

---

# Policy Gradient for Rectangular Robust Markov Decision Processes

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Policy gradient methods have become a standard for training reinforcement learning  
2 agents in a scalable and efficient manner. However, they do not account for  
3 transition uncertainty, whereas learning robust policies can be computationally  
4 expensive. In this paper, we introduce robust policy gradient (RPG), a policy-  
5 based method that efficiently solves rectangular robust Markov decision processes  
6 (MDPs). We provide a closed-form expression for the worst occupation measure.  
7 Incidentally, we find that the worst kernel is a rank-one perturbation of the nominal.  
8 Combining the worst occupation measure with a robust Q-value estimation yields  
9 an explicit form of the robust gradient. Our resulting RPG can be estimated from  
10 data with the same time complexity as its non-robust equivalent. Hence, it relieves  
11 the computational burden of convex optimization problems required for training  
12 robust policies by current policy gradient approaches.

## 13 1 Introduction

14 Markov decision processes (MDPs) provide an analytical framework to solve sequential decision-  
15 making problems and seek the best performance in a fixed environment. Since the resulting policy can  
16 be highly sensitive to parameter values [16], the robust MDP setting alternatively maximizes return  
17 under the worst scenario, thus yielding robustness to uncertain environments [18, 10]. In practice, the  
18 robust MDP paradigm quantifies the level of uncertainty through a set  $\mathcal{U}$  determining the possible  
19 range of model perturbations. Then, a policy is said to be robust-optimal if it reaches maximal  
20 performance under the most adversarial model within the uncertainty set. Developing efficient solvers  
21 for robust MDPs is of great interest, as it can lead to behavior policies with generalization guarantees  
22 [31].

23 If not computationally expensive, robust MDPs can be strongly NP-hard [30]. Thus, to preserve  
24 tractability, we commonly assume that  $\mathcal{U}$  is convex and  $s$ -rectangular, i.e.,  $\mathcal{U} = \times_{s \in \mathcal{S}} \mathcal{U}_s$  [18, 10, 30].  
25 The latter assumption means that the overall uncertainty should be designed independently for  
26 each state. Further simplification may consider  $(s, a)$ -rectangular uncertainty sets of the form  
27  $\mathcal{U} = \times_{(s,a) \in \mathcal{X}} \mathcal{U}_{(s,a)}$ , albeit this naturally leads to more conservative strategies. In any case, planning  
28 in robust MDPs can be computationally costly, as it involves successive max-min problems [7, 1, 30].  
29 To address this issue, the works [3, 12] have established an equivalence between robustness and  
30 regularization in reinforcement learning (RL) in order to derive efficient robust planning methods  
31 for  $s$  and  $(s, a)$ -rectangular robust MDPs. Indeed, it appears that resorting to proper regularization  
32 instead of solving a minimization problem can yield robust behavior without requiring the polynomial  
33 time complexity of convex optimization problems [3].

34 Alternatively to planning, policy gradient algorithms (PG) directly learn an optimal policy by applying  
35 gradient steps towards better performance [22]. Due to its scalability, ease of implementation, and  
36 adaptability to many different settings such as model-free and continuous state-action spaces [11, 21],

37 PG has become the workhorse of RL. Although regularization techniques such as max-entropy [6] or  
 38 Tsallis [13] have shown robust behavior without impairing computational cost, they only account for  
 39 adversarial reward [2, 5, 3]. Differently, robust PG formulations (RPG) formulations aim to address  
 40 uncertainty to reward *and* transition functions.

41 Despite their ability to propel robust behavior, RPG methods that target robust optimal policies  
 42 are still rare in the RL literature. The global convergence of RPG established in [14, 27] further  
 43 motivates us to come up with a practical method for estimating the gradients. In fact, [14] occult the  
 44 estimation part, as they assume full access to the policy gradient. Differently, the solution proposed  
 45 in [27] requires solving convex optimization problems to find the worst model, which represents  
 46 a time complexity of  $O(S^4 A \log \epsilon^{-1})$  for  $(s, a)$ -rectangular, or  $O(S^4 A^3 \log \epsilon^{-1})$  for  $s$ -rectangular  
 47 uncertainty sets [27, Sec. 4.1]. These worst kernel and reward models are needed to compute RPG  
 48 using the policy gradient theorem [23]. Other approaches that elicit an expression for RPG rely on a  
 49 specific type of uncertainty set such as reward uncertainty with known kernel [3],  $r$ -contaminated  
 50 kernel with known reward [29], or  $(s, a)$ -rectangular uncertainty [14], whereas we aim to tackle more  
 51 general robust MDPs.

52 In this work, we introduce an RPG method for both  $s$  and  $(s, a)$ -rectangular ball-constrained un-  
 53 certainty sets, with similar complexity as non-robust PG. Our approach provides a closed-form  
 54 expression of RPG without relying on an oracle while applying to the most common robust MDPs.  
 55 To this end, we derive the worst reward and transition functions, thus revealing the adversarial nature  
 56 of the corresponding uncertainty set. Surprisingly, we also find that the worst kernel is a rank-one  
 57 perturbation of the nominal kernel. Leveraging this rank-one perturbation enables us to derive a  
 58 robust occupation measure. We concurrently propose an alternative definition of the robust Q-value  
 59 together with an efficient way to estimate it. Combining these results enables us to obtain RPG in  
 60 closed form. Our resulting RPG update requires  $O(S^2 A \log \epsilon^{-1})$  computations, thus showing similar  
 61 time complexity as non-robust PG.

62 To summarize our contributions: (i) We establish the worst reward and transition models in closed-  
 63 form; (ii) We show that the worst-case transition function is a rank-one perturbation of the nominal;  
 64 (iii) We introduce alternative robust Q-values that can be evaluated through efficient Bellman recursion  
 65 while retrieving the robust value function; (iv) We establish an expression of RPG that can be  
 66 estimated with similar time complexity as non-robust PG. Experiments show that our RPG speeds up  
 67 state-of-the-art robust PG updates by 2 orders of magnitude.

## 68 2 Related work

69 Although some previous works use gradient methods to learn robust policies, they seek empirical  
 70 robustness to adversarial behavior rather than robust MDP solutions [19, 26, 4]. In that sense, our  
 71 study differs from adversarial RL as we explicitly optimize the max-min objective to find a robust  
 72 optimal policy. Accordingly, the risk-averse approach focuses on the *internal uncertainty* due to the  
 73 stochasticity of the system, whereas robust RL addresses the *external uncertainty* of the system’s  
 74 dynamics. As a result, common risk-averse objectives can be reformulated as robust problems with  
 75 specific uncertainty sets [24].

76 Previous studies that did aim to derive robust policy-based methods are [3, 29, 27]. These are  
 77 summarized in Table 1, which also displays the complexity of existing approaches. [3] established  
 78 RPG for  $s$ -rectangular reward-robust MDPs, i.e., robust MDPs with uncertain reward but given  
 79 kernel. Although it applies to general norms, their result does not account for transition perturbation.  
 80 Differently, in [29], the authors introduced RPG for  $r$ -contaminated MDPs, i.e., robust MDPs with  
 81 uncertainty set  $\mathcal{U} := \{R_0\} \times [(1-r)P_0 + r\Delta_S^{S \times A}]$ . Although it has similar complexity as non-robust  
 82 PG, by construction, their setting is limited to  $(s, a)$ -rectangularity with known reward and mixed  
 83 transition. As such, the proof techniques in [29] are tailor-made to the  $r$ -contamination framework  
 84 and do not apply to more general robust MDPs. In fact, we remark that the  $r$ -contamination setting is  
 85 equivalent to the action robustness approach introduced in [26], which emphasizes its limitation to  
 86 action perturbation. Differently, our RPG holds whenever the worst kernel is a rank-one perturbation  
 87 of the nominal transition function (see Lemma 4.4).

88 To address generic robust MDPs, [27] recently introduced RPG for general uncertainty sets. Their  
 89 gradient update has a complexity of  $O(S^6 A^4 \epsilon^{-4})$ , which is more expensive than non-robust PG by a  
 90 factor of  $S^4 A^3 \epsilon^{-4}$ . They additionally assume access to an oracle gradient of the robust return with

91 respect to the transition model. Avoiding this oracle assumption naturally leads to even higher time  
 92 complexity. At the same time, the two works [14, 27] guarantee global convergence of projected  
 93 robust gradient iterates, thus establishing the potential promise of RPG. In fact, equipped with RPG  
 94 convergence, the remaining challenge in making it practical is to efficiently estimate the gradient.  
 95 This represents the main focus of our study: We aim to explicit an RPG method that generalizes  
 96 existing results on specific uncertainty sets [3, 29] while holding for  $s$ -rectangular robust MDPs.

Table 1: Time complexity of RPG update according to the type of uncertainty set. For conciseness, the displayed complexity hides logarithmic factors in  $A$  and  $S$ . Our RPG method has the same complexity as non-robust PG while it generalizes other RPG methods with similar efficiency.

UNCERTAINTY SET $\mathcal{U}$	TIME COMPLEXITY	REFERENCE
$\{R_0\} \times \{P_0\}$	$S^2 A \log \epsilon^{-1}$	[23]
$\{R_0\} \times [(1-r)P_0 + r\Delta_S^{S \times A}]$	$S^2 A \log \epsilon^{-1}$	[29]
$(s, a)$ -rectangular ball $\mathcal{U}_p^{sa}$	$S^2 A \log \epsilon^{-1}$	<b>This work</b>
$(s, a)$ -rectangular, convex $\mathcal{U}^{sa}$	$S^4 A \log \epsilon^{-1}$	Convex optimization
$s$ -rectangular ball $\mathcal{U}_p^s$	$S^2 A \log \epsilon^{-1}$	<b>This work</b>
$s$ -rectangular ball $(R_0 + \mathcal{R}_p^s) \times \{P_0\}$	$S^2 A \log \epsilon^{-1}$	[3]
$s$ -rectangular, convex $\mathcal{U}^s$	$S^4 A^3 \log \epsilon^{-1}$	Convex optimization
$s$ -rectangular, convex $\mathcal{U}^s$	$S^6 A^4 \epsilon^{-4}$	[27]
$s$ -rectangular, non-convex $\mathcal{U}^s$	NP-hard	[30]
Non-rectangular, convex $\mathcal{U}$	NP-hard	[30]

### 97 3 Preliminaries

98 **Notation:** We denote the cardinal of an arbitrary finite set  $\mathcal{Z}$  by  $|\mathcal{Z}|$ . Given two real functions  
 99  $\mathbf{a}, \mathbf{b} : \mathcal{Z} \rightarrow \mathbb{R}$ , their inner product is  $\langle \mathbf{a}, \mathbf{b} \rangle_{\mathcal{Z}} := \sum_{z \in \mathcal{Z}} \mathbf{a}(z) \mathbf{b}(z)$ , which induces the  $\ell_2$ -norm  
 100  $\|\mathbf{a}\|_2 := \sqrt{\langle \mathbf{a}, \mathbf{a} \rangle_{\mathcal{Z}}}$ . More generally, the  $\ell_p$ -norm of  $\mathbf{a}$  is denoted by  $\|\mathbf{a}\|_p$  whose conjugate norm  
 101 is  $\|\mathbf{a}\|_q := \max_{\|\mathbf{b}\|_p \leq 1} \langle \mathbf{a}, \mathbf{b} \rangle_{\mathcal{Z}}$  with  $q^{-1} = 1 - p^{-1}$ . The vector of all zeros (resp. all ones)  
 102 with appropriate dimensions is denoted by  $\mathbf{0}$  (resp.  $\mathbf{1}$ ), and the probability simplex over  $\mathcal{Z}$  by  
 103  $\Delta_{\mathcal{Z}} := \{\mathbf{a} : \mathcal{Z} \rightarrow \mathbb{R}_+ \mid \langle \mathbf{a}, \mathbf{1} \rangle_{\mathcal{Z}} = 1\}$ . Finally,  $I_{\mathcal{Z}}$  designates the identity matrix in  $\mathbb{R}^{\mathcal{Z} \times \mathcal{Z}}$ .

