac{N \cdot s}{m}$), k is the spring constant ($\frac{N}{m}$) and F is the external force. Let $\omega_0 = \sqrt{\frac{k}{m}}$ and $\zeta = \frac{c}{2\sqrt{km}}$ be the natural frequency and the damping ratio, respectively. Notice that, since in the following $0 < \zeta < 1$, the free system ($F = 0$) is under-damped and will oscillate with a frequency equal to $\omega_d = \omega_0\sqrt{1 - \zeta^2}$.

The continuous state space ($\mathcal{S} \in \mathbb{R}^2$) is defined by the position and velocity of the mass ($s = [x, \dot{x}]^T$) and the continuous scalar action a defines the external force applied to the system ($F = b \cdot a$). The control decision is performed every 0.1s, with the goal to bring and keep the mass to the desired state $s_g = [0.5, 0]^T$. At every time step the agent receives a reward that is proportional to the distance from the goal position x_g and to the control action:

$$r_{t+1} = -w_d |x_t - x_g| - w_a |b \cdot a_t|,$$

where, in all the experiments, we set b , w_d , w_a to 10, 1 and 1/20, respectively. Discount factor γ is set to 0.9. Note that the Lipschitz constant associated to the reward function $L_{\mathcal{R}}$ is 1.

Since the problem is continuous we exploited a Gaussian policy model

$$\pi^\theta(a|s) = \mathcal{N}(\phi(s)^T \theta, \sigma), \tag{24}$$

where $\theta \in \mathbb{R}^d$ and $\phi : \mathcal{S} \rightarrow \mathbb{R}^d$ are the basis functions. The standard deviation σ has been fixed to 0.5 in all the experiments. Three polynomial basis functions have been used to approximate the policy mean action: $\phi(s) = [1, x, \dot{x}]^T$. The initial setting is $x_0 = Unif(\{-1, 1\})$, $\dot{x}_0 = 0$ and $\theta_0 = \mathbf{0}$. The learning task was performed using GPOMDP with optimal baseline (Peters and Schaal 2008b), exploiting 500 episodes by 50 steps for each gradient estimate.

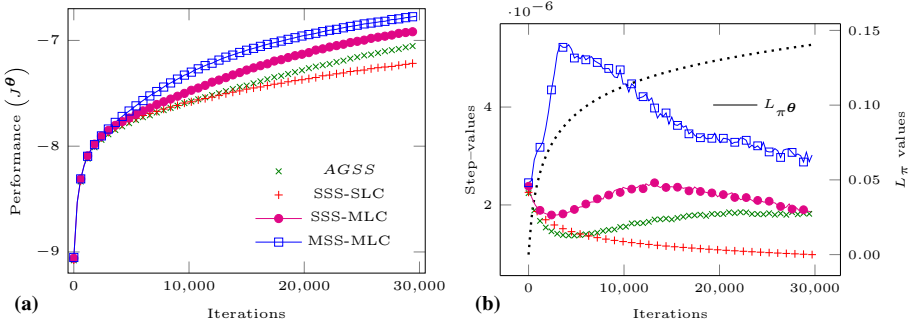


Fig. 2 Mass-Spring-Damper domain with configuration c_1 . **a** The score J^θ_μ as a function of iterations. **b** The step sizes over iterations. For MSS-MLC algorithm, the L_2 -norm of the vector α is drawn. Results are averaged over 100 runs. For sake of readability, the graphs do not display the (very small) confidence intervals

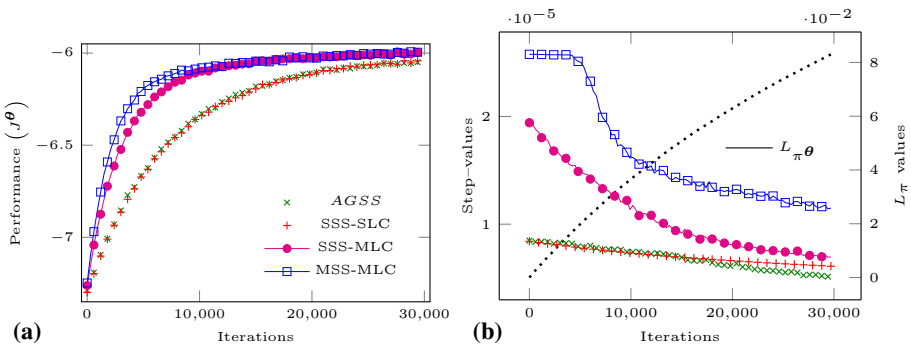


Fig. 3 Algorithm behaviours for the Mass-Spring-Damper domain with configuration c_2 . The same structure of Fig. 2 is used here

