# Towards Faster Global Convergence of Robust Policy Gradient Methods

**Anonymous authors**
Paper under double-blind review

## Abstract

We establish the global convergence of the policy gradient method for robust Markov Decision Processes (MDPs) under the assumption that the robust return is smooth with respect to a policy. Despite restrictive, such smoothness assumption is satisfied in many interesting settings such as reward-robust MDPs. We also obtain iteration complexity comparable to non-robust MDPs that is significantly faster than existing rates for robust MDPs.

## 1 Introduction

Robust MDPs model decision making problems in which the environment is uncertain or partially known (Nilim & Ghaoui, 2005; Iyengar, 2005a). The robust optimal policies are not only resilient to the changes in the environmental parameters, but also better generalizable (Mannor et al., 2007; Xu & Mannor, 2010). Most of the work on the topic are value-based methods (Lutter et al., 2021; Pinto et al., 2017; Ho et al., 2020; Badrinath & Kalathil, 2021; Behzadian et al., 2021; Goyal & Grand-Clément, 2018; Mannor et al., 2016; Wang et al., 2022c;a; Wang & Zou, 2021; Kumar et al., 2022), and only few are policy-based methods (Wang & Zou, 2022b; Kumar et al., 2023; Wang et al., 2022b). Moreover, up to our knowledge, there are only two works that establish global convergence of policy gradient methods: with iteration complexity of $O(SA\epsilon^{-3})$ for R-contamination uncertainty sets (Wang & Zou, 2022a) and $O(S^4A^2\epsilon^{-4})$ for general kernel uncertainty sets Wang et al. (2022b), where $S$ and $A$ are the number of states and the number of actions respectively, and $\epsilon$ is the target difference between optimal robust return and our achieved robust return.

The state-of-the-art iteration complexity for global convergence of policy gradient for non-robust MDPs is $O(SA\epsilon^{-1})$, however the proofs are quite technical and complicated (Agarwal et al., 2020; Mei et al., 2020; Xiao, 2022). Robustness requirement adds another layer of complexity to the problem, leading to a significant degradation in the iteration complexity (Wang & Zou, 2022a; Wang et al., 2022b).

We simplify and extend the machinery for non-robust MDPs (Xiao, 2022) to the robust MDPs, and achieve the same iteration complexity as for non-robust MDPs, i.e., $O(SA\epsilon^{-1})$. This complexity is achieved under the assumption that the robust return is $L$-smooth which holds in many useful scenarios (Gadot et al., 2023) (see Appendix B.2 for more examples).

Table 1: Iteration Complexity for Global Convergence of Policy Gradient Methods

| Robust MDPs | Complexity | Remark |
|---|---|---|
| Non-Robust | $O(SA\epsilon^{-1})$ | Xiao (2022) |
| $(s,a)$ rectangular R-Contamination robust MDPs | $O(SA\epsilon^{-3})$ | Wang & Zou (2022a) |
| General $L$-smooth robust MDPs | $O(SL\epsilon^{-1})$ | Ours |
| General kernel robust MDPs | $O(S^4A^2\epsilon^{-4})$ | Wang et al. (2022b) |

## 2 Main

A robust Markov decision process (MDP) is a tuple $(\mathcal{S}, \mathcal{A}, \gamma, \mu, \mathcal{U})$ such that $\mathcal{S}$ and $\mathcal{A}$ are finite state and action spaces respectively, $\gamma \in [0,1)$ is a discount factor and $\mu \in \Delta_{\mathcal{S}}$ the initial state

distribution. Denoting $\mathcal{X} := \mathcal{S} \times \mathcal{A}$, the couple $(P, R) \in \mathcal{U}$ corresponds to the MDP model with $P : \mathcal{X} \to \Delta_{\mathcal{S}}$ being a transition kernel and $R : \mathcal{X} \to \mathbb{R}$ a reward function. A policy $\pi : \mathcal{S} \to \Delta_{\mathcal{A}}$ maps each state to a probability distribution over $\mathcal{A}$, and we denote by $\Pi$ the set of possible policies. For any policy $\pi \in \Pi$, $R^{\pi} \in \mathbb{R}^{\mathcal{S}}$ is the expected immediate reward defined by $R^{\pi}(s) := \langle \pi_s, R(s, \cdot) \rangle_{\mathcal{A}}, \quad \forall s \in \mathcal{S}$, where $\pi_s$ is a shorthand for $\pi(\cdot|s)$. We similarly define the stochastic matrix induced by $\pi$ as $P^{\pi}(s'|s) := \langle \pi_s, P(s'|s, \cdot) \rangle_{\mathcal{A}}, \quad \forall s, s' \in \mathcal{S}$. Our goal is to maximize the robust discounted return $\min_{\pi \in \Pi} \rho_{\mathcal{U}}^{\pi}$ over the set of policies $\Pi$, with

$$\rho_{\mathcal{U}}^{\pi} := \min_{(P,R) \in \mathcal{U}} \rho_{(P,R)}^{\pi}, \quad \text{where return } \rho_{(P,R)}^{\pi} := \mu^T (I - \gamma P^{\pi})^{-1} R^{\pi}. \tag{1}$$