#### 104 3.1 Markov Decision Processes

A Markov decision process (MDP) is a tuple  $(\mathcal{S}, \mathcal{A}, \gamma, \mu, P, R)$  such that  $\mathcal{S}$  and  $\mathcal{A}$  are finite state and action spaces of cardinal  $S$  and  $A$  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  $(P, R)$  corresponds to the MDP model with  $P : \mathcal{X} \rightarrow \Delta_{\mathcal{S}}$  being a transition kernel and  $R : \mathcal{X} \rightarrow \mathbb{R}$  a reward function. A policy  $\pi : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$  maps each state to a probability distribution over  $\mathcal{A}$ , and we denote by  $\Pi$  the set of such functions. For any policy  $\pi \in \Pi$ ,  $R^\pi \in \mathbb{R}^{\mathcal{S}}$  is the expected immediate reward defined as  $R^\pi(s) := \langle \pi_s, R(s, \cdot) \rangle_{\mathcal{A}}$ ,  $\forall s \in \mathcal{S}$ , where  $\pi_s$  is a shorthand for  $\pi(\cdot | s)$ . We similarly define the stochastic matrix induced by  $\pi$  as  $P^\pi(s' | s) := \langle \pi_s, P(s' | s, \cdot) \rangle_{\mathcal{A}}$ ,  $\forall s, s' \in \mathcal{S}$ , and extend the occupation measure to an arbitrary initial vector  $k \in \mathbb{R}^{\mathcal{S}}$  by defining

$$d_{P,k}^\pi := k^\top (I_{\mathcal{S}} - \gamma P^\pi)^{-1}.$$

The performance measure we aim to maximize is the value function  $v_{(P,R)}^\pi := (I_{\mathcal{S}} - \gamma P^\pi)^{-1} R^\pi$ , or alternatively, the return  $\rho_{(P,R)}^\pi := \langle \mu, v_{(P,R)}^\pi \rangle_{\mathcal{S}}$ . We denote the optimal value function (resp. optimal return) by  $v_{(P,R)}^* = \max_{\pi \in \Pi} v_{(P,R)}^\pi$  (resp.  $\rho_{(P,R)}^* = \langle \mu, v_{(P,R)}^* \rangle_{\mathcal{S}}$ ). It can be obtained using Bellman operators, which are defined as  $T_{(P,R)}^\pi v := R^\pi + \gamma P^\pi v$  and  $T_{(P,R)}^* v := \max_{\pi \in \Pi} T_{(P,R)}^\pi v$ ,  $\forall v \in \mathbb{R}^{\mathcal{S}}$ , respectively [20]. For any vector  $v \in \mathbb{R}^{\mathcal{S}}$ , we associate its Q-function  $Q \in \mathbb{R}^{\mathcal{X}}$  such that

$$Q(s, a) = r(s, a) + \gamma \langle P(\cdot | s, a), v \rangle_{\mathcal{S}}, \quad \forall (s, a) \in \mathcal{X}.$$

With a slight abuse of notation, we can similarly define a Bellman operator over Q-values as

$$T_{(P,R)}^\pi Q(s, a) := r(s, a) + \gamma \sum_{(s', a') \in \mathcal{X}} P(s' | s, a) \pi_{s'}(a') Q(s', a'), \quad \forall (s, a) \in \mathcal{X}.$$

### 105 3.2 Robust Markov Decision Processes

106 In a robust MDP setting, we assume that  $(P, r) \in \mathcal{U}$  and aim to maximize return under the worst model  
 107 from the set. We denote the robust performance of a policy  $\pi \in \Pi$  by  $\rho_{\mathcal{U}}^{\pi} := \min_{(P,R) \in \mathcal{U}} \rho_{(P,R)}^{\pi}$ . It  
 108 is maximal when it reaches  $\rho_{\mathcal{U}}^* := \max_{\pi \in \Pi} \rho_{\mathcal{U}}^{\pi}$  at an optimal robust policy  $\pi_{\mathcal{U}}^* \in \arg \max_{\pi \in \Pi} \rho_{\mathcal{U}}^{\pi}$ .  
 109 When considering the robust value function  $v_{\mathcal{U}}^{\pi} := \min_{(P,R) \in \mathcal{U}} v_{(P,R)}^{\pi}$ , we further need to assume  
 110 that  $\mathcal{U}$  is convex and rectangular so that an optimal robust policy realizing  $v_{\mathcal{U}}^* := \max_{\pi} v_{\mathcal{U}}^{\pi}$  can  
 111 be computed in polynomial time [30]. We thus assume  $\mathcal{U}$  to be convex and rectangular in the  
 112 remainder of this work. Specifically, we denote an  $(s, a)$ -rectangular uncertainty set by  $\mathcal{U}^{sa} :=$   
 113  $\times_{(s,a) \in \mathcal{X}} (\mathcal{P}_{(s,a)}, \mathcal{R}_{(s,a)})$ . It represents a particular case of  $s$ -rectangular uncertainty which we  
 114 similarly denote by  $\mathcal{U}^s := \times_{s \in \mathcal{S}} (\mathcal{P}_s, \mathcal{R}_s)$ . In both cases, there exists an optimal robust policy that is  
 115 stationary, although all optimal ones may be stochastic [30].

116 Similarly to non-robust MDPs, robust MDPs can be solved through Bellman recursion. Indeed, the  
 117 robust value function  $v_{\mathcal{U}}^{\pi}$  (resp., optimal robust value function  $v_{\mathcal{U}}^*$ ) is known to be the unique fixed  
 118 point of the  $\gamma$ -contracting robust Bellman operator  $T_{\mathcal{U}}^{\pi} v := \min_{(P,R) \in \mathcal{U}} T_{(P,R)}^{\pi} v$  (resp., the optimal  
 119 robust Bellman operator  $T_{\mathcal{U}}^* v := \max_{\pi \in \Pi} T_{\mathcal{U}}^{\pi} v$ ), both defined for any  $v \in \mathbb{R}^{\mathcal{S}}$ . Although this ensures  
 120 linear convergence of robust value iteration, the evaluation of each Bellman operator can still be  
 121 prohibitive for practical use.

#### 122 3.2.1 Ball Constrained Uncertainty set

123 To facilitate the computation of robust Bellman updates, we consider uncertainty sets that are centered  
 124 around a nominal model  $(P_0, R_0)$ , i.e., of the form  $\mathcal{U} = (P_0, R_0) + (\mathcal{P}, \mathcal{R})$ , and constrained according  
 125 to  $\ell_p$ -norm balls [3, 12, 7, 1]. In the  $(s, a)$ -rectangular case, the corresponding uncertainty set is  
 126 denoted by  $\mathcal{U}_p^{sa} := \mathcal{R}_p^{sa} \times \mathcal{P}_p^{sa} = \times_{(s,a) \in \mathcal{X}} (\mathcal{P}_{(s,a)}, \mathcal{R}_{(s,a)})$  where for any  $(s, a) \in \mathcal{X}$ ,

$$\mathcal{R}_{(s,a)} = \{r \in \mathbb{R} \mid \|r\|_p \leq \alpha_{s,a}\}, \quad \text{and} \quad \mathcal{P}_{(s,a)} = \{p \in \mathbb{R}^{\mathcal{S}} \mid \langle p, \mathbf{1} \rangle_{\mathcal{S}} = 0, \|p\|_p \leq \beta_{s,a}\}.$$

127 Similarly, an  $s$ -rectangular norm-constrained uncertainty is denoted by  $\mathcal{U}_p^s := \times_{s \in \mathcal{S}} (\mathcal{P}_s, \mathcal{R}_s)$  where  
 128 for any  $s \in \mathcal{S}$ ,

$$\mathcal{R}_s = \{r \in \mathbb{R}^{\mathcal{A}} \mid \|r\|_p \leq \alpha_s\}, \quad \text{and} \quad \mathcal{P}_s = \{p \in \mathbb{R}^{\mathcal{X}} \mid \langle p(\cdot, a), \mathbf{1} \rangle_{\mathcal{S}} = 0 \quad \forall a \in \mathcal{A}, \|p\|_p \leq \beta_s\}.$$

129 In both cases, the noise radius  $\beta$  should be small enough so that transition kernels of the form  $P_0 + \mathcal{P}$   
 130 are well defined. This normed ball structure on the uncertainty sets enables us to compute robust  
 131 Bellman updates with similar time complexity as non-robust ones using regularization [3, 12].

132 First, define the generalized variance function and the mean function as

$$\kappa_q(v) = \min_{w \in \mathbb{R}} \|v - w\mathbf{1}\|_q, \quad \omega_q(v) \in \arg \min_{w \in \mathbb{R}} \|v - w\mathbf{1}\|_q,$$

133 respectively, where  $q$  is the conjugate value of  $p$  (see Tab. 2 for their closed-form expression when  
 134  $q \in \{1, 2, \infty\}$ ). Then, we can efficiently evaluate robust value functions by regularizing a standard  
 135 Bellman operator instead of solving a minimization. We formalize this below.

136 **Proposition 3.1.** ([12, Thm. 2-3].) *For any policy  $\pi \in \Pi$  and any rectangular  $\ell_p$ -ball-constraint*  
 137 *uncertainty set, the robust Bellman operator is equivalent to its regularized form:*

$$(T_{\mathcal{U}}^{\pi} v)(s) = T_{(P_0, R_0)}^{\pi} v(s) + \Omega_q(\alpha, \beta, v),$$

138 where  $\Omega_q(\alpha, \beta, v) := -\langle \pi_s, \alpha_{s,\cdot} + \gamma \kappa_q(v) \beta_{s,\cdot}^p \rangle_{\mathcal{A}}$  for  $(s, a)$ -rectangular uncertainty  $\mathcal{U}_p^{sa}$ , and  
 139  $\Omega_q(\alpha, \beta, v) := -(\alpha_s + \gamma \beta_s \kappa_q(v)) \|\pi_s\|_q$  for  $s$ -rectangular uncertainty  $\mathcal{U}_p^s$ .

140 In the following, we leverage the regularized formulation of robust value functions to explicitly derive  
 141 RPG for rectangular  $\ell_p$ -ball uncertainty sets.

#### 142 3.2.2 Robust Gradient Method

143 Since the robust return can be non-differentiable, we need to follow the projected sub-gradient ascent  
 144 rule in order to optimize the robust return, namely, update  $\pi_{k+1} := \mathbf{proj}_{\Pi}(\pi_k + \eta \partial_{\pi} \rho_{\mathcal{U}}^{\pi_k})$  where

$$\partial_{\pi} \rho_{\mathcal{U}}^{\pi} := \nabla_{\pi} \rho_{(P,R)}^{\pi} \Big|_{(P,R)=(P_{\mathcal{U}}^{\pi}, R_{\mathcal{U}}^{\pi})}, \quad (1)$$

Table 2: Expressions of the  $q$ -mean, the  $q$ -variance, and its gradient. We assume that the vector  $v$  is sorted, i.e.,  $v(s_i) \geq v(s_{i+1}), \forall i \in \{1, 2, \dots, S\}$ , and denote  $n_l := \lfloor (S+1)/2 \rfloor, n_u := \lceil (S+1)/2 \rceil$ .

	$\omega_q(v)$	$\kappa_q(v)$	$\nabla_v \kappa_q(v)$
$q$	$\arg \min_{w \in \mathbb{R}} \ v - w\mathbf{1}\ _q$	$\min_{\omega \in \mathbb{R}} \ v - \omega\mathbf{1}\ _q$	$\frac{\partial \kappa_q(v)}{\partial v(s_i)}$
$\infty$	$\frac{v(s_1) + v(s_S)}{2}$	$\frac{v(s_1) - v(s_S)}{2}$	$\begin{cases} \frac{1}{2} & \text{if } i = 1 \\ -\frac{1}{2} & \text{if } i = S \\ 0 & \text{o.w.} \end{cases}$
$2$	$\frac{\sum_{i=1}^S v(s_i)}{S}$	$\sqrt{\sum_{i=1}^S (v(s_i) - \omega_2(v))^2}$	$\frac{v(s_i) - \omega_2(v)}{\kappa_2(v)}$
$1$	$\frac{v(s_{n_l}) + v(s_{n_u})}{2}$	$\sum_{i=1}^{n_l} (v(s_i) - v(s_{S-i}))$	$\begin{cases} 1 & \text{if } i < n_l \\ -1 & \text{if } i > n_u \\ 0 & \text{o.w.} \end{cases}$

145  $\eta$  is the learning rate,  $\mathbf{proj}_\Pi$  denotes the orthogonal projection on  $\Pi$ , and  $(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)$  is the worst model  
 146 associated with  $\pi \in \Pi$  and  $\mathcal{U}$ , i.e.,  $(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi) \in \arg \inf_{(P,R) \in \mathcal{U}} \rho_{(P,R)}^\pi$ .

147 Given oracle access to sub-gradient  $\partial \rho_{\mathcal{U}}^\pi$ , projected gradient ascent converges to an  $\epsilon$ -optimal policy  
 148  $\pi_{\mathcal{U}}^*$ . Moreover, under similar conditions as in the non-robust setting, projected gradient ascent holds  
 149 an iteration complexity of  $O(S^4 A^2 \epsilon^{-4})$  [27]. Yet, the sub-gradient in (1) is generally intractable,  
 150 particularly because general convex uncertainty sets may yield NP-hard complexity. Instead, we  
 151 propose to focus on ball-constrained uncertainty sets in order to efficiently compute RPG updates.

## 152 4 Towards RPG: Expressing the worst quantities

153 In this section, we provide all the ingredients needed for deriving RPG. Before diving into the  
 154 gradient expression, we first settle on the general framework of policy gradient. Secondly, in Sec. 4.1,  
 155 we focus on expressing the worst model according to the nominal explicitly. Surprisingly, we find  
 156 that the worst transition kernel is a rank-one perturbation of the nominal. This finding enables us  
 157 to derive the robust occupancy measure, i.e., the occupation measure of the worst kernel in Sec. 4.2.  
 158 As a last piece, in Sec. 4.3, we propose an alternative definition of robust Q-value and show that  
 159 it can be estimated from a specific Bellman recursion.

160 Consider again the projected gradient ascent rule:

$$\pi_{k+1} := \mathbf{proj}_\Pi(\pi_k + \eta \partial_\pi \rho_{\mathcal{U}}^{\pi_k}).$$

161 By definition of the sub-gradient in (1) and applying the standard PG theorem [23], it holds that:

$$\partial_\pi \rho_{\mathcal{U}}^\pi = \sum_{(s,a) \in \mathcal{X}} d_{\mathcal{U}}^\pi(s) Q_{\mathcal{U}}^\pi(s,a) \nabla \pi_s(a),$$

162 where  $Q_{\mathcal{U}}^\pi := Q_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}^\pi$  is the Q-value associated with the worst-case model, and  $d_{\mathcal{U}}^\pi := d_{P_{\mathcal{U}}^\pi}^\pi$  the  
 163 occupation measure of the worst transition kernel. In fact, for the uncertainty sets we focus on in  
 164 this work, the worst Q-value  $Q_{\mathcal{U}}^\pi$  retrieves the common definition of robust Q-value [18, 25] (see  
 165 the appendix for a detailed discussion). Therefore, for conciseness and with a slight abuse, we shall  
 166 designate  $Q_{\mathcal{U}}^\pi$  by the robust Q-value, and  $d_{\mathcal{U}}^\pi$  by the robust occupation measure. The remaining  
 167 question is how to compute these quantities and in particular, can we efficiently find the worst  
 168 parameters  $(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)$ ? The following part of our study aims to address these questions.