As done before, the objective is to show how the changes in the parameters affect the algorithm behaviour. In this scenario we concentrate the attention on the Lipschitz continuity of the model by changing the $L_{\mathcal{P}}$ constant. We defined two configurations: $c_1 = \{m = 0.5, k = 15, c = 0.001\}$ and $c_2 = \{m = 0.1, k = 15, c = 5\}$ with $L_{\mathcal{P}}$ -constants equal to 0.79 and 0.91, respectively. In both the configurations the position, velocity and action were limited in $[-1, 1]$, $[-2, 2]$ and $[-3, 3]$, respectively.

Figures 2a and 3a show that the knowledge of the individual Lipschitz constants can be successfully exploited in order to speed up the learning process. In particular, the algorithms that exploit the individual Lipschitz constants (SSS-MLC and MSS-MLC) outperform both SSS-SLC and AGSS algorithms in all the experiments because they are able to exploit larger learning step for each gradient component. Notice that in the c_1 configuration AGSS learns faster than SSS-SLC, while in the c_2 configuration, characterized by a smaller $L_{\mathcal{P}}$ value, the performances of the two algorithms get similar. As expected, MSS-MLC outperforms SSS-MLC, thanks to the possibility of updating the policy parameters outside the gradient direction. In particular, the advantage of MSS-MLC on SSS-MLC increases as the angle between the gradient direction and the vector of the Lipschitz constants gets larger. When such angle is zero, the two algorithms lead to the same result.

6.3 Ship steering with water current

In this scenario we adapted a standard control problem where the task is to steer a ship, which is cruising at constant speed, to a goal in minimum time (Rosenstein and Barto 2004). The task is made not trivial by the presence of different water currents around the goal. The continuous state and action spaces are described by the 2-dimensional ship position ($\mathcal{S} \in \mathbb{R}^2$) and the scalar heading ($\mathcal{A} \in \mathbb{R}$), respectively. The following equations describes the continuous motion of the ship:

$$\dot{x} = C (\cos (\omega + \omega_0) - y), \quad \dot{y} = C \sin (\omega + \omega_0) \omega$$

where $s = [x, y]^T$ denotes the state, C is the ship velocity and ω is the ship heading. Notice that the water current influences only the horizontal coordinate x , and its magnitude is proportional to the vertical position y . We changed the problem in order to minimize the travelled distance, by defining the immediate reward as the minimum distance from the goal region:

$$r_{t+1} = \min(g_r - \|s_t - s_g\|_2, 0),$$

where $s_g = [0, 0]^T$ is the goal position and $g_r = 0.2Km$ is the goal radius. The effects of a real ship inertia and resistance to water are not modelled, i.e., there is no time lag between changes in the desired heading and the actual one. Control decisions were taken every $25s$, with an immediate change of the ship heading. The ship starts at $s_0 = [2.5, -1]^T$.

The policy model is a standard Gaussian model described in Eq. (24) with fixed standard deviation $\sigma = 0.6$. Since the optimal policies for the objectives are not linear in the state variable, a Gaussian radial basis approximation was used:

$$\phi(s) = \left\{ \mathcal{N}(s; c_i, \Sigma_i) \right\}_{i=1}^3.$$

The shape and position of the Gaussian basis are defined by the following parameters:

$$\begin{aligned} c_1 &= [2.5, 0]^T, & \Sigma_1 &= \text{diag}(10, 5) \\ c_2 &= [1.25, 1]^T, & \Sigma_2 &= \text{diag}(3, 1) \\ c_3 &= [0, 0]^T, & \Sigma_3 &= \text{diag}(10, 5) \end{aligned}$$

State and action variables are bounded: $x \in [-2, 6]$, $y \in [-2, 2]$ and $a \in [-\pi, \pi]$. As done in the previous domains, we exploited GPOMDP with baseline with 500 episodes for the gradient estimate. Due to the low discount factor ($\gamma = 0.6$), we performed only 30 steps for each episode.