However, the above optimization is proven to be Strongly NP-Hard for general convex uncertainty set (Wolfram Wiesemann, 2012). Hereafter, we assume the uncertainty set $\mathcal{U}$ to be convex, that allows robust return $\rho_{\mathcal{U}}^{\pi}$ to be differentiable and $L$-smooth w.r.t. policy $\pi$ with $L > 0$, that is, $\left| \rho_{\mathcal{U}}^{\pi'} - \rho_{\mathcal{U}}^{\pi'} - \langle \nabla_{\pi} \rho_{\mathcal{U}}^{\pi}, \pi' - \pi \rangle \right| \leq \frac{L}{2} \|\pi' - \pi\|^2, \ \forall \pi', \pi \in \Pi$. In particular, the $L_p$-ball reward uncertainty set guarantees smoothness of the robust return (Gadot et al., 2023), more cases are discussed in Appendix B.2. To optimize the robust return, we rely on the projected gradient ascent: $\pi_{k+1} := \mathbf{proj}_{\Pi}(\pi_k + \frac{1}{L} \nabla_{\pi} \rho_{\mathcal{U}}^{\pi_k})$, where $\frac{1}{L}$ is the learning rate and $\mathbf{proj}_{\Pi}$ denotes the orthogonal projection onto set $\Pi$. For this policy gradient procedure, we get the following convergence result:

**Theorem 1 (Global optimality).** *For all iterations $k \geq 1$, it holds that:*

$$\rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^{\pi_k} \leq \frac{8LC_{\text{PL}}^2 \mathbf{diam}(\Pi)^2 (\rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^{\pi_0})}{k},$$

*where $C_{\text{PL}} \leq \frac{1}{\min_s \mu(s)}$ is a problem dependent constant, that we define properly in Appendix C.3.*

*Proof.* For the full proof see Appendix C. Below, we briefly describe the proof sketch. The main challenge of the proof is to extend the Gradient Domination property to the robust problem (sometimes called PL condition). This condition was proven to hold for non-robust MDPs (Xiao, 2022). Despite it might not hold for general robust problems, we prove it for smooth robust problems:

**Lemma 1 (Gradient Domination lemma).** *For any policy $\pi \in \Pi$, its sub-optimality is bounded by its policy gradient as $\rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^{\pi} \leq C_{\text{PL}} \max_{\pi' \in \Pi} \langle \pi' - \pi, \nabla \rho_{\mathcal{U}}^{\pi} \rangle$.*

For the proof see Appendix C.3. The rest of the proof mainly follows Xiao (2022), but all results are carefully extended for the robust case. In particular, from combining Lemma 1 with the cohesive bond lemma (Appendix C.4), we get directly that $\rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^{\pi_{k+1}} \leq C_{PL} \max_{\pi' \in \Pi} \langle \nabla \rho_{\mathcal{U}}^{\pi_{k+1}}, \pi' - \pi_{k+1} \rangle \leq 2C_{PL} L \|\pi_{k+1} - \pi_k\| \mathbf{diam}(\Pi)$, that provides us the lower bound on $\|\pi_{k+1} - \pi_k\|$. In contrast, sufficient increase lemma (Appendix C.2) provides and upper bound on $\|\pi_{k+1} - \pi_k\|$: $\rho_{\mathcal{U}}^{\pi_{k+1}} - \rho_{\mathcal{U}}^{\pi_k} \geq \frac{L}{2} \|\pi_{k+1} - \pi_k\|^2, \forall k$. Combining these upper and lower bounds, we get the following sub-optimality recursion: $a_{k+1}^2 + a_{k+1} \leq a_k$, for $a_k := \frac{\rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^{\pi_k}}{8LC_{PL}^2 \mathbf{diam}(\Pi)^2}$ (Appendix C.5), which in turn implies the convergence. Indeed, note that, the sub-optimality recursion $a_{k+1} - a_k \leq -a_{k+1}^2$ corresponds to an ordinary differential equation $\frac{da}{dk} \leq -a^2$, whose solution is $a(k) \leq \frac{1}{k + \frac{1}{a(0)}} \leq \frac{1}{k}$.

This relation intuitively indicates an $O(\frac{1}{\epsilon})$ iteration complexity for achieving an $\epsilon$-optimal solution. □

## 3 DISCUSSION

We establish global convergence for robust MDPs with $L$-smooth robust return with respect to a policy, and proved its iteration complexity $O(SA\epsilon^{-1})$, which is much faster than existing state-of-the-art complexity $O(S^4 A^2 \epsilon^{-4})$ for robust MDPs (Wang et al., 2022b). Moreover, our proof trivially yields a simpler and more intuitive proof for non-robust MDPs by taking a single environment uncertainty set.

REFERENCES

Alekh Agarwal, Sham M. Kakade, Jason D. Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift, 2020.