169 Given an uncertainty set  $\mathcal{U}$ , let first define the normalized and balanced robust value function as:

$$u_{\mathcal{U}}^\pi(s) := \frac{\text{sign}(v_{\mathcal{U}}^\pi(s) - \omega_q(v_{\mathcal{U}}^\pi)) \|v_{\mathcal{U}}^\pi(s) - \omega_q(v_{\mathcal{U}}^\pi)\|^{q-1}}{\kappa_q(v_{\mathcal{U}}^\pi)^{q-1}}. \quad (2)$$

170 By construction, it has zero mean and unit norm, i.e.,  $\langle u_{\mathcal{U}}^{\pi}, \mathbf{1} \rangle_{\mathcal{S}} = 0$  and  $\|u_{\mathcal{U}}^{\pi}\|_p = 1$ . In fact, as  
 171 stated in the result below,  $u_{\mathcal{U}}^{\pi}$  is the gradient of the  $q$ -variance function, and correlates with the  
 172 (unnormalized, unbalanced) robust value function according to the same  $q$ -variance.

173 **Proposition 4.1.** *For any policy  $\pi \in \Pi$  and  $\ell_p$ -ball rectangular uncertainty set, the following holds:*

$$\begin{aligned} u_{\mathcal{U}}^{\pi} &= \nabla_v \kappa_q(v) \Big|_{v=v_{\mathcal{U}}^{\pi}}, \\ \langle u_{\mathcal{U}}^{\pi}, v_{\mathcal{U}}^{\pi} \rangle &= \kappa_q(v_{\mathcal{U}}^{\pi}). \end{aligned}$$

#### 174 4.1 Worst Kernel and Reward

175 In the following results, we explicit the relationship between the nominal and the worst-case model  
 176 for  $(s, a)$  and  $s$ -rectangular  $\ell_p$ -balls. We will then leverage this relationship to compute the robust  
 177 Q-values and the robust occupation measure, both necessary for RPG.

178 **Theorem 4.2** ( $(s, a)$ -rectangular case). *Given uncertainty set  $\mathcal{U} = \mathcal{U}_p^{s,a}$  and any policy  $\pi \in \Pi$ , the  
 179 worst model is related to the nominal one through:*

$$R_{\mathcal{U}}^{\pi}(s, a) = R_0(s, a) - \alpha_{s,a} \quad \text{and} \quad P_{\mathcal{U}}^{\pi}(\cdot|s, a) = P_0(\cdot|s, a) - \beta_{s,a} u_{\mathcal{U}}^{\pi}.$$

Based on Thm. 4.2, it follows that in the  $(s, a)$ -rectangular case, the worst reward function is  
 independent of the employed policy. As we establish in Thm. 4.3 below, this no longer applies under  
 $s$ -rectangularity. In either case, the worst kernel is policy-dependent, discouraging the system to  
 move toward high-rewarding states and directing it to low-rewarding ones instead. Surprisingly, the  
 vector penalty  $u_{\mathcal{U}}^{\pi} \in \mathbb{R}^{\mathcal{S}}$  additionally illustrates that the worst kernel is a rank-one perturbation of the  
 nominal. Indeed, considering the stochastic matrix induced by any policy  $\pi \in \Pi$ , we have

$$[P_{\mathcal{U}}^{\pi} - P_0^{\pi}](s'|s) = - \left( \sum_{a \in \mathcal{A}} \beta_{s,a} \pi_s(a) \right) u_{\mathcal{U}}^{\pi}(s'), \quad \forall s \in \mathcal{S},$$

180 so that the perturbation matrix  $P_{\mathcal{U}}^{\pi} - P_0^{\pi}$  is of rank one. In the sequel, we will leverage this finding to  
 181 compute the robust occupation measure.

182 **Theorem 4.3** ( $s$ -rectangular case). *Given uncertainty set  $\mathcal{U} = \mathcal{U}_p^s$  and any policy  $\pi \in \Pi$ , the worst  
 183 model is related to the nominal one through:*

$$R_{\mathcal{U}}^{\pi}(s, a) = R_0(s, a) - \alpha_s \left( \frac{\pi_s(a)}{\|\pi_s\|_q} \right)^{q-1} \quad \text{and} \quad P_{\mathcal{U}}^{\pi}(\cdot|s, a) = P_0(\cdot|s, a) - \beta_s u_{\mathcal{U}}^{\pi} \left( \frac{\pi_s(a)}{\|\pi_s\|_q} \right)^{q-1}.$$

184 Similarly to the  $(s, a)$ -case, the adversarial kernel is a rank-one perturbation of the nominal. Yet, an  
 185 extra dependence on the policy through the coefficient  $\left( \frac{\pi_s(a)}{\|\pi_s\|_q} \right)^{q-1}$  appears in the  $s$ -case, affecting  
 186 both the worst reward and the worst kernel. Intuitively, it means that the worst model cannot be  
 187 chosen independently for each action, but must instead depend on the agent's policy. This further  
 188 explains why optimal policies can all be stochastic in  $s$ -rectangular robust MDPs [30].

189 Thms. 4.2 and 4.3 enable us to derive the worst MDP model in closed form with time complexity  
 190  $O(S^2 A \log \epsilon^{-1})$ , up to logarithmic factors (please see the appendix for a detailed discussion). It thus  
 191 holds the same complexity as non-robust value iteration, since we additionally need to compute the  
 192 value function to derive its corresponding regularizer [3, 12]. On the other hand, if we employ convex  
 193 optimization using value methods instead, obtaining the worst model requires a time complexity  
 194 of  $O(S^4 A \log \epsilon^{-1})$  in the  $(s, a)$ -rectangular case, and  $O(S^4 A^3 \log \epsilon^{-1})$  in the  $s$ -rectangular case  
 195 [27][Sec. 4.1].

#### 196 4.2 Robust Occupation Measure

197 We finally derive the robust occupation measure using nominal values, which will lead to an explicit  
 198 RPG. Although intractable in general, we show that focusing on ball-constrained uncertainty enables  
 199 deriving the robust occupation matrix efficiently from the (nominal) occupation measure. We first  
 200 establish the lemma below, which leverages the fact that the worst transition function is a rank-one  
 201 perturbation of the nominal and represents our core contribution.

202 **Lemma 4.4.** Let  $b, k \in \mathbb{R}^S$  and  $P_0, P_1 \in (\Delta_S)^S$  two transition matrices. If  $P_1 = P_0 - bk^\top$ , i.e.,  $P_1$   
 203 is a rank-one perturbation of  $P_0$ , then their occupation matrices  $D_i := (I - \gamma P_i)^{-1}$ ,  $i = 0, 1$  are  
 204 related through:

$$D_1 = D_0 - \gamma \frac{D_0 b k^\top D_0}{(1 + \gamma k^\top D_0 b)}.$$

205 Combining Thms. 4.2 and 4.3 with the above lemma, we obtain the robust occupation in terms of the  
 206 nominal, as stated in Thm. 4.6 below. Prior to this, we introduce the notion of *expected transition*  
 207 *uncertainty* below.

208 **Definition 4.5.** Let  $\mathcal{U}$  a rectangular  $\ell_p$ -ball-constrained uncertainty set of transition radius  $\beta$ . For  
 209 any policy  $\pi \in \Pi$ , the expected transition uncertainty at any state  $s \in \mathcal{S}$  is given by  $\beta_s^\pi :=$   
 210  $\sum_{a \in \mathcal{A}} \pi_s(a) \beta_{s,a}$  if  $\mathcal{U} = \mathcal{U}_p^{s^a}$ , and  $\beta_s^\pi := \beta_s \|\pi_s\|_q$  if  $\mathcal{U} = \mathcal{U}_p^s$ .

211 **Theorem 4.6.** For any rectangular  $\ell_p$ -ball-constrained uncertainty and  $\pi \in \Pi$ , it holds that:

$$d_{\mathcal{U}, \mu}^\pi = d_{P_0, \mu}^\pi - \gamma \frac{\langle d_{P_0, \mu}^\pi, \beta^\pi \rangle_{\mathcal{S}}}{1 + \gamma \langle d_{P_0, u_{\mathcal{U}}^\pi}^\pi, \beta^\pi \rangle_{\mathcal{S}}} d_{P_0, u_{\mathcal{U}}^\pi}^\pi. \quad (3)$$

212 Thm. 4.6 explicitly highlights the relationship between the robust occupation measure and the nominal  
 213 one. Thus, according to Eq. (3), the standard non-robust occupation measure in the first term needs  
 214 to be penalized by another one,  $d_{P_0, u_{\mathcal{U}}^\pi}^\pi = (u_{\mathcal{U}}^\pi)^\top (I_S - \gamma P_0^\pi)^{-1}$ , to obtain the robust occupation  
 215 measure. Recall that  $u_{\mathcal{U}}^\pi$  is the balanced-scaled value function determined by  $\pi \in \Pi$  and uncertainty  
 216 set  $\mathcal{U}$ . Thus, the penalty term  $d_{P_0, u_{\mathcal{U}}^\pi}^\pi$  tends to zero if all coordinates of the robust value function  
 217 vector converge to the same value.

218 Nonetheless, our expression (3) does present some challenges. First, the occupation measure appear-  
 219 ing in the correction term indicates that instead of taking a fixed initial state distribution, we should  
 220 start from a *varying* and *signed* measure represented by the balanced value function. Although it  
 221 suggests putting more weight on worst-performing states, obtaining a non-biased estimator for this  
 222 occupancy measure remains unclear in model-free learning.

### 223 4.3 Robust Q-values

In this section, we focus on the last element needed for RPG and aim to estimate the robust Q-  
 value denoted previously by  $Q_{\mathcal{U}}^\pi := Q_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}$ . Define its associated value function as  $v_{\mathcal{U}}^\pi(s) =$   
 $\langle \pi_s, Q_{\mathcal{U}}^\pi(s, \cdot) \rangle$ ,  $\forall s \in \mathcal{S}, \pi \in \Pi$ . Based on standard Bellman recursion, it thus holds that:

$$Q_{\mathcal{U}}^\pi(s, a) = R_{\mathcal{U}}^\pi(s, a) + \gamma \langle P_{\mathcal{U}}^\pi(\cdot | s, a), v_{\mathcal{U}}^\pi \rangle_{\mathcal{S}}, \quad \forall (s, a) \in \mathcal{X}, \pi \in \Pi,$$

224 while  $Q_{\mathcal{U}}^\pi$  is the unique fixed point of the  $\gamma$ -contracting operator

$$(\mathcal{L}_{\mathcal{U}}^\pi Q)(s, a) := T_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}^\pi Q(s, a), \quad \forall Q \in \mathbb{R}^{\mathcal{X}}. \quad (4)$$

225 The relations above hold for general uncertainty sets, provided that we have access to the worst model.  
 226 The  $s$ -rectangularity assumption additionally enables us to retrieve the robust value function using  
 227 the Bellman operator above [30]. Concretely, we have:  $v_{\mathcal{U}}^\pi = \min_{(P, R) \in \mathcal{U}} v_{(P, R)}^\pi = v_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}^\pi$ .

228 The following result derives a regularized operator equivalent to  $\mathcal{L}_{\mathcal{U}}^\pi$ , which results in an efficient  
 229 iteration method to compute the robust Q-value.

230 **Proposition 4.7.** The Bellman operator  $\mathcal{L}_{\mathcal{U}}^\pi$  defined in Eq. (4) is equivalent to:

$$(\mathcal{L}_{\mathcal{U}}^\pi Q)(s, a) = T_{(P_0, R_0)}^\pi Q(s, a) + \Omega'_q(\alpha_{s,a}, \beta_{s,a}, v),$$

231 where  $v(s) := \langle \pi_s, Q(s, \cdot) \rangle_{\mathcal{A}}$ ,  $\Omega'_q(\alpha, \beta, v) := -(\alpha_{s,a} + \gamma \beta_{s,a} \kappa_q(v))$  for  $(s, a)$ -rectangular uncer-  
 232 tainty  $\mathcal{U}_p^{s^a}$ , and  $\Omega'_q(\alpha, \beta, v) := -\left(\frac{\pi_s(a)}{\|\pi_s\|_q}\right)^{q-1} (\alpha_s + \gamma \beta_s \kappa_q(v))$  for  $s$ -rectangular  $\mathcal{U}_p^s$ .

## 233 5 Robust Policy Gradient

234 We are now able to derive an RPG by combining our previous results. Notably, unlike previous works  
 235 that need to sample next-state transitions based on all models from the uncertainty set [19, 15, 4],  
 236 here, we only need the nominal kernel to get the occupation measures.

237 **Theorem 5.1 (RPG).** *For any rectangular  $\ell_p$ -ball-constrained uncertainty, the robust policy gradient*  
 238 *is given by:*

$$\partial_\pi \rho_{\mathcal{U}}^\pi = \sum_{(s,a) \in \mathcal{X}} (d_{P_0, \mu}^\pi(s) - c^\pi(s)) Q_{\mathcal{U}}^\pi(s, a) \nabla \pi_s(a), \quad (5)$$

239 *where*

$$c^\pi(s) := \frac{\gamma \langle d_{P_0, \mu}^\pi, \beta^\pi \rangle_{\mathcal{S}}}{1 + \gamma \langle d_{P_0, u_{\mathcal{U}}^\pi}, \beta^\pi \rangle_{\mathcal{S}}} d_{P_0, u_{\mathcal{U}}^\pi}^\pi(s), \quad \forall s \in \mathcal{S}.$$

240 Thm. 5.1 is a direct application of non-robust PG, as its proof simply consists in plugging Eq. (3)  
 241 into the standard PG expression  $\partial_\pi \rho_{\mathcal{U}}^\pi = \sum_{(s,a) \in \mathcal{X}} d_{\mathcal{U}, \mu}^\pi(s) Q_{\mathcal{U}}^\pi(s, a) \nabla \pi_s(a)$ . We obtain a regular  
 242 PG in the first term, with the robust Q-value instead of the non-robust one, plus a correction term  
 243  $c^\pi$  resulting from taking the occupation measure of the worst kernel instead of the nominal. Unlike  
 244 previous work that uses policy regularization to achieve empirical robustness in PG methods [2, 9],  
 245 Thm. 5.1 establishes an RPG that accounts for transition uncertainty and targets a robust optimal  
 246 policy.