We developed two different scenarios that differ for the reference angle ω_0 : $\pi/2$ and 0, respectively. Since the ship heading points to north when $\omega_0 = \pi/2$ and $\omega = 0$, the learning process is easier in the former configuration than in the latter one. Notice that the initial policy, that is the same in both the configurations, is defined by zero-weights, that is, in the first iteration $\omega \sim \mathcal{N}(0, 0.6)$. In contrast, in the second configuration the agent must learn to steer of 90° , moving the ship heading from east to north, but this may be critical under Lipschitz assumptions. We will deeply discuss the second scenario later in the section.

Concerning the first configuration ($\omega_0 = \pi/2$), we have tested two different ship velocities ($C_1 = 0.01Km/s$ and $C_2 = 0.005Km/s$) that lead to $L_{\mathcal{P}}$ constants of 1.13 and 1.06. The Lipschitz constant $L_{\mathcal{R}}$ is 1. As the reader may notice from the $L_{\mathcal{P}}$ constants, the transition model is no more a contraction as in the previous domains. While the overall performance of Lipschitz algorithms does not seem to be affected by the change in the model transitions,

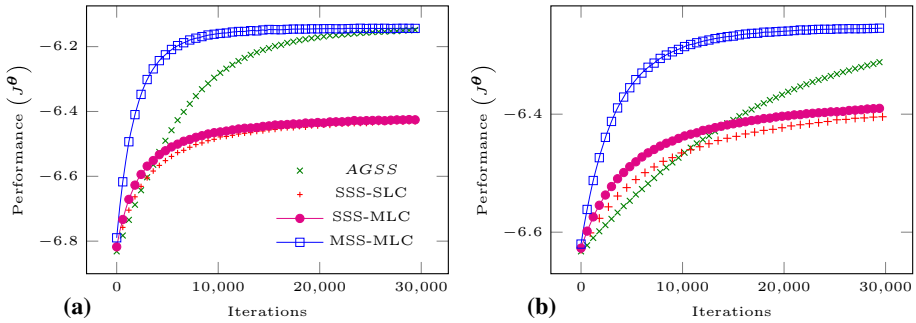


Fig. 4 Score J_μ^θ as a function of iterations. *Error bars* represent the samples standard deviation. The underlying domain consists in the ship steering with reference angle $\omega_0 = \pi/2$ and ship speed of $0.01Km/s$ (a) and $0.005Km/s$ (b), respectively

the *AGSS* algorithm is no more able to mimic the performance of the best algorithms. This is a direct consequence of the smaller gradient estimates due to slow dynamics. Moreover, as shown in Fig. 4b, the *AGSS* algorithm is negatively affected by changes in the transition model. In particular, the performance gap between the best algorithms and *AGSS* increases as the model becomes smoother because the latter approach is not able to exploit such information.

Concerning the second scenario ($\omega_0 = 0$), the *AGSS* algorithm reaches a better performance than all the Lipschitz algorithms when $C = 0.01Km/s$, see Fig. 5a. As mentioned before, the main issue is the high coefficient required on the first component of the basis functions in order to move the shift heading from east to north in the initial phase of the trial. Although the final policy is quite smooth (the change in the action value is slow and continuous), in the initial phase of the learning process, the first component of the policy is the term that guarantees the highest improvement, thus, it is characterized by large gradient values and parameter steps, that lead to high rate of increase of the Lipschitz constant. As the parameter increases, the policy becomes less Lipschitz (i.e., the L_π increases) and the system becomes less contractive, as shown in Fig. 5b. As a consequence, the learning step decreases over iterations at the same rate of increase of the Lipschitz constants.

However, when the ship moves at slow speed, the Lipschitz algorithm are able to exploit the improved smoothness of the system through larger step sizes. On the other hand, the *AGSS* algorithm exhibits a surprisingly slow learning behavior. The reasons appear clear comparing Figs. 5b, d. The improved Lipschitz continuity of the model directly influences the computation of the step size, leading to larger movements in the initial phase of the learning. At the same time, the rate of increase of the Lipschitz constant of the policy is slower than in the previous setting. Notice that, no changes are visible in the trend of the step sizes shown by the *AGSS* algorithm that is almost constant over iterations.