Kishan Panaganti Badrinath and Dileep Kalathil. Robust reinforcement learning using least squares policy iteration with provable performance guarantees. In *International Conference on Machine Learning*, pp. 511–520. PMLR, 2021.

Bahram Behzadian, Marek Petrik, and Chin Pang Ho. Fast algorithms for $l_\infty$-constrained s-rectangular robust MDPs. *Advances in Neural Information Processing Systems*, 34:25982–25992, 2021.

Jalaj Bhandari and Daniel Russo. Global optimality guarantees for policy gradient methods, 2019. URL https://arxiv.org/abs/1906.01786.

Uri Gadot, Esther Derman, Navdeep Kumar, Maxence Mohamed Elfatihi, Kfir Levy, and Shie Mannor. Solving non-rectangular reward-robust mdps via frequency regularization, 2023.

Vineet Goyal and Julien Grand-Clément. Robust markov decision process: Beyond rectangularity, 2018. URL https://arxiv.org/abs/1811.00215.

Chin Pang Ho, Marek Petrik, and Wolfram Wiesemann. Partial policy iteration for l1-robust markov decision processes, 2020. URL https://arxiv.org/abs/2006.09484.

Garud N. Iyengar. Robust dynamic programming. *Mathematics of Operations Research*, 30(2): 257–280, May 2005a. ISSN 1526-5471. doi: 10.1287/moor.1040.0129. URL http://dx.doi.org/10.1287/MOOR.1040.0129.

Garud N Iyengar. Robust dynamic programming. *Mathematics of Operations Research*, 30(2): 257–280, 2005b.

Vijay Konda and John Tsitsiklis. Actor-critic algorithms. In S. Solla, T. Leen, and K. Müller (eds.), *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 1999. URL https://proceedings.neurips.cc/paper/1999/file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf.

Navdeep Kumar, Kfir Levy, Kaixin Wang, and Shie Mannor. Efficient policy iteration for robust markov decision processes via regularization, 2022. URL https://arxiv.org/abs/2205.14327.

Navdeep Kumar, Esther Derman, Matthieu Geist, Kfir Levy, and Shie Mannor. Policy gradient for s-rectangular robust markov decision processes, 2023. URL https://arxiv.org/abs/2301.13589.

Yan Li, Tuo Zhao, and Guanghui Lan. First-order policy optimization for robust markov decision process, 2022. URL https://arxiv.org/abs/2209.10579.

Michael Lutter, Shie Mannor, Jan Peters, Dieter Fox, and Animesh Garg. Robust value iteration for continuous control tasks. *Robotics: Science and Systems*, 2021.

Shie Mannor, Duncan Simester, Peng Sun, and John N Tsitsiklis. Bias and variance approximation in value function estimates. *Management Science*, 53(2):308–322, 2007.

Shie Mannor, Ofir Mebel, and Huan Xu. Robust mdps with k-rectangular uncertainty. *Math. Oper. Res.*, 41(4):1484–1509, nov 2016. ISSN 0364-765X.

Jincheng Mei, Chenjun Xiao, Csaba Szepesvari, and Dale Schuurmans. On the global convergence rates of softmax policy gradient methods. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 6820–6829. PMLR, 13–18 Jul 2020. URL https://proceedings.mlr.press/v119/mei20b.html.

Paul Milgrom and Ilya Segal. Envelope theorems for arbitrary choice sets. *Econometrica*, 70: 583–601, 02 2002. doi: 10.1111/1468-0262.00296.

Arnab Nilim and Laurent El Ghaoui. Robust control of markov decision processes with uncertain transition matrices. *Oper. Res.*, 53:780–798, 2005.

Lerrel Pinto, James Davidson, Rahul Sukthankar, and Abhinav Gupta. Robust adversarial reinforcement learning. In *International Conference on Machine Learning*, pp. 2817–2826. PMLR, 2017.

John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust region policy optimization, 2015. URL https://arxiv.org/abs/1502.05477.

Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018. URL http://incompleteideas.net/book/the-book-2nd.html.

Richard S Sutton, David A McAllester, Satinder P Singh, Yishay Mansour, et al. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems*, volume 99, pp. 1057–1063. Citeseer, 1999.

Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In S. Solla, T. Leen, and K. Müller (eds.), *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 2000. URL https://proceedings.neurips.cc/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf.

Kaixin Wang, Navdeep Kumar, Kuangqi Zhou, Bryan Hooi, Jiashi Feng, and Shie Mannor. The geometry of robust value functions. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 22727–22751. PMLR, 17–23 Jul 2022a. URL https://proceedings.mlr.press/v162/wang22k.html.

Qiuhao Wang, Chin Pang Ho, and Marek Petrik. On the convergence of policy gradient in robust mdps, 2022b. URL https://arxiv.org/abs/2212.10439.

Yue Wang and Shaofeng Zou. Online robust reinforcement learning with model uncertainty, 2021. URL https://arxiv.org/abs/2109.14523.

Yue Wang and Shaofeng Zou. Policy gradient method for robust reinforcement learning, 2022a.