## 247 5.1 Complexity Analysis

248 A major concern in solving robust MDPs is time complexity [30]. Similarly, it is of major importance  
 249 to assess the additional time required for computing an RPG update, compared to its non-robust  
 250 variant. Although previous work has analyzed the convergence rate of RPG to a global optimum  
 251 [27], it assumes access to an oracle gradient, thus occulting the computational concerns raised from  
 252 gradient estimation. In fact, the NP-hardness of non-rectangular and/or non-convex robust MDPs  
 253 [30] already indicates that their resulting RPG can be intractable.

254 To compute RPG in Thm. 5.1, we first need to evaluate the robust Q-value. Based on Lemma 4.7 and  
 255 the Bellman operators introduced there, our evaluation method involves an additional estimation of  
 256 the variance function  $\kappa_p$ . According to [12], this takes logarithmic time at most, using binary search.  
 257 As to the compensation term  $c^\pi$  in Eq. (11), it requires computing occupancy measures with respect  
 258 to two different initial vectors, namely the balanced value function and the initial distribution. Thus,  
 259 the computational cost for estimating  $c^\pi$  is the same as estimating a non-robust occupancy measure.  
 260 Tab. 1 summarizes the complexity of different approaches while a detailed discussion can be found in  
 261 the appendix. We refer to [27][Sec. 4.1] for the complexity of RPG based on convex optimization.

262 **Generalization to arbitrary norms.** Until now, we have focused on  $\ell_p$ -norm for concreteness.  
 263 However, the above results apply to any norm  $\|\cdot\|$ , at least if the uncertainty set is  $(s, a)$ -rectangular,  
 264 in which case the variance function changes to  $\kappa(v) := \min_{\|c\| \leq 1, \mathbf{1}^\top c = 0} \langle c, v \rangle$  and the balanced value  
 265 to  $\arg \min_{\|c\| \leq 1, \mathbf{1}^\top c = 0} \langle c, v \rangle$ . The rank-one perturbation structure of the worst kernel is preserved, so  
 266 the robust occupation measure can be obtained similarly using Lemma 4.4. The  $s$ -rectangular is more  
 267 involved. We defer its discussion to the appendix and leave its complete derivation for future work.

## 268 6 Experiments

269 In order to test the effectiveness of our RPG update, we evaluate its increased time complexity  
 270 relative to non-robust PG. In the following experiments, we randomly generate nominal models for  
 271 arbitrary state-action space sizes. Each experiment was averaged over 100 runs. We refer the reader  
 272 to the appendix for more details on the radius levels and other implementation choices.

273 We first focus on  $\ell_1$ -robust MDPs to compare our RPG with a convex optimization approach.  
 274 Specifically, we consider a robust PG with an optimization solver, which we designate by LP-RPG.  
 275 Indeed, recall that  $\ell_1$ -ball-constraints induce a linear program (LP) rather than a more general convex  
 276 optimization problem. Therefore, to compute the robust value function for a given policy, we  
 277 iteratively evaluate the robust Bellman operator using LP [27, Section 4.1]. Using this approximated  
 278 value function, we can compute the worst value parameters to apply PG theorem by [23] and deduce  
 279 an LP-based robust PG update. Differently, our RPG method relies on the regularized formulation  
 280 of robust value iteration proposed in [3, 12], from which we deduce the normalized-balanced value  
 281 function as in Eq. (10). We finally apply Thm. 4.6 to compute the robust occupation measure, and  
 282 Prop. 4.7 to obtain the robust Q-value.

283 Tab. 3 displays the results obtained for the two alternative methods described above. In all experiments,  
 284 the standard deviation was typically 2-10% so we omitted it for brevity. As can be seen in Tab. 3,  
 285 LP-RPG does not scale well compared to RPG, whereas RPG has similar time complexity as PG.  
 286 Notably, the running time of  $s$ -rectangular LP-RPG scales much better with the space size than its  
 287  $(s, a)$ -rectangular equivalent, which confirms the theoretical complexities from Tab. 1. Yet, since  
 288 these methods were time-consuming, we repeated these for a few runs only. In fact, LP-RPG is more  
 289 expensive than RPG by 1-3 orders of magnitude, which illustrates its inefficiency. We emphasize  
 290 that here, we only focused on  $\ell_1$ -robust MDPs to leverage LP solvers in robust policy evaluation. We  
 291 expect the computational cost of LP-RPG to scale even more poorly for other  $\ell_p$ -robust MDPs that  
 292 involve polynomial time-consuming convex programs.

Table 3: Comparison of the relative running time between RPG and the convex optimization approach (here, LP). Our method is faster than LP-based updates by 1 to 3 orders of magnitude.

		$\  \{(P_0, R_0)\} \ $		$U_1^{sa}$		$U_1^s$	
S	A	PG	RPG	LP-RPG	RPG	LP-RPG	
10	10	1	1.4	326	1.4	77	
30	10	1	1.4	351	1.4	109	
50	10	1	1.4	408	1.4	159	
100	20	1	1.5	469	1.3	268	
500	50	1	1.3	925	1.3	5343	

293 We further compare our RPG to non-robust PG on different  $\ell_p$ -balls. Tab. 4 confirms the comparable  
 294 time complexity of RPG to non-robust PG, thus demonstrating the effectiveness of our method. We  
 295 note that for  $p \in \{1, 2, \infty\}$ , the corresponding regularization quantities can be computed in closed  
 296 form, whereas they involve a binary search for other values [12]. We thus get a slight running-time  
 297 increase for  $p \in \{5, 10\}$ .

Table 4: Relative running time for computing RPG under different types of uncertainty sets.

S	A	$\  \{(P_0, R_0)\} \ $	$U_2^{sa}$	$U_2^s$	$U_5^{sa}$	$U_5^s$	$U_{10}^{sa}$	$U_{10}^s$	$U_{\infty}^{sa}$	$U_{\infty}^s$
10	10	1	1.5	1.5	4.9	4.7	4.7	4.9	1.5	1.6
30	10	1	1.4	1.5	4.2	4.3	4.2	4.0	1.4	1.4
50	10	1	1.5	1.4	4.5	4.1	4.0	4.0	1.4	1.4
100	20	1	1.4	1.3	2.6	2.5	2.5	2.4	1.3	1.2
500	50	1	1.2	1.2	1.7	1.7	1.7	1.7	1.2	1.3

## 298 7 Discussion

299 This paper introduced an explicit expression of RPG for rectangular robust MDPs. Our approach  
 300 involved auxiliary results such as deriving the worst model in closed form and showing that it is a  
 301 rank-one perturbation of the nominal kernel. The resulting RPG extends vanilla PG with additional  
 302 correction terms that can be derived in closed form as well. Thus, the computational time of RPG is  
 303 similar to its non-robust variant.

304 A key assumption that would be interesting to relax is the normed-ball structure of the uncertainty  
 305 sets considered in this study. Indeed, since the proofs of our technical results rely on norm properties,  
 306 it is still unclear if and how RPG can generalize to metric-based or  $f$ -divergence uncertainty sets.  
 307 The latter type of uncertainty can be particularly useful for data-driven settings, as the radius can be  
 308 chosen according to cross-validation or statistical bounds [8]. Another compelling direction would be  
 309 to explore other variants of RPG using mirror descent or natural policy gradient and examine their  
 310 compatibility with deep architectures, which would further demonstrate the practical efficiency of our  
 311 RPG method.

## 312 References

313 [1] Bahram Behzadian, Marek Petrik, and Chin Pang Ho. Fast algorithms for  $\ell_{\infty}$ -constrained  
 314  $s$ -rectangular robust MDPs. *Advances in Neural Information Processing Systems*, 34:25982–

- 315 25992, 2021.
- 316 [2] Rob Brekelmans, Tim Genewein, Jordi Grau-Moya, Grégoire Delétang, Markus Kunesch, Shane  
317 Legg, and Pedro Ortega. Your policy regularizer is secretly an adversary. *Transactions on*  
318 *Machine Learning Research (TMLR)*, 2022.
- 319 [3] Esther Derman, Matthieu Geist, and Shie Mannor. Twice regularized MDPs and the equivalence  
320 between robustness and regularization. *Advances in Neural Information Processing Systems*,  
321 34:22274–22287, 2021.
- 322 [4] Esther Derman, Daniel J. Mankowitz, Timothy A. Mann, and Shie Mannor. Soft-robust actor-  
323 critic policy-gradient. *AUAI press for Association for Uncertainty in Artificial Intelligence*,  
324 pages 208–218, 2018.
- 325 [5] Benjamin Eysenbach and Sergey Levine. Maximum entropy RL (provably) solves some robust  
326 RL problems. *International Conference on Learning Representations*, 2022.
- 327 [6] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-  
328 policy maximum entropy deep reinforcement learning with a stochastic actor. In *International*  
329 *Conference on Machine Learning*, pages 1861–1870. PMLR, 2018.
- 330 [7] Chin Pang Ho, Marek Petrik, and Wolfram Wiesemann. Partial policy iteration for  $l_1$ -robust  
331 Markov decision processes. *J. Mach. Learn. Res.*, 22:275–1, 2021.
- 332 [8] Chin Pang Ho, Marek Petrik, and Wolfram Wiesemann. Robust  $\phi$ -divergence MDPs. *arXiv*  
333 *preprint arXiv:2205.14202*, 2022.
- 334 [9] Hisham Husain, Kamil Ciosek, and Ryota Tomioka. Regularized policies are reward robust. In  
335 *International Conference on Artificial Intelligence and Statistics*, pages 64–72. PMLR, 2021.
- 336 [10] Garud N Iyengar. Robust dynamic programming. *Mathematics of Operations Research*,  
337 30(2):257–280, 2005.
- 338 [11] Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning.  
339 In *International Conference on Machine Learning*. Citeseer, 2002.
- 340 [12] Navdeep Kumar, Kfir Levy, Kaixin Wang, and Shie Mannor. Efficient policy iteration for robust  
341 markov decision processes via regularization, 2022.
- 342 [13] Kyungjae Lee, Sungjoon Choi, and Songhwai Oh. Sparse Markov decision processes with  
343 causal sparse Tsallis entropy regularization for reinforcement learning. *IEEE Robotics and*  
344 *Automation Letters*, 3(3):1466–1473, 2018.
- 345 [14] Yan Li, Tuo Zhao, and Guanghui Lan. First-order policy optimization for robust markov  
346 decision process. *arXiv preprint arXiv:2209.10579*, 2022.
- 347 [15] Daniel Mankowitz, Timothy Mann, Pierre-Luc Bacon, Doina Precup, and Shie Mannor. Learning  
348 robust options. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32,  
349 2018.
- 350 [16] Shie Mannor, Duncan Simester, Peng Sun, and John N Tsitsiklis. Bias and variance approxima-  
351 tion in value function estimates. *Management Science*, 53(2):308–322, 2007.
- 352 [17] Paul Milgrom and Ilya Segal. Envelope theorems for arbitrary choice sets. *Econometrica*,  
353 70:583–601, 02 2002.
- 354 [18] Arnab Nilim and Laurent El Ghaoui. Robust control of Markov decision processes with  
355 uncertain transition matrices. *Operations Research*, 53(5):780–798, 2005.
- 356 [19] Lerrel Pinto, James Davidson, Rahul Sukthankar, and Abhinav Gupta. Robust adversarial  
357 reinforcement learning. In *International Conference on Machine Learning*, pages 2817–2826.  
358 PMLR, 2017.
- 359 [20] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*.  
360 John Wiley & Sons, 2014.

- 361 [21] John Schulman, Xi Chen, and Pieter Abbeel. Equivalence between policy gradients and soft  
362 q-learning, 2017.
- 363 [22] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT  
364 Press, second edition, 2018.
- 365 [23] Richard S Sutton, David A McAllester, Satinder P Singh, Yishay Mansour, et al. Policy gradient  
366 methods for reinforcement learning with function approximation. In *Advances in Neural  
367 Information Processing Systems*, volume 99, pages 1057–1063. Citeseer, 1999.
- 368 [24] Aviv Tamar, Yinlam Chow, Mohammad Ghavamzadeh, and Shie Mannor. Policy gradient for  
369 coherent risk measures. *Advances in neural information processing systems*, 28, 2015.
- 370 [25] Aviv Tamar, Shie Mannor, and Huan Xu. Scaling up robust MDPs using function approximation.  
371 In *International Conference on Machine Learning*, pages 181–189. PMLR, 2014.
- 372 [26] Chen Tessler, Yonathan Efroni, and Shie Mannor. Action robust reinforcement learning and  
373 applications in continuous control. In *International Conference on Machine Learning*, pages  
374 6215–6224. PMLR, 2019.
- 375 [27] Qiu hao Wang, Chin Pang Ho, and Marek Petrik. On the convergence of policy gradient in  
376 robust mdps. *arXiv preprint arXiv:2212.10439*, 2022.
- 377 [28] Yue Wang, Fei Miao, and Shaofeng Zou. Robust constrained reinforcement learning. *arXiv  
378 preprint arXiv:2209.06866*, 2022.
- 379 [29] Yue Wang and Shaofeng Zou. Policy gradient method for robust reinforcement learning. In  
380 *International Conference on Machine Learning*, pages 23484–23526. PMLR, 2022.
- 381 [30] Wolfram Wiesemann, Daniel Kuhn, and Berç Rustem. Robust Markov decision processes.  
382 *Mathematics of Operations Research*, 38(1):153–183, 2013.
- 383 [31] Huan Xu and Shie Mannor. Robustness and generalization. *Machine learning*, 86:391–423,  
384 2012.