7 Conclusions

In this paper we have studied how to automatically set the step-size parameter in policy gradient algorithms under assumptions on the Lipschitz continuity of the MDP and the policy model. We have shown that under such continuity assumptions both the expected return and the policy gradient are Lipschitz in the policy-parameter space and we have derived Lipschitz constants for each component of the gradient. On the basis of such constants, we have

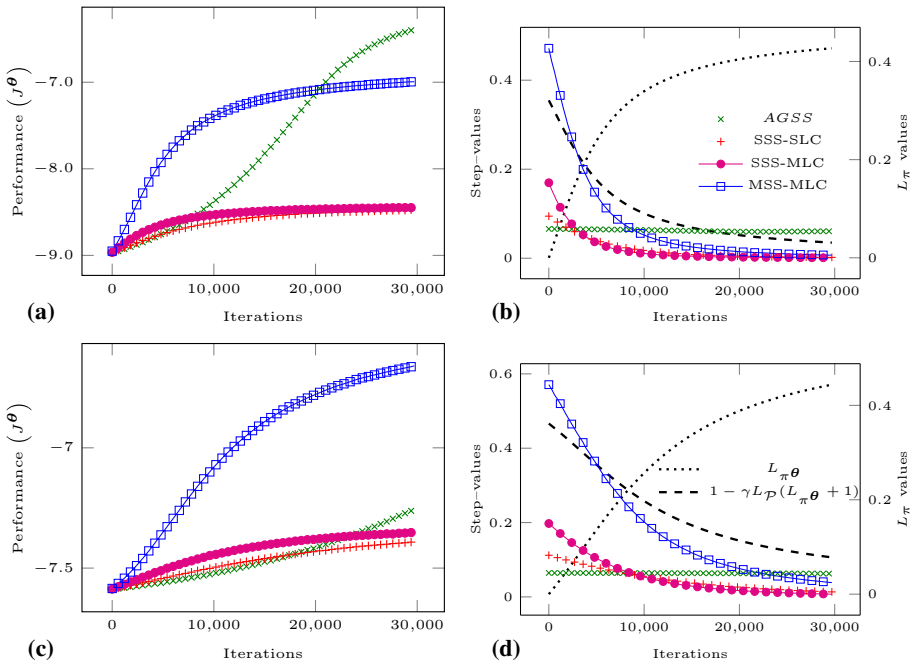


Fig. 5 Algorithm behaviors in the ship domain with $\omega_0 = 0$. While the upper figures refer to the first configuration ($C = 0.01 Km/s$), lower figures report the results for $C = 0.005 Km/s$. **a, c** The performance J^θ_μ as function of the iterations. **b, d** The step sizes (L_2 -norm of the vector for MSS-MLC algorithm) and changes in the Lipschitz continuity properties over iterations. In particular, the L_π constant and the contractive property of the system ($1 - \gamma L_{\mathcal{P}}(1 + L_{\pi}\theta)$) are shown

proposed to update the policy parameters according to three different learning strategies that guarantee a performance improvement at each iteration. As shown by empirical evaluations, when the MDP has strong continuity properties (i.e., small Lipschitz constants), the proposed approach can take advantage of this by making larger step sizes than the ones made by related approaches that do not exploit such information.

The main drawback of the proposed approach is that it can be applied only when the state transition dynamics of the MDP are a contraction (i.e., $\gamma L_{\mathcal{P}}(1 + L_{\pi}\theta) < 1$). To overcome such limitation, it would be interesting to study the effect of replacing the Kantorovich distance with other metrics between distributions. Another way to remove the limitation consists in combining the proposed approach with the AGSS (Pirota et al. 2013) (e.g., by simply taking at each iteration the step size that guarantees the largest improvement). Finally, future research could address how to use the proposed approach in problems where the Lipschitz constants are unknown and need to be estimated from data.

Appendices

Proof of Proposition 1

Given an $L_{\mathcal{R}}$ -LC MDP, for any pair of stationary policies corresponding to parameters θ and $\hat{\theta}$, the absolute difference between the performance of policy π^θ and policy $\pi^{\hat{\theta}}$ can be bounded as follows:

$$\begin{aligned}
 \left| J_{\mu}^{\theta} - \widehat{J}_{\mu}^{\widehat{\theta}} \right| &= \frac{1}{1 - \gamma} \left| \mathbb{E}_{(s,a) \sim \zeta_{\mu}^{\theta}} [\mathcal{R}(s, a)] - \mathbb{E}_{(s,a) \sim \zeta_{\mu}^{\widehat{\theta}}} [\mathcal{R}(s, a)] \right| \\
 &= \frac{1}{1 - \gamma} \left| \iint_{\mathcal{S}\mathcal{A}} \left(\zeta_{\mu}^{\theta}(s, a) - \zeta_{\mu}^{\widehat{\theta}}(s, a) \right) \mathcal{R}(s, a) \, d\text{ads} \right| \\
 &\leq \frac{L_{\mathcal{R}}}{1 - \gamma} \sup_f \left\{ \left| \iint_{\mathcal{S}\mathcal{A}} f \, d \left(\zeta_{\mu}^{\theta} - \zeta_{\mu}^{\widehat{\theta}} \right) \right| : \|f\|_L \leq 1 \right\} \\
 &= \frac{L_{\mathcal{R}}}{1 - \gamma} \mathcal{K} \left(\zeta_{\mu}^{\theta}, \zeta_{\mu}^{\widehat{\theta}} \right)
 \end{aligned}$$

□

Proof of Lemma 5

Under Assumptions 1, 2 and 3, the i -th component of η^{θ} is $L_{\eta^{\theta}}^i$ -LC w.r.t. the state-action space, that is: $\forall (s, \widehat{s}, a, \widehat{a}) \in \mathcal{S}^2 \times \mathcal{A}^2$,

$$\begin{aligned}
 \left| \eta_i^{\theta}(s, a) - \eta_i^{\theta}(\widehat{s}, \widehat{a}) \right| &= \left| \nabla_{\theta_i} \log \pi^{\theta}(a|s) Q^{\theta}(s, a) - \nabla_{\theta_i} \log \pi^{\theta}(\widehat{a}|\widehat{s}) Q^{\theta}(\widehat{s}, \widehat{a}) \right| \\
 &= \left| \left(\nabla_{\theta_i} \log \pi^{\theta}(a|s) - \nabla_{\theta_i} \log \pi^{\theta}(\widehat{a}|\widehat{s}) \right) Q^{\theta}(s, a) \right. \\
 &\quad \left. + \left(Q^{\theta}(s, a) - Q^{\theta}(\widehat{s}, \widehat{a}) \right) \nabla_{\theta_i} \log \pi^{\theta}(\widehat{a}|\widehat{s}) \right| \\
 &\leq \sup_{s', a'} \left| Q^{\theta}(s', a') \right| \left| \nabla_{\theta_i} \log \pi^{\theta}(a|s) - \nabla_{\theta_i} \log \pi^{\theta}(\widehat{a}|\widehat{s}) \right| \\
 &\quad + \sup_{s', a'} \left| \nabla_{\theta_i} \log \pi^{\theta}(a'|s') \right| \left| Q^{\theta}(s, a) - Q^{\theta}(\widehat{s}, \widehat{a}) \right| \\
 &\leq \left(\frac{R}{1 - \gamma} L_{\nabla \log \pi^{\theta}}^i + M_{\theta}^i L_{Q^{\theta}} \right) d_{\mathcal{S}\mathcal{A}}((s, a), (\widehat{s}, \widehat{a})).
 \end{aligned}$$

□

Proof of Lemma 6

For any pair of stationary policies corresponding to θ and $\widehat{\theta}$, the absolute difference of the i -th component of functions η^{θ} and $\eta^{\widehat{\theta}}$ is upper bounded by: $\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \forall (\theta, \widehat{\theta}) \in \Theta^2$,

$$\begin{aligned}
 \left| \eta_i^{\theta}(s, a) - \widehat{\eta}_i^{\widehat{\theta}}(s, a) \right| &= \left| \nabla_{\theta_i} \log \pi^{\theta}(a|s) Q^{\theta}(s, a) - \nabla_{\theta_i} \log \pi^{\widehat{\theta}}(a|s) \widehat{Q}^{\widehat{\theta}}(s, a) \right| \\
 &= \left| \left(\nabla_{\theta_i} \log \pi^{\theta}(a|s) - \nabla_{\theta_i} \log \pi^{\widehat{\theta}}(a|s) \right) Q^{\theta}(s, a) \right. \\
 &\quad \left. + \left(Q^{\theta}(s, a) - \widehat{Q}^{\widehat{\theta}}(s, a) \right) \nabla_{\theta_i} \log \pi^{\theta}(a|s) \right| \\
 &\leq \sup_{s', a'} \left| Q^{\theta}(s', a') \right| \left| \nabla_{\theta_i} \log \pi^{\theta}(a|s) - \nabla_{\theta_i} \log \pi^{\widehat{\theta}}(a|s) \right| \\
 &\quad + \sup_{s', a'} \left| \nabla_{\theta_i} \log \pi^{\theta}(a'|s') \right| \left| Q^{\theta}(s, a) - \widehat{Q}^{\widehat{\theta}}(s, a) \right|
 \end{aligned}$$