Yue Wang and Shaofeng Zou. Policy gradient method for robust reinforcement learning. *International Conference on Machine Learning*, 162:23484–23526, 2022b.

Yue Wang, Fei Miao, and Shaofeng Zou. Robust constrained reinforcement learning. *arXiv preprint arXiv:2209.06866*, 2022c.

Berç Rustem Wolfram Wiesemann, Daniel Kuhn. Robust markov decision processes. *Mathematics of Operations Research 38(1):153-183*, 2012.

Lin Xiao. On the convergence rates of policy gradient methods, 2022. URL https://arxiv.org/abs/2201.07443.

Huan Xu and Shie Mannor. Robustness and generalization, 2010. URL https://arxiv.org/abs/1005.2243.

## A APPENDIX

### RELATED WORK

**Non-Robust MDPs**. Policy gradient is derived in Sutton et al. (2000) for non-robust MDPs which is widely used in practice with many variants (Schulman et al., 2015; Konda & Tsitsiklis, 1999; Sutton & Barto, 2018). Recently, there have been global convergences guarantees results Agarwal et al. (2020); Bhandari & Russo (2019) with an iteration complexity $O(1/\epsilon)$ for finding $\epsilon$-optimal policy Xiao (2022).

$(s, a)$**-rectangular R-Contamination Robust MDPs**. The paper Wang & Zou (2022a) derives policy gradient for R-rectangular robust MDPs complexity $O(S^2 A \log(\frac{1}{\epsilon}))$ similar to non-robust MDPs. Further, it establishes global convergence policy gradient with an iteration complexity $O(1/\epsilon^3)$ for finding $\epsilon$-optimal policy assuming oracle policy gradient.

**General $(s, a)$-rectangular Robust MDPs** The paper Li et al. (2022) establishes global convergence for robust mirror policy decent for $(s, a)$-rectangular robust MDPs in general with an iteration complexity $O(1/\epsilon)$ and $O(log(1/\epsilon))$ for finding $\epsilon$-optimal policy, with two increasing-stepsize schemes. However, it assumes the oracle access to policy gradient.

**General Robust MDPs** The paper Wang et al. (2022b) establishes global convergence for Double-Loop Robust Policy Gradient for general robust MDPs with an iteration complexity $O(1/\epsilon^4)$ for finding $\epsilon$-optimal policy, assuming the oracle access to policy gradient. Solving the policy gradient upto $\epsilon$ tolerance via value methods that takes $(s, a)$-rectangular and $s$-rectangular case with complexity of $O(S^4 A \log(1/\epsilon))$ and $O(S^4 A^3 \log(1/\epsilon))$ respectively using convex optimizations tools. Our techniques are completely different than this work.

## B  SMOOTHNESS

### B.1  BACKGROUND

We go over the some definitions and background results helpful for our discussion. Let value function and Q-value function be defined as follows:

$$v^{\pi}_{(P,R)} := (I - \gamma P^{\pi})^{-1}, \qquad Q^{\pi}_{(P,R)}(s,a) := R(s,a) + \gamma \sum_{s'} P(s'|s,a) v^{\pi}_{(P,R)}(s').$$

Following Kumar et al. (2023), let the occupation measure for policy $\pi$ with initial vector $u \in \mathbb{R}^{\mathcal{S}}$ be defined by

$$d^{\pi}_{P,u} := u^T (I - \gamma P^{\pi})^{-1}.$$

Let $P^{\pi}_{\mathcal{U}}, R^{\pi}_{\mathcal{U}}$ be the worst kernel and reward function under policy $\pi$, that is:

$$(P^{\pi}_{\mathcal{U}}, R^{\pi}_{\mathcal{U}}) \in \operatorname*{arg\,min}_{(P,R) \in \mathcal{U}} \rho^{\pi}_{(P,R)}.$$

Similarly to Kumar et al. (2023), we define the robust values as follows

$$d^{\pi}_{\mathcal{U}} := d^{\pi}_{P^{\pi}_{\mathcal{U}}}, \qquad v^{\pi}_{\mathcal{U}} = v^{\pi}_{(P^{\pi}_{\mathcal{U}}, R^{\pi}_{\mathcal{U}})}, \qquad Q^{\pi}_{\mathcal{U}} = Q^{\pi}_{(P^{\pi}_{\mathcal{U}}, R^{\pi}_{\mathcal{U}})}.$$