385	<b>Contents</b>	
386	<b>A Balanced and Normed Vectors</b>	<b>13</b>
387	A.1 Proof of Proposition 4.1 . . . . .	14
388	<b>B Worst Kernel and Reward</b>	<b>14</b>
389	B.1 Proof of Theorem 4.2 . . . . .	14
390	B.2 Proof of Theorem 4.3 . . . . .	15
391	<b>C Occupation Matrix</b>	<b>16</b>
392	C.1 Proof of Lemma 4.4 . . . . .	16
393	C.2 Proof of Theorem 4.6 . . . . .	17
394	<b>D Robust Q-value</b>	<b>17</b>
395	D.1 Basic Properties . . . . .	17
396	D.2 Evaluation . . . . .	18
397	D.3 Convergence . . . . .	19
398	<b>E Robust Policy Gradient</b>	<b>20</b>
399	<b>F Complexity Analysis</b>	<b>20</b>
400	F.1 Helper results for Robust MDPs . . . . .	21
401	F.2 Computing the Policy gradient. . . . .	22
402	<b>G Generalization to arbitrary norms</b>	<b>23</b>
403	G.1 $s$ -rectangular robust MDPs. . . . .	23
404	G.2 $s$ -rectangular Case . . . . .	23
405	G.3 Generalization to non-norms . . . . .	24
406	<b>H Experiments</b>	<b>24</b>
407	H.1 RPG by LP . . . . .	24

## 408 A Balanced and Normed Vectors

409 In this section, we lay down some basic properties of  $p$ -normalized-balanced vectors.

410 First recall the  $p$ -variance and the  $p$ -mean defined as:

$$\kappa_p(v) = \min_{\omega \in \mathbb{R}} \|v - \omega \mathbf{1}\|_p, \quad \omega_p(v) = \arg \min_{\omega \in \mathbb{R}} \|v - \omega \mathbf{1}\|_p.$$

411 Given any  $v \in \mathbb{R}^S$ , let also the  $p$ -balanced-normalized function:

$$u_p(v)(s) := \text{SIGN}(v(s) - \omega_p(v)) \left( \frac{|v(s) - \omega_p(v)|}{\kappa_p(v)} \right)^{p-1}, \quad \forall v \in \mathbb{R}^S, s \in S.$$

412 According to [12][Sec. 16.1, Lemma 1], the following holds:

$$\kappa_q(v) = -\frac{1}{\epsilon} \left[ \min_{\|c\|_p \leq \epsilon, \langle c, \mathbf{1} \rangle = 0} \langle c, v \rangle \right], \quad (6)$$

413 namely, the  $p$ -variance function is the optimal value of a linear optimization under kernel noise  
414 constraint. The result below further characterizes the solution to the above problem.

415 **Lemma A.1.** *The vector defined as  $c^* := -\epsilon u_q(v)$  is an optimal solution to the optimization problem*

$$\min_{\|c\|_p \leq \epsilon, \langle c, \mathbf{1} \rangle = 0} \langle c, v \rangle.$$

416 *Proof.* It suffices to show that  $c^*$  satisfies both constraints  $\|c^*\|_p \leq \epsilon$  and  $\langle c^*, \mathbf{1} \rangle = 0$ , and that it  
417 reaches optimal value, i.e.,  $-\frac{1}{\epsilon} \langle c^*, v \rangle = \kappa_q(v)$ . We thus compute:

$$\begin{aligned} \|c^*\|_p &= \left( \sum_{s \in S} |c^*(s)|^p \right)^{\frac{1}{p}} \\ &= \left( \sum_{s \in S} \left| -\epsilon \text{SIGN}(v(s) - \omega_q(v)) \left( \frac{|v(s) - \omega_q(v)|}{\kappa_q(v)} \right)^{q-1} \right|^p \right)^{\frac{1}{p}} \\ &= \left( \left( \frac{\epsilon}{\kappa_q(v)^{q-1}} \right)^p \sum_{s \in S} \left| \left( \frac{|v(s) - \omega_q(v)|}{\kappa_q(v)} \right)^{q-1} \right|^p \right)^{\frac{1}{p}} \\ &= \frac{\epsilon}{\kappa_q(v)^{q-1}} \left( \sum_{s \in S} |v(s) - \omega_q(v)|^{(q-1)p} \right)^{\frac{1}{p}} \\ &= \frac{\epsilon}{\kappa_q(v)^{q-1}} \left( \sum_{s \in S} |v(s) - \omega_q(v)|^q \right)^{\frac{1}{p}} && \text{(By assumption, } \frac{p+q}{pq} = 1) \\ &= \frac{\epsilon}{\kappa_q(v)^{q-1}} \kappa_q(v)^{\frac{q}{p}} && \text{(By definition, } \kappa_q(v) = \|v - \omega_q \mathbf{1}\|_q) \\ &= \epsilon, && \left( \frac{q}{p} - (q-1) = \frac{q-pq+p}{p} = 0 \right) \end{aligned}$$

418 so the norm constraint is satisfied. We check the noise constraint by computing:

$$\begin{aligned} \sum_{s \in S} c^*(s) &= \sum_{s \in S} -\epsilon \text{SIGN}(v(s) - \omega_q(v)) \left( \frac{|v(s) - \omega_q(v)|}{\kappa_q(v)} \right)^{q-1} \\ &= \frac{-\epsilon}{\kappa_q(v)^{q-1}} \sum_{s \in S} \text{SIGN}(v(s) - \omega_q(v)) |v(s) - \omega_q(v)|^{q-1}. \end{aligned}$$

419 Now, considering the real function  $\varphi : w \rightarrow \|v - w \mathbf{1}\|_q$  and taking its derivative, we remark the  
420 proportional relation:

$$\sum_{s \in S} c^*(s) = C \cdot \varphi'(\omega_q(v)),$$

421 where  $C \in \mathbb{R}$  is the proportionality coefficient. By construction,  $\omega_q(v)$  is a minimizer of  $\varphi$ , so we  
 422 must have  $\varphi'(\omega_q(v)) = 0$  and  $c^*$  satisfies the noise constraint.

423 We finally show that  $c^*$  reaches the optimal value:

$$\begin{aligned}
 -\frac{1}{\epsilon} \langle c^*, v \rangle &= -\frac{1}{\epsilon} \langle c^*, v - \omega_q(v) \mathbf{1} \rangle && (\langle c^*, \mathbf{1} \rangle_{\mathcal{S}} = 0) \\
 &= \sum_{s \in \mathcal{S}} \frac{|v(s) - \omega_q(v)|^q}{\kappa_q(v)^{q-1}} && \text{(Putting the value of } c^*) \\
 &= \frac{\kappa_q(v)^q}{\kappa_q(v)^{q-1}} && (\kappa_q(v) = \|v - \omega_q \mathbf{1}\|_q) \\
 &= \kappa_q(v).
 \end{aligned}$$

424

□

## 425 A.1 Proof of Proposition 4.1

426 **Proposition.** For any policy  $\pi \in \Pi$  and  $\ell_p$ -ball rectangular uncertainty set, the following holds:

$$\begin{aligned}
 u_{\mathcal{U}}^{\pi} &= \nabla_v \kappa_q(v) \Big|_{v=v_{\mathcal{U}}^{\pi}}, \\
 \langle u_{\mathcal{U}}^{\pi}, v_{\mathcal{U}}^{\pi} \rangle &= \kappa_q(v_{\mathcal{U}}^{\pi}).
 \end{aligned}$$

427 *Proof.* The second claim directly follows from Lemma A.1 applied to  $v := v_{\mathcal{U}}^{\pi}$ , so that by optimality,  
 428  $\kappa_q(v_{\mathcal{U}}^{\pi}) = \langle u_{\mathcal{U}}^{\pi}, v_{\mathcal{U}}^{\pi} \rangle$ . For the first claim, we take the gradient of  $\kappa_p(v) := \min_{w \in \mathbb{R}} \|v - w \mathbf{1}\|_p$  w.r.t.  $v$   
 429 using the envelope theorem [17]. Then, the  $p$ -balanced-normalized vector  $u_p(v)$  is a sub-gradient of  
 430  $\kappa_p(v)$ , that is,

$$u_p(v) = \nabla \kappa_q(v),$$

431 which we apply to  $v := v_{\mathcal{U}}^{\pi}$ .

□

432 We have the additional properties below:

- 433 • The variance function  $\kappa_q$  is translation-invariant in all-ones directions, i.e., for all  $\omega \in$   
 434  $\mathbb{R}$ ,  $\kappa_q(v) = \kappa_q(v + \omega \mathbf{1})$ . As a result,  $\langle \nabla \kappa_q(v), \mathbf{1} \rangle_{\mathcal{S}} = 0$ .
- 435 • The balanced-normalized vector  $u_p(v)$  has unit norm, i.e.,  $\|u_p(v)\|_p = 1$  by Lemma A.1.

## 436 B Worst Kernel and Reward

437 Here we present the proofs for the worst/adversarial kernel and reward function characterization.

### 438 B.1 Proof of Theorem 4.2

439 **Theorem** ( $(s, a)$ -rectangular case). Given uncertainty set  $\mathcal{U} = \mathcal{U}_p^{sa}$  and any policy  $\pi \in \Pi$ , the worst  
 440 model is related to the nominal one through:

$$R_{\mathcal{U}}^{\pi}(s, a) = R_0(s, a) - \alpha_{s,a} \quad \text{and} \quad P_{\mathcal{U}}^{\pi}(\cdot | s, a) = P_0(\cdot | s, a) - \beta_{s,a} u_{\mathcal{U}}^{\pi}.$$

441 *Proof.* By definition,

$$(P_{\mathcal{U}_p^{sa}}^{\pi}, R_{\mathcal{U}_p^{sa}}^{\pi}) \in \arg \min_{(P, R) \in \mathcal{U}_p^{sa}} T_{(P, R)}^{\pi} v_{\mathcal{U}_p^{sa}}^{\pi}.$$

442 Additionally, since  $\mathcal{U}_p^{sa} = (R_0 + \mathcal{R}) \times (P_0 + \mathcal{P})$ , it results that:

$$(R_{\mathcal{U}_p^{sa}}^{\pi}, P_{\mathcal{U}_p^{sa}}^{\pi}) = (P_0 + P^*, R_0 + R^*)$$

443 where

$$(P^*, R^*) \in \arg \min_{(P, R) \in \mathcal{P} \times \mathcal{R}} T_{(P, R)}^{\pi} v_{\mathcal{U}_p^{sa}}^{\pi}.$$

444 By the  $(s, a)$ -rectangularity assumption, we get that for all  $(s, a) \in \mathcal{X}$ ,

$$(P^*(\cdot|s, a), R^*(s, a)) \in \arg \min_{(p_{s,a}, r_{s,a}) \in \mathcal{P}_{s,a} \times \mathcal{R}_{s,a}} \left\{ r_{s,a} + \gamma \sum_{s' \in \mathcal{S}} p_{s,a}(s') v_{\mathcal{U}_p^\pi}^\pi(s') \right\}$$

445 It is clear from the above that the worst reward is independent of policy  $\pi$ . Thus, by the ball constraint,  
446 it is given by

$$R^*(s, a) = -\alpha_{s,a}, \quad \forall (s, a) \in \mathcal{X}.$$

447 Differently, the worst kernel depends on the value function which itself depends on the policy. It is  
448 given by

$$P^*(\cdot|s, a) = \arg \min_{p_{s,a} \in \mathcal{P}_{sa}} \left\{ \sum_{s' \in \mathcal{S}} p_{s,a}(s') v_{\mathcal{U}_p^\pi}^\pi(s') \right\}, \quad \forall (s, a) \in \mathcal{X}.$$

449 The optimization is of the form

$$\arg \min_{\|c\|_p \leq \beta, \langle c, \mathbf{1} \rangle = 0} \langle c, v \rangle,$$

450 so by Lemma A.1,

$$P^*(s'|s, a) = -\beta_{s,a} \text{SIGN} \left( v_{\mathcal{U}_p^\pi}^\pi(s') - \omega_q(v_{\mathcal{U}_p^\pi}^\pi) \right) \frac{|v_{\mathcal{U}_p^\pi}^\pi(s') - \omega_q(v_{\mathcal{U}_p^\pi}^\pi)|^{q-1}}{\kappa_q(v)^{q-1}}.$$

451 As a result, we proved that for all  $(s, a) \in \mathcal{X}$ ,  $R_{\mathcal{U}_p^\pi}^\pi(s, a) = R_0(s, a) - \alpha_{s,a}$  and

$$P_{\mathcal{U}_p^\pi}^\pi(s'|s, a) = P_0(s'|s, a) - \beta_{s,a} \text{SIGN} \left( v_{\mathcal{U}_p^\pi}^\pi(s') - \omega_q(v_{\mathcal{U}_p^\pi}^\pi) \right) \frac{|v_{\mathcal{U}_p^\pi}^\pi(s') - \omega_q(v_{\mathcal{U}_p^\pi}^\pi)|^{q-1}}{\kappa_q(v)^{q-1}}.$$

452

□

## 453 B.2 Proof of Theorem 4.3

454 **Theorem** ( $s$ -rectangular case). *Given uncertainty set  $\mathcal{U} = \mathcal{U}_p^s$  and any policy  $\pi \in \Pi$ , the worst model*  
455 *is related to the nominal one through:*

$$R_{\mathcal{U}}^\pi(s, a) = R_0(s, a) - \alpha_s \left( \frac{\pi_s(a)}{\|\pi_s\|_q} \right)^{q-1} \quad \text{and} \quad P_{\mathcal{U}}^\pi(\cdot|s, a) = P_0(\cdot|s, a) - \beta_s u_{\mathcal{U}}^\pi \left( \frac{\pi_s(a)}{\|\pi_s\|_q} \right)^{q-1}.$$