Proof follows from the application of Assumption 3 and from the bound on the supremum norm of the Q -function. □

Proof of Lemma 7

Under Assumptions 1, 2 and 3, the i -th component of η is upper-bounded by: $\forall(s, a) \in S \times \mathcal{A}, \forall(\theta, \hat{\theta}) \in \Theta^2,$

$$\begin{aligned} \left| \eta_i^\theta(s, a) - \hat{\eta}_i^\theta(s, a) \right| &\leq \frac{R}{1-\gamma} \left| \nabla_{\theta_i} \log \pi^\theta(a|s) - \nabla_{\theta_i} \log \pi^{\hat{\theta}}(a|s) \right| + M_\theta^i \left| Q^\theta(s, a) - Q^{\hat{\theta}}(s, a) \right| \\ &\leq \left(\frac{R}{1-\gamma} L_{\nabla \log \pi}^i(\theta) + M_\theta^i L_Q(\theta) \right) d_\Theta(\theta, \hat{\theta}), \end{aligned}$$

where inequalities follow from the combination of Lemma 6 with Theorem 2. □

Proof of Theorem 3

Under Assumptions 1, 2, and 3, the i -th component of the gradient $\nabla_\theta J$ of the expected return is $L_{\nabla J}^i(\theta)$ -PLC, that is: $\forall(\theta, \hat{\theta}) \in \Theta^2,$

$$\begin{aligned} \left| \nabla_{\theta_i} J_\mu^\theta - \nabla_{\theta_i} J_\mu^{\hat{\theta}} \right| &\leq \frac{1}{1-\gamma} \left| \mathbb{E}_{(s,a) \sim \zeta_\mu^\theta} \left[\eta_i^\theta(s, a) \right] - \mathbb{E}_{(s,a) \sim \zeta_\mu^{\hat{\theta}}} \left[\eta_i^\theta(s, a) \right] \right| \\ &\quad + \frac{1}{1-\gamma} \left| \mathbb{E}_{(s,a) \sim \zeta_\mu^{\hat{\theta}}} \left[\eta_i^\theta(s, a) - \hat{\eta}_i^\theta(s, a) \right] \right| \\ &= \frac{L_{\eta^\theta}^i}{1-\gamma} \left| \iint_{S \times \mathcal{A}} \frac{\eta_i^\theta}{L_{\eta^\theta}^i} d(\zeta_\mu^\theta - \zeta_\mu^{\hat{\theta}}) \right| + \frac{1}{1-\gamma} \left| \mathbb{E}_{(s,a) \sim \zeta_\mu^{\hat{\theta}}} \left[\eta_i^\theta(s, a) - \hat{\eta}_i^\theta(s, a) \right] \right| \\ &\leq \frac{L_{\eta^\theta}^i}{1-\gamma} \mathcal{K}(\zeta_\mu^\theta, \zeta_\mu^{\hat{\theta}}) + \frac{L_\eta^i(\theta)}{1-\gamma} d_\Theta(\theta, \hat{\theta}) \\ &\leq \frac{1}{1-\gamma} \left(L_{\eta^\theta}^i L_\zeta(\theta) + L_\eta^i(\theta) \right) d_\Theta(\theta, \hat{\theta}). \end{aligned}$$

The proof follows by combining Theorem 5 and Theorem 7 with Proposition 2. □

Proof of concavity of function (22)

For sake of readability we report here Eq. (22):

$$\Delta J_\mu^\theta(\mathbf{A}) = J_\mu^{\theta_{t+1}} - J_\mu^{\theta_t} \geq \nabla_\theta J_\mu^{\theta_t \top} \mathbf{A} \nabla_\theta J_\mu^{\theta_t} - \frac{1}{2} \left\| \mathbf{A} \nabla_\theta J_\mu^{\theta_t} \right\|_p \mathbf{L}_{\nabla J}(\theta_t)^\top \mathbf{A} \left| \nabla_\theta J_\mu^{\theta_t} \right|.$$