Let $\mathcal{R}^{\mathrm{s}}_p, \mathcal{P}^{\mathrm{s}}_p$ denote $\mathrm{s}$-rectangular $L_p$-bounded reward noise uncertainty set and $\mathrm{s}$-rectangular $L_p$-bounded kernel noise uncertainty set respectively, as they are decomposable over states, that is,

$$\mathcal{R}^{\mathrm{s}}_p = \otimes_{s \in \mathcal{S}} R_s, \quad \text{and} \quad \mathcal{P}^{\mathrm{s}}_p = \otimes_{s \in \mathcal{S}} P_s,$$

where $\|R_s\|_p \leq \alpha_s, \|P_s\|_p \leq \beta_s$ and $\sum_{s'} P_s(s'|a) = 0 \forall s, a$ (see Wolfram Wiesemann (2012)). Similarly, $\mathcal{R}^{\mathrm{sa}}, \mathcal{P}^{\mathrm{sa}}$ denotes $\mathrm{sa}$-rectangular reward uncertainty set and kernel respectively, as they are decomposable over states-actions, that is,

$$\mathcal{R}^{\mathrm{sa}}_p = \otimes_{s \in \mathcal{S}, a \in \mathcal{A}} R_{sa}, \quad \text{and} \quad \mathcal{P}^{\mathrm{sa}}_p = \otimes_{s \in \mathcal{S}, a \in \mathcal{A}} P_{sa},$$

where $\|R_{sa}\|_p \leq \alpha_{sa}, \|P_{sa}\|_p \leq \beta_{sa}$ and $\sum_{s'} P_{sa}(s') = 0 \forall s, a$ (see Iyengar (2005b); Nilim & Ghaoui (2005)).

Let $q$ be the Holder's conjugate of $p \in [1, \infty)$, that satisfying,

$$\frac{1}{p} + \frac{1}{q} = 1.$$

**Proposition 1.** *The non-robust return $\rho^{\pi}_{P,R}$ is $\frac{A}{(1-\gamma)^2}$-smooth for $\|R\|_{\infty} \leq 1$.*

*Proof.*  See Agarwal et al. (2020).  □

### B.2 SMOOTHNESS RESULTS

Now we are ready to state our first smoothness result of this section.

**Proposition 2.** *For the non-rectangular reward uncertainty set $\mathcal{U} = \{P_0\} \otimes \{R_0 + \mathcal{R}_p\}$, the robust return $\rho_{\mathcal{U}}^\pi$ is $2(q-1)\frac{(SA)^{\frac{q+1}{q}} A^2}{(1-\gamma)^4} + 2\gamma \frac{(SA)^{\frac{1}{q}} A}{(1-\gamma)^3}$-smooth, where $\mathcal{R}_p := \{R \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \mid \|R\|_p \leq \alpha\}$.*

*Proof.* Proved in Gadot et al. (2023). □

**Lemma 2.** *$\rho_{\mathcal{U}}^\pi$ is $\frac{\|R\|_\infty + \|\alpha\|_\infty}{(1-\gamma)^3}$-smooth for the uncertainty set $\mathcal{U} = \{P_0\} \times \{R_0 + \mathcal{R}_p^{sa}\}$, where $\|\alpha\|_\infty = \max_{s,a} \alpha_{sa}$.*

*Proof.* From Kumar et al. (2022), we have

$$\rho_{\mathcal{U}}^\pi = \rho_{(P_0,R_0)}^\pi - \sum_{s,a} d_{P_0}^\pi(s)\pi(a|s)\alpha_{sa}. \tag{2}$$

Now we know $\langle d_P^\pi, R^\pi \rangle$ is $\frac{\|R\|_\infty}{(1-\gamma)^3}$-smooth Agarwal et al. (2020). Hence the robust return $\rho_{\mathcal{U}}^\pi$ is $\frac{\|R\|_\infty + \|\alpha\|_\infty}{(1-\gamma)^3}$ - smooth. □

**Lemma 3.** *$\rho_{\mathcal{U}}^\pi$ is smooth for the uncertainty set $\mathcal{U} = \{P_0\} \times \{R_0 + \mathcal{R}_p^s\}$ for $p \in (1, \infty)$.*

*Proof.* From Kumar et al. (2022), we have

$$\rho_{\mathcal{U}}^\pi = \rho_{(P_0,R_0)}^\pi - \sum_s d_{P_0}^\pi(s)\|\pi_s\|_q \alpha_s. \tag{3}$$

It is clear that $\|\pi_s\|_q$ is smooth in $\pi$ for all $p \in (1, \infty)$. Furthermore, we know $d_{P_0}^\pi(s)$ is also smooth, as $d_{P_0}^\pi(s) = \rho_{P_0,R}^\pi$ where $R(s',a') = \mathbf{1}((s',a') = (s,a))$, and we know $\rho_{P,R}^\pi$ is smooth for all $P$ and all $R$ Agarwal et al. (2020). This establishes the smoothness of the robust return in our case, as it is algebraic function of smooth functions. □

We conclude that there are various setting under which the robust return is smooth, the our global convergence analysis for policy gradient methods applies.

## C GLOBAL CONVERGENCE

### C.1 ASSUMPTIONS

**Assumption 1.** *We assume the uncertainty set $\mathcal{U}$ is convex and compact.*

The above assumption is very mild that is satisfied in most of the settings.

**Assumption 2.** *[Smoothness] The function $\rho_{\mathcal{U}}^\pi$ is L-smooth function, that is*

$$\left| \rho_{\mathcal{U}}^{\pi'} - \rho_{\mathcal{U}}^\pi - \langle \nabla \rho_{\mathcal{U}}^\pi, \pi' - \pi \rangle \right| \leq \frac{L}{2}\|\pi' - \pi\|^2, \qquad \forall \pi', \pi \in \Pi. \tag{4}$$

The assumption doesn't hold for general uncertainty set, however it may hold for many useful uncertainty sets.