456 *Proof.* By definition,

$$(P_{\mathcal{U}_p^\pi}^\pi, R_{\mathcal{U}_p^\pi}^\pi) \in \arg \min_{(P, R) \in \mathcal{U}_p^\pi} T_{(P, R)}^\pi v_{\mathcal{U}_p^\pi}^\pi,$$

457 and since  $\mathcal{U}_p^s = (R_0 + \mathcal{R}) \times (P_0 + \mathcal{P})$ , we have

$$(P_{\mathcal{U}_p^s}^\pi, R_{\mathcal{U}_p^s}^\pi) = (P_0 + P^*, R_0 + R^*)$$

458 where

$$(P^*, R^*) \in \arg \min_{(P, R) \in \mathcal{P} \times \mathcal{R}} T_{(P, R)}^\pi v_{\mathcal{U}_p^\pi}^\pi.$$

459 By the  $s$ -rectangularity assumption, we get that for all  $s \in \mathcal{S}$

$$(P^*(\cdot|s, \cdot), R^*(s, \cdot)) = \arg \min_{(p_s, r_s) \in \mathcal{P}_s \times \mathcal{R}_s} \sum_{a \in \mathcal{A}} \pi_s(a) \left\{ r_{s,a} + \gamma \sum_{s' \in \mathcal{S}} p_{s,a}(s') v_{\mathcal{U}_p^\pi}^\pi(s') \right\}.$$

460 Here, the worst reward does depend on policy  $\pi$  and is given by

$$R^*(s, a) = -\alpha_s \frac{\pi_s(a)^{q-1}}{\sum_a \pi_s(a)^{q-1}}, \quad \forall (s, a) \in \mathcal{X}.$$

461 As for the worst kernel, it depends both on the value function and the policy. It is given by

$$P^*(\cdot|s, \cdot) = \arg \min_{p_s \in \mathcal{P}_s} \left\{ \sum_{a \in \mathcal{A}} \pi_s(a) \sum_{s' \in \mathcal{S}} p_{s,a}(s') v_{\mathcal{U}_p^\pi}^\pi(s') \right\}.$$

462 The optimization of interest is of the form

$$\min_{\|c_a\|_p \leq \beta_s, \langle c_a, \mathbf{1} \rangle = 0, a \in \mathcal{A}} \left\{ \sum_{a' \in \mathcal{A}} \pi_s(a') \langle c_{a'}, v \rangle \right\},$$

463 which is equivalent to the following two-fold minimization:

$$\min_{\sum_{a \in \mathcal{A}} (\beta_{s,a})^p \leq (\beta_s)^p} \min_{\|c_a\|_p \leq \beta_s, \langle c_a, \mathbf{1} \rangle = 0, a \in \mathcal{A}} \left\{ \sum_{a' \in \mathcal{A}} \pi_s(a') \langle c_{a'}, v \rangle \right\}.$$

464 Thus, rewriting the problem in our context,

$$\begin{aligned} & \min_{\sum_a (\beta_{s,a})^p \leq (\beta_s)^p} \min_{\|p_{sa}\|_p \leq \beta_s, \sum_{s'} p_{sa}(s') = 0} \sum_a \pi_s(a) \langle p_{s,a}, v \rangle \\ &= \min_{\sum_a (\beta_{s,a})^p \leq (\beta_s)^p} \sum_a \pi_s(a) \min_{\|p_{sa}\|_p \leq \beta_s, \sum_{s'} p_{sa}(s') = 0} \langle p_{s,a}, v \rangle \\ &= \min_{\sum_a (\beta_{sa})^p \leq (\beta_s)^p} \sum_a \pi_s(a) (-\beta_{sa} \kappa_q(v)) \quad (\text{By Lemma A.1}) \\ &= -\kappa_q(v) \max_{\sum_a (\beta_{sa})^p \leq (\beta_s)^p} \sum_a \pi_s(a) \beta_{sa}. \end{aligned}$$

465 Computing the optimal  $\beta$  above, the optimization is now the same as in the  $(s, a)$ -rectangular case.

466 Hence, we have

$$P^*(s'|s, a) = -\beta_s \frac{\pi_s(a)^{q-1}}{\|\pi_s\|_q^{q-1}} \text{SIGN}(v_{\mathcal{U}_p^\pi}^\pi(s') - \omega_q(v_{\mathcal{U}_p^\pi}^\pi)) \frac{|v_{\mathcal{U}_p^\pi}^\pi(s') - \omega_q(v_{\mathcal{U}_p^\pi}^\pi)|^{q-1}}{\kappa_q(v)^{q-1}},$$

467 which ends the proof by definition of the balanced value function  $u_{\mathcal{U}}^\pi$ .  $\square$

## 468 C Occupation Matrix

### 469 C.1 Proof of Lemma 4.4

470 **Lemma.** Let  $b, k \in \mathbb{R}^S$  and  $P_0, P_1 \in (\Delta_S)^S$  two transition matrices. If  $P_1 = P_0 - bk^\top$ , i.e.,  $P_1$   
471 is a rank-one perturbation of  $P_0$ , then their occupation matrices  $D_i := (I - \gamma P_i)^{-1}$ ,  $i = 0, 1$  are  
472 related through:

$$D_1 = D_0 - \gamma \frac{D_0 b k^\top D_0}{(1 + \gamma k^\top D_0 b)}.$$

473 *Proof.* By definition,  $D_1 = (I_S - \gamma P_1)^{-1}$  so it follows that:

$$\begin{aligned} & (I_S - \gamma P_1) D_1 = I_S \\ & \iff I_S + \gamma P_1 D_1 = D_1 \\ & \iff I_S + \gamma (P_0 - b k^\top) D_1 = D_1 \quad (\text{By assumption, } P_1 = P_0 - b k^\top) \\ & \iff I_S - \gamma b k^\top D_1 = (I_S - \gamma P_0) D_1 \\ & \iff (I_S - \gamma P_0)^{-1} (I_S - \gamma b k^\top D_1) = D_1 \quad (\text{Multiplying both sides by } (I_S - \gamma P_0)^{-1}) \\ & \iff D_0 (I_S - \gamma b k^\top D_1) = D_1 \quad (\text{By definition, } D_0 = (I_S - \gamma P_0)^{-1}) \\ & \iff D_0 - \gamma D_0 b k^\top D_1 = D_1. \quad (7) \end{aligned}$$

474 Now, multiplying both sides by  $k$  and noticing that  $k^\top D_0 b$  is a scalar we get

$$\begin{aligned} & k^\top D_0 - \gamma k^\top D_0 b k^\top D_1 = k^\top D_1 \\ \iff & k^\top D_0 = (1 + \gamma k^\top D_0 b) k^\top D_1 \\ \iff & k^\top D_1 = \frac{k^\top D_0}{(1 + \gamma k^\top D_0 b)}. \end{aligned} \quad (8)$$

475 Combining Eqs. (7) and (8) thus yields:

$$D_1 = D_0 - \gamma \frac{D_0 b k^\top D_0}{(1 + \gamma k^\top D_0 b)},$$

476 which concludes the proof.  $\square$

## 477 C.2 Proof of Theorem 4.6

478 **Theorem.** For any rectangular  $\ell_p$ -ball-constrained uncertainty and  $\pi \in \Pi$ , it holds that:

$$d_{\mathcal{U}, \mu}^\pi = d_{P_0, \mu}^\pi - \gamma \frac{\langle d_{P_0, \mu}^\pi, \beta^\pi \rangle_{\mathcal{S}}}{1 + \gamma \langle d_{P_0, u_{\mathcal{U}}^\pi}, \beta^\pi \rangle_{\mathcal{S}}} d_{P_0, u_{\mathcal{U}}^\pi}^\pi.$$

479 *Proof.* From Thms. 4.2 and 4.3, it holds that:

$$P_{\mathcal{U}}^\pi(s'|s) = P_0^\pi(s'|s) - \beta_s^\pi u_{\mathcal{U}}^\pi(s'), \quad \forall s, s' \in \mathcal{S}.$$

480 Therefore, setting  $P_1 := P_{\mathcal{U}}^\pi$ ,  $P_0 := P_0^\pi$ ,  $b := \beta^\pi$  and  $k := u_{\mathcal{U}}^\pi$ , we can apply Lemma 4.4 and  
481 relate the corresponding occupation matrices. Additionally multiplying both sides of the relation by  
482  $\mu^\top \in \mathbb{R}^{1 \times \mathcal{S}}$  yields the desired result.  $\square$

## 483 D Robust Q-value

### 484 D.1 Basic Properties

485 In the literature, robust Q-values are defined in various ways that turn out to be conflicting for  $s$   
486 but non- $(s, a)$  rectangular uncertainty sets. In this section, we propose to define the robust Q-value  
487 solely based on the worst model. Define the robust Q-value, the robust value function, and the robust  
488 occupation respectively as:

$$Q_{\mathcal{U}}^\pi := Q_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}, \quad d_{\mathcal{U}}^\pi := d_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}, \quad v_{\mathcal{U}}^\pi := v_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}.$$

489 For  $s$ -rectangular uncertainty sets (in particular, for  $(s, a)$ -rectangular), the above definition of robust  
490 value function coincides with the common one, i.e.,  $v_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}^\pi = \min_{(P, R) \in \mathcal{U}} v_{(P, R)}^\pi$  [30]. If the  
491 uncertainty set is additionally  $(s, a)$ -rectangular (as in [28] or [3, 12]), the above definition of robust  
492 Q-value also coincides with the common one because then,

$$Q_{\mathcal{U}^{\text{sa}}}^\pi(s, a) = \min_{(P, R) \in \mathcal{U}^{\text{sa}}} \left( R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) v_{\mathcal{U}^{\text{sa}}}^\pi(s') \right), \quad \forall (s, a) \in \mathcal{X}.$$

493 Getting back to our own definition, robust Q-value and value functions are related through:

$$\begin{aligned} v_{\mathcal{U}}^\pi(s) &= \langle \pi_s, Q_{\mathcal{U}}^\pi(s, \cdot) \rangle_{\mathcal{A}} \\ Q_{\mathcal{U}}^\pi(s, a) &= R_{\mathcal{U}}^\pi(s, a) + \gamma \sum_{s' \in \mathcal{S}} \pi_s(a) P_{\mathcal{U}}^\pi(s'|s, a) v_{\mathcal{U}}^\pi(s'), \end{aligned}$$

494 as both quantities are defined based on worst kernel and reward, i.e.,  $Q_{\mathcal{U}}^\pi := Q_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}$  and  $v_{\mathcal{U}}^\pi :=$   
495  $v_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}$ .

496 Given an optimal robust policy  $\pi_{\mathcal{U}}^*$ , we further use  $P_{\mathcal{U}}^*, R_{\mathcal{U}}^*, v_{\mathcal{U}}^*, Q_{\mathcal{U}}^*, d_{\mathcal{U}}^*$  as a shorthand for  
497  $P_{\mathcal{U}}^{\pi_{\mathcal{U}}^*}, R_{\mathcal{U}}^{\pi_{\mathcal{U}}^*}, v_{\mathcal{U}}^{\pi_{\mathcal{U}}^*}, Q_{\mathcal{U}}^{\pi_{\mathcal{U}}^*}, d_{\mathcal{U}}^{\pi_{\mathcal{U}}^*}$  respectively. For  $(s, a)$ -rectangular uncertainty set  $\mathcal{U}^{\text{sa}}$ , the optimal value  
498 function is the best optimal Q-value, that is

$$v_{\mathcal{U}^{\text{sa}}}^*(s) = \max_{a \in \mathcal{A}} Q_{\mathcal{U}^{\text{sa}}}^*(s, a), \quad \forall s \in \mathcal{S}.$$

499 because an optimal policy deterministically takes the action with the highest Q-value [18, 10]. This  
 500 does no longer hold for  $s$ -rectangular or coupled uncertainty sets, as there, an optimal policy may be  
 501 stochastic [30]. Still, based on Thms. 4.2 and 4.3, we get the Bellman recursion below.