We need to prove that $\Delta J_\mu^\theta(\mathbf{A})$ is a concave function:

$$\begin{aligned} \Delta J_\mu^\theta(\lambda \mathbf{A} + (1-\lambda)\mathbf{B}) &= \nabla_\theta J_\mu^{\theta_t \top} (\lambda \mathbf{A} + (1-\lambda)\mathbf{B}) \nabla_\theta J_\mu^{\theta_t} \\ &\quad - \frac{1}{2} \left\| (\lambda \mathbf{A} + (1-\lambda)\mathbf{B}) \nabla_\theta J_\mu^{\theta_t} \right\|_p \mathbf{L}_{\nabla J}(\theta_t)^\top (\lambda \mathbf{A} + (1-\lambda)\mathbf{B}) \left| \nabla_\theta J_\mu^{\theta_t} \right| \\ &\geq \nabla_\theta J_\mu^{\theta_t \top} (\lambda \mathbf{A} + (1-\lambda)\mathbf{B}) \nabla_\theta J_\mu^{\theta_t} \\ &\quad - \frac{1}{2} \left(\lambda \left\| \mathbf{A} \nabla_\theta J_\mu^{\theta_t} \right\|_p + (1-\lambda) \left\| \mathbf{B} \nabla_\theta J_\mu^{\theta_t} \right\|_p \right) \\ &\quad \mathbf{L}_{\nabla J}(\theta_t)^\top (\lambda \mathbf{A} + (1-\lambda)\mathbf{B}) \left| \nabla_\theta J_\mu^{\theta_t} \right| \end{aligned}$$

$$\begin{aligned}
 &= \lambda \left(\nabla_{\theta} J_{\mu}^{\theta_t \top} \mathbf{A} \nabla_{\theta} J_{\mu}^{\theta_t} - \frac{\lambda}{2} \left\| \mathbf{A} \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_p \mathbf{L}_{\nabla J}(\theta_t)^{\top} \mathbf{A} \left| \nabla_{\theta} J_{\mu}^{\theta_t} \right| \right) \\
 &+ (1 - \lambda) \left(\nabla_{\theta} J_{\mu}^{\theta_t \top} \mathbf{B} \nabla_{\theta} J_{\mu}^{\theta_t} - \frac{1 - \lambda}{2} \left\| \mathbf{B} \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_p \mathbf{L}_{\nabla J}(\theta_t)^{\top} \mathbf{B} \left| \nabla_{\theta} J_{\mu}^{\theta_t} \right| \right) - \frac{1}{2} \lambda \\
 &(1 - \lambda) \left(\left\| \mathbf{A} \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_p \mathbf{L}_{\nabla J}(\theta_t)^{\top} \mathbf{B} \left| \nabla_{\theta} J_{\mu}^{\theta_t} \right| + \left\| \mathbf{B} \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_p \mathbf{L}_{\nabla J}(\theta_t)^{\top} \mathbf{A} \left| \nabla_{\theta} J_{\mu}^{\theta_t} \right| \right) \\
 &\geq \lambda \left(\nabla_{\theta} J_{\mu}^{\theta_t \top} \mathbf{A} \nabla_{\theta} J_{\mu}^{\theta_t} - \frac{\lambda}{2} \left\| \mathbf{A} \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_p \mathbf{L}_{\nabla J}(\theta_t)^{\top} \mathbf{A} \left| \nabla_{\theta} J_{\mu}^{\theta_t} \right| \right) \\
 &+ (1 - \lambda) \left(\nabla_{\theta} J_{\mu}^{\theta_t \top} \mathbf{B} \nabla_{\theta} J_{\mu}^{\theta_t} - \frac{1 - \lambda}{2} \left\| \mathbf{B} \nabla_{\theta} J_{\mu}^{\theta_t} \right\|_p \mathbf{L}_{\nabla J}(\theta_t)^{\top} \mathbf{B} \left| \nabla_{\theta} J_{\mu}^{\theta_t} \right| \right) \\
 &\geq \lambda \Delta J_{\mu}^{\theta}(\mathbf{A}) + (1 - \lambda) \Delta J_{\mu}^{\theta}(\mathbf{B}) \tag{25}
 \end{aligned}$$

where (25) follows from the Minkowski inequality. □

References

Amari, S., & Douglas, S. (1998). Why natural gradient? In: Acoustics, Speech and Signal Processing, 1998. Proceedings of the 1998 IEEE international conference on, vol 2, pp. 1213–1216 vol. 2, doi:[10.1109/ICASSP.1998.675489](https://doi.org/10.1109/ICASSP.1998.675489).