Policy udate is done via projected gradient ascent as:

$$\pi_{k+1} := \mathbf{proj}_\Pi(\pi_k + \frac{1}{L}\nabla_\pi \rho_{\mathcal{U}}^{\pi_k}), \tag{5}$$

where $\frac{1}{L}$ is the learning rate and $\mathbf{proj}_\Pi$ denotes the orthogonal projection onto set $\Pi$.

**Lemma 4.** *(Convex Projection Lemma) For any convex set $\mathbb{X} \subseteq \mathbb{R}^d$, any point $a \in \mathbb{X}$, and any update direction $u \in \mathbb{R}^d$, let $b = proj_{\mathbb{X}}(a + u)$ be the projection of $a + u$ onto $\mathbb{X}$, then we have*

1. $\langle u, b - a \rangle \geq \|b - a\|_2^2$.

2. $\langle c - b, u - (b - a) \rangle \leq 0, \qquad \forall c \in \mathcal{X}$.

*Proof.* Follows trivially from the geometry (see figure 1), and the fact that the hyperplane separates a convex set from a point not in the set. $\qquad\square$



Figure 1: Convex Projection

**Remark 1.** *For simplicity of the notation, we use $\rho_k$ for $\rho_{\mathcal{U}}^{\pi_k}$.*

### C.2 MONOTONE IMPROVEMENT OF THE ROBUST RETURN

Now present 'Sufficient Increase Lemma' ensures monotone improvement of the robust return, that uses only convexity of the projection set $\Pi$, smoothness of the robust return.

**Lemma 5.** *[Sufficient Increase Lemma] Gradient ascent ensures the monotone improvement in the robust return. Precisely,*

$$\rho_{k+1} - \rho_k \geq \frac{L}{2}\|\pi_{k+1} - \pi_k\|^2, \qquad \forall k.$$

*Proof.* From smoothness of the robust return, we have

$$\rho_{k+1} \geq \rho_k + \langle \nabla \rho_k, \pi_{k+1} - \pi_k \rangle - \frac{L}{2}\|\pi_{k+1} - \pi_k\|^2, \tag{6}$$

$$= \rho_k + L\langle \frac{1}{L}\nabla \rho_k, \pi_{k+1} - \pi_k \rangle - \frac{L}{2}\|\pi_{k+1} - \pi_k\|^2, \tag{7}$$

$$\geq \rho_k + L\|\pi_{k+1} - \pi_k\|^2 - \frac{L}{2}\|\pi_{k+1} - \pi_k\|^2, \qquad \text{(from Lemma 4)}. \tag{8}$$

This ends the proof. $\qquad\square$

Note that the above sufficient increase lemma does not require the robust return to be concave. Further, the above lemma is enough to ensure iterates $\{\rho_{\mathcal{U}}^{\pi_k}\}$ converge to some some value $\hat{\rho}$, as the iterates forms monotonically increasing sequence. However, it doesn't imply the $\hat{\rho}$ is global maxima or local maxima for that matter. This just implies, the iterates $\rho_{\mathcal{U}}^{\pi_k}$ keeps on increasing until the gradient $G(\pi_k)$ doesn't diminish to zero.

## C.3 GRADIENT DOMINATION

Hence, for the global optimality, we need second part, to ensure that the norm of the gradient vanishes only when the sub-optimality does.

**Lemma 6 (Gradient Domination lemma).** *For any policy $\pi \in \Pi$, its sub-optimality is bounded by its policy gradient as*

$$\rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^\pi \leq C_{\mathrm{PL}} \max_{\pi' \in \Pi} \langle \pi' - \pi, \nabla \rho_{\mathcal{U}}^\pi \rangle,$$

*where $C_{\mathrm{PL}} := \max_{(\pi,s) \in \Pi \times \mathcal{S}} \frac{d_{P_{\mathcal{U}}^\pi}^{\pi_{P_{\mathcal{U}}^\pi}^*}(s)}{d_{P_{\mathcal{U}}^\pi}^{\pi}(s)}.$*

*Proof.* Recall, $(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi) \in \arg\min_{(P,R) \in \mathcal{U}} \rho_{(P,R)}^\pi$ is the worst parameters w.r.t. policy $\pi$. And $\pi_{(P,R)}^* \in \arg\max_\pi \rho_{(P,R)}^\pi$ be the best policy for dynamics $(P, R)$. Then we have

$$\rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^\pi = \max_{\pi'} \min_{(P,R) \in \mathcal{U}} \rho_{(P,R)}^{\pi'} - \rho_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}^\pi, \qquad \text{(by definition)}, \tag{9}$$

$$\leq \max_{\pi'} \rho_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}^{\pi'} - \rho_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}^\pi, \qquad \text{(by definition of min-operator)}, \tag{10}$$

$$\text{(Using non-robust Performance Difference Lemma Agarwal et al. (2020), we get)} \tag{11}$$