502 **Proposition D.1.** *Let an  $\ell_p$ -ball constrained uncertainty set. Then, for all  $(s, a) \in \mathcal{X}$  and  $\pi \in \Pi$ , the*  
 503 *robust Q-value satisfies the following recursion in the  $(s, a)$  and  $s$ -rectangular case respectively:*

$$\begin{aligned} Q_{\mathcal{U}_p^{sa}}^\pi(s, a) &= T_{(P_0, R_0)}^\pi Q_{\mathcal{U}_p^{sa}}^\pi(s, a) - \alpha_{sa} - \gamma \beta_{sa} \kappa_q(v_{\mathcal{U}_p^{sa}}^\pi), \\ Q_{\mathcal{U}_p^s}^\pi(s, a) &= T_{(P_0, R_0)}^\pi Q_{\mathcal{U}_p^s}^\pi(s, a) - \left( \frac{\pi_s(a)}{\|\pi_s\|_q} \right)^{q-1} \left( \alpha_s + \gamma \beta_s \kappa_q(v_{\mathcal{U}_p^s}^\pi) \right). \end{aligned}$$

504 *Proof.* We give proof for the  $(s, a)$ -rectangular case only. The  $s$ -rectangular case follows the exact  
 505 same lines except that it uses Thm. 4.3 instead of Thm. 4.2. We have:

$$\begin{aligned} Q_{\mathcal{U}}^\pi(s, a) &= Q_{(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)}^\pi(s, a) && \text{(By definition)} \\ &= R_{\mathcal{U}}^\pi(s, a) + \sum_{s' \in \mathcal{S}} P_{\mathcal{U}}^\pi(s'|s, a) v_{\mathcal{U}_p^{sa}}^\pi(s') \\ &= R_0(s, a) - \alpha_{sa} + \gamma \sum_{s' \in \mathcal{S}} \left( P_0(s'|s, a) - \beta_{sa} u_{\mathcal{U}_p^{sa}}^\pi(s') \right) v_{\mathcal{U}_p^{sa}}^\pi(s') && \text{(By Thm. 4.2)} \\ &= R_0(s, a) - \alpha_{sa} + \gamma \sum_{s' \in \mathcal{S}} P_0(s'|s, a) v_{\mathcal{U}_p^{sa}}^\pi(s') - \gamma \beta_{sa} \kappa_q(v_{\mathcal{U}_p^{sa}}^\pi) && \text{(2d statement of Prop. 4.1)} \\ &= R_0(s, a) + \gamma \sum_{s', a'} P_0(s'|s, a) \pi_{s'}(a') Q_{\mathcal{U}_p^{sa}}^\pi(s', a') - \alpha_{sa} - \gamma \beta_{sa} \kappa_q(v_{\mathcal{U}_p^{sa}}^\pi) \\ &= T_{(P_0, R_0)}^\pi Q_{\mathcal{U}_p^{sa}}^\pi(s, a) - \alpha_{sa} - \gamma \beta_{sa} \kappa_q(v_{\mathcal{U}_p^{sa}}^\pi). \end{aligned}$$

506

□

507 The above recursion applies the standard Bellman operator on robust Q-values. We can similarly  
 508 apply it on the robust value function (itself can be computed efficiently based on [3, 12]).

509 **Corollary D.2.** *Let an  $\ell_p$ -ball constrained uncertainty set. Then, for all  $(s, a) \in \mathcal{X}$  and  $\pi \in \Pi$ , the*  
 510 *robust Q-value satisfies the following recursion in the  $(s, a)$  and  $s$ -rectangular case respectively:*

$$\begin{aligned} Q_{\mathcal{U}_p^{sa}}^\pi(s, a) &= R_0(s, a) + \gamma \sum_{s'} P_0(s'|s, a) v_{\mathcal{U}_p^{sa}}^\pi(s') - \alpha_{sa} - \gamma \beta_{sa} \kappa_q(v_{\mathcal{U}_p^{sa}}^\pi), \\ Q_{\mathcal{U}_p^s}^\pi(s, a) &= R_0(s, a) + \gamma \sum_{s'} P_0(s'|s, a) v_{\mathcal{U}_p^s}^\pi(s') - \left( \frac{\pi_s(a)}{\|\pi_s\|_q} \right)^{q-1} \left( \alpha_s + \gamma \beta_s \kappa_q(v_{\mathcal{U}_p^s}^\pi) \right). \end{aligned}$$

## 511 D.2 Evaluation

512 Based on the Bellman recursion above, we now derive robust Q-learning equations to learn a robust  
 513 Q-value. Precisely, we investigate if the linear operator below is contracting and can be evaluated  
 514 efficiently:

$$(\mathcal{L}_{\mathcal{U}}^\pi Q)(s, a) := R_{\mathcal{U}, v}^\pi(s, a) + \gamma \sum_{(s', a') \in \mathcal{X}} P_{\mathcal{U}, v}^\pi(s'|s, a) \pi_{s'}(a') Q(s', a'), \quad \forall Q \in \mathbb{R}^{\mathcal{X}}, \quad (9)$$

515 where  $(P_{\mathcal{U}, v}^\pi, R_{\mathcal{U}, v}^\pi) \in \arg \min_{(P, R) \in \mathcal{U}} T_{(P, R)}^\pi v$  and  $v(s) = \langle \pi_s, Q(s, \cdot) \rangle_{\mathcal{A}}$ ,  $\forall s \in \mathcal{S}$ .

516 **Proposition D.3.** *Consider an  $\ell_p$ -ball constrained uncertainty set. Then, for all  $Q \in \mathbb{R}^{\mathcal{X}}$  and  $\pi \in \Pi$ ,*  
 517 *the operator  $\mathcal{L}^\pi$  can be evaluated as:*

$$\begin{aligned} (\mathcal{L}_{\mathcal{U}_p^{sa}}^\pi Q)(s, a) &= T_{(P_0, R_0)}^\pi Q(s, a) - \alpha_{sa} - \gamma \beta_{sa} \kappa_q(v), \\ (\mathcal{L}_{\mathcal{U}_p^s}^\pi Q)(s, a) &= T_{(P_0, R_0)}^\pi Q(s, a) - \left( \frac{\pi_s(a)}{\|\pi_s\|_q} \right)^{q-1} \left( \alpha_s + \gamma \beta_s \kappa_q(v) \right), \end{aligned}$$

518 where for all  $Q \in \mathbb{R}^{\mathcal{X}}$ , its corresponding value is  $v(s) := \langle \pi_s, Q(s, \cdot) \rangle$ ,  $\forall s \in \mathcal{S}$ .

519 *Proof.* We give proof for the  $(s, a)$ -rectangular case only. The  $s$ -rectangular case follows the exact  
 520 same lines except that we take the worst model for  $s$ -rectangular balls. By definition,

$$\begin{aligned}
 (\mathcal{L}_{\mathcal{U}_p^{\text{sa}}}^\pi Q)(s, a) &= \min_{(P, R) \in \mathcal{U}_p^{\text{sa}}} \left\{ R(s, a) + \gamma \sum_{(s', a') \in \mathcal{X}} P(s'|s, a) \pi_{s'}(a') Q(s', a') \right\} \\
 &= \min_{R \in \mathcal{R}_p^{\text{sa}}} R(s, a) + \gamma \min_{P \in \mathcal{P}_p^{\text{sa}}} \left\{ \sum_{s' \in \mathcal{S}} P(s'|s, a) v(s') \right\} \\
 &= R_0(s, a) - \alpha_{s, a} + \gamma \sum_{s' \in \mathcal{S}} P_0(s'|s, a) v(s') - \beta_{s, a} \kappa_q(v) \quad (\text{By [12]}) \\
 &= R_0(s, a) - \alpha_{s, a} + \gamma \sum_{(s', a') \in \mathcal{X}} P_0(s'|s, a) \pi_{s'}(a') Q(s', a') - \beta_{s, a} \kappa_q(v) \\
 &= T_{(P_0, R_0)}^\pi Q(s, a) - \alpha_{s, a} - \beta_{s, a} \kappa_q(v).
 \end{aligned}$$

521

□

### 522 D.3 Convergence

523 In the rest of this section, we focus on  $\ell$ -ball constrained uncertainty sets of the form  $\mathcal{U}_p^{\text{sa}}$  or  $\mathcal{U}_p^{\text{s}}$ . Let  
 524 our Q-value iteration  $Q_{n+1} := \mathcal{L}_{\mathcal{U}}^\pi Q_n$ , and denote  $v_n(s) = \langle \pi_s, Q_n(s, \cdot) \rangle_{\mathcal{A}}, \forall s \in \mathcal{S}, n \in \mathbb{N}$ .

525 **Proposition D.4.** For all  $Q \in \mathbb{R}^{\mathcal{X}}$ , denote  $v(s) := \langle \pi_s, Q(s, \cdot) \rangle_{\mathcal{A}}, \forall s \in \mathcal{S}$ . Then, for any policy  
 526  $\pi \in \Pi$ , the Q-value iteration defined according to  $Q_{n+1} = \mathcal{L}_{\mathcal{U}}^\pi Q_n$  induces

$$v_{n+1} := \mathcal{T}_{\mathcal{U}}^\pi v_n.$$

527 *Proof.* By construction, for all  $s \in \mathcal{S}$  we have

$$\begin{aligned}
 v_{n+1}(s) &= \langle \pi_s, Q_{n+1}(s, \cdot) \rangle_{\mathcal{A}} \\
 &= \langle \pi_s, (\mathcal{L}_{\mathcal{U}}^\pi Q_n)(s, \cdot) \rangle_{\mathcal{A}} \\
 &= \sum_{a \in \mathcal{A}} \pi_s(a) \left[ R_{\mathcal{U}, v_n}^\pi(s, a) + \gamma \sum_{(s', a') \in \mathcal{X}} P_{\mathcal{U}, v_n}^\pi(s'|s, a) \pi_{s'}(a') Q_n(s', a') \right] \quad (\text{By Eq. 9}) \\
 &= \sum_{a \in \mathcal{A}} \pi_s(a) \left[ R_{\mathcal{U}, v_n}^\pi(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_{\mathcal{U}, v_n}^\pi(s'|s, a) v_n(s') \right] \quad (\text{By definition of } v_n) \\
 &= (\mathcal{T}_{\mathcal{U}}^\pi v_n)(s),
 \end{aligned}$$

528 where the last equality holds because  $(P_{\mathcal{U}, v_n}^\pi, R_{\mathcal{U}, v_n}^\pi) \in \arg \min_{(P, R) \in \mathcal{U}} \mathcal{T}_{(P, R)}^\pi v_n$ . □

529 As a result of the above proposition, the value iteration induced by our Q-value iteration rule converges  
 530 linearly to the robust value function, i.e.,  $\|v_n - v_{\mathcal{U}}^\pi\|_\infty \leq \gamma^n \|v_0\|_\infty$ . Therefore, Q-value iterates  
 531 converge to a fixed point. Precisely,  $v_n \rightarrow_n v_{\mathcal{U}}^\pi$  implies that  $(P_{\mathcal{U}, v_n}^\pi, R_{\mathcal{U}, v_n}^\pi) \rightarrow_n (P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi)$ , which in  
 532 turn implies that  $Q_n \rightarrow_n Q_{\mathcal{U}}^\pi$ . The result below further characterizes the convergence rate.

533 **Proposition D.5.** The Q-value iteration  $Q_{n+1} := \mathcal{L}_{\mathcal{U}}^\pi Q_n$  converges linearly to  $Q_{\mathcal{U}}^\pi$  for uncertainty  
 534 set  $\mathcal{U} = \mathcal{U}_p^{\text{sa}}, \mathcal{U}_p^{\text{s}}$  for every policy  $\pi \in \Pi$ .

*Proof.*

$$\begin{aligned}
 \|Q_{n+1} - Q_{\mathcal{U}}^\pi\|_\infty &= \|R_{\mathcal{U}, v_n}^\pi + \gamma P_{\mathcal{U}, v_n}^\pi v_n - R_{\mathcal{U}}^\pi + \gamma P_{\mathcal{U}}^\pi v_{\mathcal{U}}^\pi\|_\infty, \\
 &= \gamma \|P_{\mathcal{U}, v_n}^\pi v_n - P_{\mathcal{U}}^\pi v_{\mathcal{U}}^\pi\|_\infty, \quad (\text{as } R_{\mathcal{U}, v}^\pi = R_{\mathcal{U}}^\pi, \quad \forall v), \\
 &= \gamma \|(P_0 - B^\pi u_n) v_n - (P_0 - B^\pi u_{\mathcal{U}}^\pi) v_{\mathcal{U}}^\pi\|_\infty, \quad (\text{as } R_{\mathcal{U}, v}^\pi = R_{\mathcal{U}}^\pi, \quad \forall v),
 \end{aligned}$$

535 where  $B^\pi(s, a) = \beta_{s, a}$  for  $\mathcal{U} = \mathcal{U}_p^{\text{sa}}$  and  $B^\pi(s, a) = \beta_s \left( \frac{\pi_s(a)}{\|\pi_s\|_q} \right)^{q-1}$  for  $\mathcal{U} = \mathcal{U}_p^{\text{s}}$ . The equality  
 536 comes from the worst kernel characterization. This implies

$$\begin{aligned}
\|Q_{n+1} - Q_{\mathcal{U}}^{\pi}\|_{\infty} &\leq \gamma \|P_0(v_n - v_{\mathcal{U}}^{\pi})\|_{\infty} + \gamma \|B^{\pi}(u_n)^{\top} v_n - B^{\pi}(u_{\mathcal{U}}^{\pi})^{\top} v_{\mathcal{U}}^{\pi}\|_{\infty}, \\
&\leq \gamma^{n+1} \|v_0 - v_{\mathcal{U}}^{\pi}\|_{\infty} + \gamma \|(u_n)^{\top} v_n - (u_{\mathcal{U}}^{\pi})^{\top} v_{\mathcal{U}}^{\pi}\|, \quad (\text{as } B^{\pi}(s, a) \leq 1), \\
&\leq \gamma^{n+1} \|v_0 - v_{\mathcal{U}}^{\pi}\|_{\infty} + \gamma \|(u_n)^{\top} (v_n - v_{\mathcal{U}}^{\pi})\| + \gamma \|(u_n - u_{\mathcal{U}}^{\pi})^{\top} v_{\mathcal{U}}^{\pi}\|, \\
&\leq \gamma^{n+1} \|v_0 - v_{\mathcal{U}}^{\pi}\|_{\infty} + \gamma^{n+1} S \|v_0 - v_{\mathcal{U}}^{\pi}\|_{\infty} + \gamma \frac{\|\mathcal{R}\|_{\infty}}{1 - \gamma} \|u_n - u_{\mathcal{U}}^{\pi}\|_{\infty}.
\end{aligned}$$

537 Here,  $u_n, u_{\mathcal{U}}^{\pi}$  is the balanced normalized vector associated with vector  $v_n$  and  $v_{\mathcal{U}}^{\pi}$  respectively. Recall,  
538 the  $p$ -balanced normalized vector  $u$  associated with vector  $v$  is given by

$$u(s) := \frac{\text{sign}(v(s) - \omega_q(v)) \|v(s) - \omega_q(v)\|^{q-1}}{\kappa_q(v)^{q-1}}, \quad (10)$$