Ammar, H., & Taylor, M. (2012). Reinforcement learning transfer via common subspaces. *Adaptive and learning agents, lecture notes in computer science* (Vol. 7113, pp. 21–36). Berlin: Springer.

Armijo, L. (1966). Minimization of functions having lipschitz continuous first partial derivatives. *Pacific Journal of Mathematics*, 16(1), 1–3.

Baxter, J., & Bartlett, P. L. (2001). Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research*, 15, 319–350.

Bertsekas, D. P., & Shreve, S. E. (1978). *Stochastic optimal control: The discrete time case* (Vol. 139). New York: Academic Press.

Deisenroth, M. P., Neumann, G., Peters, J., et al. (2013). A survey on policy search for robotics. *Foundations and Trends in Robotics*, 2(1–2), 1–142.

Ferns, N., Panangaden, P., & Precup, D. (2005). Metrics for markov decision processes with infinite state spaces. *Proceedings of the twenty-first conference annual conference on uncertainty in artificial intelligence (UAI-05)* (pp. 201–208). Arlington, Virginia: AUAI Press.

Grondman, I., Busoniu, L., Lopes, G. A., & Babuska, R. (2012). A survey of actor-critic reinforcement learning: Standard and natural policy gradients. *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, 42(6), 1291–1307.

Hinderer, K. (2005). Lipschitz continuity of value functions in markovian decision processes. *Mathematical Methods of Operations Research*, 62(1), 3–22.

Kakade, S. (2001). A natural policy gradient. In *Advances in neural information processing systems 14* (Vol. 14, pp. 1531–1538). Vancouver, British Columbia: MIT Press.

Kober, J., & Peters, J. (2008). Policy search for motor primitives in robotics. In *Advances in neural information processing systems 21* (Vol. 21, pp. 849–856). Vancouver, British Columbia: Curran Associates, Inc.

Moré, J. J., & Thuente, D. J. (1994). Line search algorithms with guaranteed sufficient decrease. *ACM Transactions on Mathematical Software*, 20(3), 286–307.

Peters, J., & Schaal, S. (2006). Policy gradient methods for robotics. In: Intelligent robots and systems, 2006 IEEE/RSJ international conference on, (pp. 2219–2225).

Peters, J., & Schaal, S. (2008a). Natural actor-critic. *Neurocomputing*, 71(7–9), 1180–1190.

Peters, J., & Schaal, S. (2008b). Reinforcement learning of motor skills with policy gradients. *Neural Networks*, 21(4), 682–697.

Pirotta, M., Restelli, M., & Bascetta, L. (2013). Adaptive step-size for policy gradient methods. *Advances in Neural Information Processing Systems*, 26, 1394–1402.

Puterman, M. L. (1994). *Markov decision processes: Discrete stochastic dynamic programming*. New York, NY: Wiley.

Rachelson, E., & Lagoudakis, M.G. (2010). On the locality of action domination in sequential decision making. In: International symposium on artificial intelligence and mathematics.

- Robbins, H., & Monro, S. (1951). A stochastic approximation method. *The Annals of Mathematical Statistics*, 22(3), 400–407.
- Rosenstein, M. T., & Barto, A. G. (2004). Supervised actor-critic reinforcement learning. In J. Si, A. Barto, W. Powell, & D. Wunsch (Eds.), *Handbook of learning and approximate dynamic programming* (pp. 359–380). John Wiley & Sons, Inc.
- Spall, J. C. (1992). Multivariate stochastic approximation using a simultaneous perturbation gradient approximation. *Automatic Control, IEEE Transactions on*, 37(3), 332–341.
- Sutton, R.S., McAllester, D.A., Singh, S.P., & Mansour, Y. (1999). Policy gradient methods for reinforcement learning with function approximation. In: *Advances in neural information processing systems 12*, [NIPS Conference, Denver, Colorado, USA, November 29 - December 4, 1999], (pp. 1057–1063).
- Vlassis, N., Toussaint, M., Kontes, G., & Piperidis, S. (2009). Learning model-free robot control by a monte carlo em algorithm. *Autonomous Robots*, 27(2), 123–130.
- Wagner, P. (2011). A reinterpretation of the policy oscillation phenomenon in approximate policy iteration. *Advances in Neural Information Processing Systems*, 24, 2573–2581.
- Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8, 229–256.