$$= \sum_s d_{P_{\mathcal{U}}^\pi}^{\pi_{P_{\mathcal{U}}^\pi}^*}(s) \sum_a \left( \pi_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}^*(a|s) - \pi(a|s) \right) Q_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}^\pi(s, a), \tag{12}$$

$$\leq \sum_s d_{P_{\mathcal{U}}^\pi}^{\pi_{P_{\mathcal{U}}^\pi}^*}(s) \underbrace{\max_{\pi'_s} \sum_a \left( \pi'_s(a) - \pi(a|s) \right) Q_{\mathcal{U}}^\pi(s, a)}_{\geq 0} \tag{13}$$

$$= \sum_s \frac{d_{P_{\mathcal{U}}^\pi}^{\pi_{P_{\mathcal{U}}^\pi}^*}(s)}{d_{\mathcal{U}}^\pi(s)} d_{\mathcal{U}}^\pi(s) \underbrace{\max_{\pi'_s} \sum_a \left( \pi'_s(a) - \pi(a|s) \right) Q_{\mathcal{U}}^\pi(s, a)}_{\geq 0} \tag{14}$$

$$\leq \left( \max_s \frac{d_{P_{\mathcal{U}}^\pi}^{\pi_{P_{\mathcal{U}}^\pi}^*}(s)}{d_{\mathcal{U}}^\pi(s)} \right) \max_{\pi'} \sum_s d_{\mathcal{U}}^\pi(s) \sum_a \left( \pi'(a|s) - \pi(a|s) \right) Q_{\mathcal{U}}^\pi(s, a) \tag{15}$$

$$= \left( \max_s \frac{d_{P_{\mathcal{U}}^\pi}^{\pi_{P_{\mathcal{U}}^\pi}^*}(s)}{d_{\mathcal{U}}^\pi(s)} \right) \max_{\pi'} \left\langle \pi' - \pi, \nabla \rho_{\mathcal{U}}^\pi \right\rangle. \tag{16}$$

The last equality comes from the Envelope theorem Milgrom & Segal (2002) and policy gradient theorem Sutton et al. (1999),

$$\frac{\partial \rho_{\mathcal{U}}^\pi}{\partial \pi(a|s)} = \sum_s d_{P_{\mathcal{U}}^\pi}^\pi(s) Q_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}^\pi(s, a).$$

$\square$

## C.4 COHESIVE BOND

Now, we have both the parts: One that lower bounds the gradient and the other that upper bounds it. However, they are not exactly in very compatible forms, hence we require the result below that acts a cohesive bond between the two.

**Lemma 7.** *For all $\pi_k$, we have*

$$\langle \nabla \rho_{k+1}, \pi' - \pi_{k+1} \rangle \leq 2L \|\pi_{k+1} - \pi_k\| \mathbf{diam}(\Pi),$$

*where $\mathbf{diam}(C) := \max_{x,y \in C} \|x - y\|$ is the diameter of $C$.*

*Proof.* For all $x, y \in C$, we have:

$\langle \nabla \rho_{k+1}, \pi' - \pi_{k+1} \rangle$

$= \langle \nabla \rho_{k+1} - \nabla \rho_k + \nabla \rho_k, \pi' - \pi_{k+1} \rangle,$     [Subtract & add $\nabla \rho_k$]

$= \langle \nabla \rho_{k+1} - \nabla \rho_k, \pi' - \pi_{k+1} \rangle + \langle \nabla \rho_k, \pi' - \pi_{k+1} \rangle,$     [Linearity of scalar product]

$\leq \|\nabla \rho_{k+1} - \nabla \rho_k\| \|\pi' - \pi_{k+1}\| + \langle \nabla \rho_k, \pi' - \pi_{k+1} \rangle,$ Cauchy-Schwartz inequality]

$\leq L\|\pi_{k+1} - \pi_k\| \|\pi' - \pi_{k+1}\| + \langle \nabla \rho_k, \pi' - \pi_{k+1} \rangle,$     [Smoothness of robust return]

$= L\|\pi_{k+1} - \pi_k\| \|\pi' - \pi_{k+1}\| + L\langle \frac{1}{L} \nabla \rho_k - (\pi_{k+1} - \pi_k) + (\pi_{k+1} - \pi_k), \pi' - \pi_{k+1} \rangle$

$\leq 2L\|\pi_{k+1} - \pi_k\| \|\pi' - \pi_{k+1}\|, L\langle \frac{1}{L} \nabla \rho_k - (\pi_{k+1} - \pi_k), \pi' - \pi_{k+1} \rangle$     [Cauchy-Schwartz inequality]

$\leq 2L\|\pi_{k+1} - \pi_k\| \|\pi' - \pi_{k+1}\|,$     [From Lemma 4]

$\leq 2L\|\pi_{k+1} - \pi_k\| \mathbf{diam}(\Pi).$

$\square$

## C.5    PROOF OF GLOBAL CONVERGENCE

Now we are ready to prove our core result that is, the sub-optimality recursion.

**Lemma.** *Take $\eta = \frac{1}{L}$ as a learning rate. Then, the scaled sub-optimality $a_k = \frac{\rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^{\pi_k}}{8LC_{PL}^2 \mathbf{diam}(\Pi)^2}$ follows the recursion*

$$a_{k+1}^2 + a_{k+1} - a_k \leq 0.$$

*Proof.* From the PL condition proved in Lemma 1, we have

$$\rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^{\pi_{k+1}} \leq C_{PL} \max_{\pi'} \langle \pi' - \pi_{k+1}, \nabla \rho_{\mathcal{U}}^{\pi_{k+1}} \rangle \tag{17}$$

$$\leq 2L\|\pi_{k+1} - \pi_k\| \mathbf{diam}(\Pi), \quad \text{(from Lemma 7)} \tag{18}$$

$$\leq C_{\mathrm{PL}} \cdot 2\sqrt{2L(\rho_\pi^{\pi_{k+1}} - \rho_{\mathcal{U}}^{\pi_k})} \cdot \mathbf{diam}(\Pi), \quad \text{(from Lemma 5)} \tag{19}$$

Squaring both sides and adding subtracting $\rho^*$ in RHS, we get

$$\left( \rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^{\pi_{k+1}} \right)^2 \leq 8C_{\mathrm{PL}}^2 L \mathbf{diam}(\Pi)^2 \left( (\rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^{\pi_k}) + (\rho_{\mathcal{U}}^* - \rho_\pi^{\pi_{k+1}}) \right)$$

Setting $a_k := \frac{\rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^{\pi_k}}{8LC_{\mathrm{PL}}^2 \mathbf{diam}(\Pi)^2}$, the sequence $(a_k)_{k \in \mathbb{N}}$ satisfies the recursion $a_{k+1}^2 \leq a_k - a_{k+1}$. $\square$

The sub-optimality recursion derived in the theorem above, illustrates how the sub-optimality at time $k + 1$ depends at the sub-optimality at time $k$. Moreover, the sub-optimality recursion has the quadratic form and $a_k \geq 0$, hence its solution is given as

$$a_{k+1} \leq \sqrt{\frac{1}{4} + a_k} - \frac{1}{2}.$$

As a sanity check, we observe that $\sqrt{\frac{1}{4} + a} - \frac{1}{2} \leq a$ for all $a \geq 0$, implying that $(a_k)_{k \in \mathbb{N}}$ is monotonically decreasing. Further, 0 is the only non-negative fixed point of the $\sqrt{\frac{1}{4} + a} - \frac{1}{2} = a$ implying that $(a_k)_{k \in \mathbb{N}}$ monotonically decreases to 0.

Now, we investigate the convergence rate for $a_k$. Observe that if $a_0, a_k \gg 1$, then $a_{k+1} \approx \sqrt{a_k}$ and $a_k \approx (a_0)^{\frac{1}{2^k}}$. That is, the convergence rate is super-exponential! Yet, in most cases, $8LC_{\mathrm{PL}}^2 \mathbf{diam}(\Pi)^2 \gg 1$ and $\rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^{\pi_0} = 8LC_{\mathrm{PL}}^2 \mathbf{diam}(\Pi)^2 a_0 = O(1)$ is bounded so we are more interested in the case where $a_0 \ll 1$. In fact, in an MDP with a reward smaller than 1, we do have $\rho_{\mathcal{U}}^\pi = O(1)$.

In this regime, the sub-optimality recursion $a_{k+1} - a_k \leq -a_{k+1}^2$ suggests the ordinary differential equation $\frac{da}{dk} \leq -a^2$ whose solution is $a(k) \leq \frac{1}{k + \frac{1}{a(0)}} \leq \frac{1}{k}$. This intuitively indicates an $O(\frac{1}{\epsilon})$ iteration complexity for achieving an $\epsilon$-optimal solution, which we state below formally.

**Corollary** (**Global optimality**). *For all iterations $k \geq 1$, it holds that:*

$$\rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^{\pi_k} \leq \max\left(\frac{8LC_{\text{PL}}^2\mathbf{diam}(\Pi)^2}{k}, 2^{-\frac{k}{2}}\right)(\rho_{\mathcal{U}}^* - \rho_{\mathcal{U}}^{\pi_0}).$$

*Proof.* The sub-optimality recursion yields the desired result which follows directly from Xiao (2022). □

It can be seen that he exponential part is always lower than the other term. Hence, we can ignore the exponential term. Further, the diameter of the policy class $\Pi$, can be upper bounded as

$$\mathbf{diam}(\Pi)^2 = \max_{\pi,\pi}\sum_s \|\pi_s' - \pi_s\|_2^2 \leq \max_{\pi',\pi}\sum_s \|\pi_s' - \pi_s\|_1^2 \leq 4S.$$

This yields the desired result.

In addition, it can be noted that our result reduces to non-robust case, when we take uncertainty set $\mathcal{U}$ to be a singleton set. In this case, our proof is simplified and more-readable version of Agarwal et al. (2020); Xiao (2022).