539 where  $\kappa_p(v) = \min_w \|v - w\mathbf{1}\|_p$  and  $\omega_p(v) = w \|v - w\mathbf{1}\|_p$ . It is easy to see that  $\omega_p, \kappa$  are  
540 Lipschitz function in  $v$ . Hence,  $\|u_n - u_{\mathcal{U}}^{\pi}\|_{\infty} \leq CPol(\|v_n - v_{\mathcal{U}}^{\pi}\|_{\infty}) f(\kappa(v_{\mathcal{U}}^{\pi}), S, A)$ , where  $Pol, f$   
541 is some polynomial and some function. This implies,

$$\|Q_{n+1} - Q_{\mathcal{U}}^{\pi}\|_{\infty} \leq \gamma^{n+1} f(\kappa(v_{\mathcal{U}}^{\pi}), \|v_0 - v_{\mathcal{U}}^{\pi}\|_{\infty}, S, A).$$

542 This concludes our proof.  $\square$

## 543 E Robust Policy Gradient

544 **Theorem (RPG).** For any rectangular  $\ell_p$ -ball-constrained uncertainty, the robust policy gradient is  
545 given by:

$$\partial_{\pi} \rho_{\mathcal{U}}^{\pi} = \sum_{(s,a) \in \mathcal{X}} (d_{P_0, \mu}^{\pi}(s) - c^{\pi}(s)) Q_{\mathcal{U}}^{\pi}(s, a) \nabla \pi_s(a), \quad (11)$$

546 where

$$c^{\pi}(s) := \frac{\gamma \langle d_{P_0, \mu}^{\pi}, \beta^{\pi} \rangle_S}{1 + \gamma \langle d_{P_0, u_{\mathcal{U}}^{\pi}}^{\pi}, \beta^{\pi} \rangle_S} d_{P_0, u_{\mathcal{U}}^{\pi}}^{\pi}(s), \quad \forall s \in \mathcal{S}.$$

547 *Proof.* The proof directly follows from plugging the robust occupation measure of Thm. 4.6 and the  
548 robust Q-value into standard policy-gradient theorem [23].  $\square$

## 549 F Complexity Analysis

550 In this section, we study the iteration complexity to compute robust policy gradient different uncer-  
551 tainty set.

552 **Convex non-rectangular Uncertainty set.** Robust policy improvement are strongly NP Hard for  
553 non-rectangular uncertainty set, even if it is convex [30]. The policy gradient method finds global  
554 optimal given oracle access to policy gradient in polynomial time [27]. This implies computation of  
555 policy gradient must be of NP Hard.

556 **Non-robust MDPs.** For non-robust case, computation of Q-value and occupation is  $O(S^2 A \log(\epsilon^{-1}))$   
557 each, which are most costly computations. Computing the product of  $d^{\pi}, Q^{\pi}$  and  $\nabla \pi$  as in policy  
558 gradient is  $O(SA)$  operation, which insignificant. Hence, the total cost for computing policy gradient  
559 is the same as cost of Q-value. More precisely, lets approximate Q-value with  $\hat{Q}$  and occupation with  
560  $\hat{d}$ , with  $\frac{\epsilon}{SA}$  tolerance, that is  $\|Q - \hat{Q}\|_{\infty}, \|d - \hat{d}\|_{\infty} \leq \frac{\epsilon}{SA}$ . This takes  $O(S^2 A \log(SA \epsilon^{-1}))$  each.  
561 Now then, we have

$$\sum_{s,a} d'(s) Q'(s, a) \nabla \pi_s(a) = \sum_{s,a} (Q^{\pi}(s, a) + \epsilon_1(s, a)) (d^{\pi}(s, a) + \epsilon_2(s, a)) \nabla \pi_s(a)$$

562 where  $\epsilon_1(s, a) = Q(s, a) - Q^{\pi}(s, a)$  and  $\epsilon_2(s, a) = d(s, a) - d^{\pi}(s, a)$ . We know, let  $B$  be the bound  
563 on  $\|Q^{\pi}\|_{\infty}, \|d^{\pi}\|_{\infty} \leq B$ . So now we have,

$$\begin{aligned} \sum_{s,a} d'(s)Q'(s,a)\nabla\pi_s(a) &= \sum_{s,a} (Q^\pi(s,a) + \epsilon_1(s,a))(d^\pi(s,a) + \epsilon_2(s,a))\nabla\pi_s(a) \\ &= \sum_{s,a} Q^\pi(s,a)d^\pi(s,a)\nabla\pi_s(a) + O(\epsilon). \end{aligned}$$

564 So the exact complexity of policy gradient for non-robust MDP is  $O(S^2A \log(SA\epsilon^{-1}))$ .

### 565 F.1 Helper results for Robust MDPs

566 **Computing variance and mean functions.** Computing  $\kappa_p(v)$  and  $\omega_p(v)$  to  $\epsilon$  tolerance requires  
567  $O(S \log(S\epsilon^{-1}))$  and  $O(S \log(\epsilon^{-1}))$  respectively, via binary search [12].

568 **Computing occupation measure.** Let  $k \in \mathbb{R}^S$  be any vector. From definition, we have

$$d_{P,k}^\pi = \sum_{n=0}^{\infty} \gamma^n k^\top (P^\pi)^n.$$

569 We have,

$$\|d_{P,k}^\pi - \sum_{n=0}^{N-1} \gamma^n k^\top (P^\pi)^n\| = \left\| \sum_{n=N}^{\infty} \gamma^n k^\top (P^\pi)^n \right\| \leq \|k\| \sum_{n=N}^{\infty} \|\gamma P^\pi\|^n.$$

570 Since,  $\|P^\pi\| \leq 1$  as it is a stochastic matrix, then

$$\|d_{P,k}^\pi - \sum_{n=0}^{N-1} \gamma^n k^\top (P^\pi)^n\| \leq \frac{\|k\| \gamma^N \|P^\pi\|^N}{1 - \gamma \|P^\pi\|} \leq \frac{\|k\| \gamma^N}{1 - \gamma}.$$

571 This implies  $\sum_{n=0}^{N-1} \gamma^n k^\top (P^\pi)^n$  is  $O(\gamma^N)$  approximation of  $d_{P,k}^\pi$ . Now, take  $u_0 = k$  and

$$u_{n+1} := \gamma(u_n)^\top P,$$

572 then  $\sum_{n=0}^{N-1} \gamma^n k^\top (P^\pi)^n = \sum_{n=0}^{N-1} u_n$ . And each iteration take  $O(S^2)$  iterations, leading to total  
573 cost  $O(S^2N)$  for  $N$  iterations. Computing  $P^\pi$  from  $P$  is  $O(S^2A)$ . We conclude computing the  
574 occupation measure has complexity of  $O(S^2A + S^2 \log(\epsilon^{-1}))$ .

575 **Lemma F.1.** We can approximate  $d_{P,k}^\pi$  with  $\sum_{n=0}^{N-1} \gamma^n (k')^\top (P^\pi)^n$  with complexity  $O(S^2A +$   
576  $S^2 \log(\epsilon^{-1}))$  with tolerance:

$$\|d_{P,k}^\pi - \sum_{n=0}^{N-1} \gamma^n (k')^\top (P^\pi)^n\| \leq O\left(\frac{\|k\| \gamma^N + \|k - k'\|}{1 - \gamma}\right)$$

*Proof.*

$$\begin{aligned} \|d_{P,k}^\pi - \sum_{n=0}^{N-1} \gamma^n (k')^\top (P^\pi)^n\| &\leq \|d_{P,k}^\pi - \sum_{n=0}^{N-1} \gamma^n (k)^\top (P^\pi)^n\| + \left\| \sum_{n=0}^{N-1} \gamma^n (k')^\top (P^\pi)^n - \sum_{n=0}^{N-1} \gamma^n (k)^\top (P^\pi)^n \right\| \\ &\leq O\left(\frac{\|k\| \gamma^N}{1 - \gamma}\right) + \left\| \sum_{n=0}^{N-1} \gamma^n (k)^\top (P^\pi)^n - \sum_{n=0}^{N-1} \gamma^n (k')^\top (P^\pi)^n \right\| \\ &\leq O\left(\frac{\|k\| \gamma^N}{1 - \gamma}\right) + \|k - k'\| \left\| \sum_{n=0}^{N-1} \gamma^n (P^\pi)^n - \sum_{n=0}^{N-1} \gamma^n (P^\pi)^n \right\| \\ &\leq O\left(\frac{\|k\| \gamma^N}{1 - \gamma}\right) + O\left(\frac{\|k - k'\|}{1 - \gamma}\right). \end{aligned}$$

577

□

578 **Computing Q-value given value function.** Let  $v$  be  $\epsilon_1$  approximation of robust value function  $v_{\mathcal{U}}^{\pi}$ ,  
 579 that is

$$\|v - v_{\mathcal{U}}^{\pi}\|_{\infty} \leq \epsilon_1.$$

580 We want to compute Q-value using the relation:

$$\begin{aligned} Q_{\mathcal{U}}^{\pi}(s, a) &= R_{\mathcal{U}}^{\pi}(s, a) + \sum_{s', a} \pi_s(a) P_{\mathcal{U}}^{\pi}(s' | s, a) v_{\mathcal{U}}^{\pi}(s') \\ &= R_0(s, a) + \gamma \sum_{s'} P_0(s' | s, a) v_{\mathcal{U}}^{\pi}(s') - \Omega_{\mathcal{U}}(v_{\mathcal{U}}^{\pi}, \pi). \end{aligned}$$

581 where  $\Omega_{\mathcal{U}_p^s}(v_{\mathcal{U}_p^s}^{\pi}, \pi) = \frac{\pi_s(a)^{q-1}}{\|\pi_s\|_q^{q-1}} (\alpha_s + \gamma \beta_s \kappa_q(v_{\mathcal{U}_p^s}^{\pi}))$  and  $\Omega_{\mathcal{U}_p^{sa}}(v_{\mathcal{U}_p^{sa}}^{\pi}, \pi) = \alpha_{sa} + \gamma \beta_{sa} \kappa_q(v_{\mathcal{U}_p^{sa}}^{\pi})$ . Let  
 582  $Q$  be approximated from  $v$  as

$$Q(s, a) = R_0(s, a) + \gamma \sum_{s'} P_0(s' | s, a) v(s') - \Omega_{\mathcal{U}}(v, \pi).$$

583 So we have,

$$\begin{aligned} \|Q_{\mathcal{U}}^{\pi}(s, a) - Q(s, a)\|_{\infty} &= \gamma \left\| \sum_{s'} P_0(s' | s, a) (v(s') - v_{\mathcal{U}}^{\pi}(s')) \right\| + \|\Omega_{\mathcal{U}}(v, \pi) - \Omega_{\mathcal{U}}(v_{\mathcal{U}}^{\pi}, \pi)\| \\ &\leq \gamma \epsilon_1 + \|\Omega_{\mathcal{U}}(v, \pi) - \Omega_{\mathcal{U}}(v_{\mathcal{U}}^{\pi}, \pi)\| \\ &\leq \gamma \epsilon_1 + \|\beta\|_{\infty} \|\kappa_q(v) - \kappa_q(v_{\mathcal{U}}^{\pi})\| \\ &\leq \gamma \epsilon_1 + \|\beta\|_{\infty} S^{\frac{1}{q}} \epsilon_1, \quad (\text{using lemma F.2}) \\ &= O(S^{\frac{1}{q}} \epsilon_1). \end{aligned}$$

584 This implies,  $\|Q - Q_{\mathcal{U}}^{\pi}\|_{\infty} \leq O(S^{\frac{1}{q}} \epsilon_1)$ .

585 **Lemma F.2.**  $\kappa_p$  is Lipschitz function, precisely

$$\|\kappa_p(v_1) - \kappa_p(v_2)\| \leq S^{\frac{1}{p}} \|v_1 - v_2\|_{\infty} \leq S^{\frac{1}{p}} \|v_1 - v_2\|_{\infty}.$$

586 *Proof.* Let  $w_i \in \arg \min_w \|v_i - w\mathbf{1}\|_p$ , then we have

$$\begin{aligned} \|\kappa_p(v_1) - \kappa_p(v_2)\| &= \kappa_p(v_1) - \kappa_p(v_2), \quad (\text{WLOG, assuming } \kappa_p(v_1) \geq \kappa_p(v_2)) \\ &= \min_w \|v_1 - w\mathbf{1}\|_p - \|v_2 - w\mathbf{1}\|_p, \quad (\text{From definition}) \\ &\leq \|v_1 - w_2\mathbf{1}\|_p - \|v_2 - w_2\mathbf{1}\|_p, \quad (\text{From definition of min operator}) \\ &\leq \|(v_1 - w_2\mathbf{1}) - (v_2 - w_2\mathbf{1})\|_p, \quad (\text{Reverse triangle inequality}) \\ &= \|v_1 - v_2\|_p \\ &= \left[ \sum_s (v_1(s) - v_2(s))^p \right]^{\frac{1}{p}} \leq S^{\frac{1}{p}} \|v_1 - v_2\|_{\infty}. \end{aligned}$$

587 □

588 **Lemma F.3.**  $Q_{\mathcal{U}_p^{sa}}^{\pi}$  can be approximated to  $\epsilon$  tolerance with the same complexity as complexity of  
 589 computing  $v_{\mathcal{U}_p^{sa}}^{\pi}$  to  $S^{-\frac{1}{q}} \epsilon$ .

590 *Proof.* Compute value function with  $S^{-\frac{1}{q}} \epsilon$  tolerance. The rest of operations are insignificant. Rest  
 591 follows from the above. □

## 592 F.2 Computing the Policy gradient.

593 Let  $O_p^{\text{sa}}(\epsilon)$  be the complexity to compute robust value function  $v_{\mathcal{U}_p^{\text{sa}}}^{\pi}$ , upto  $\epsilon$  tolerance, see [12] for  
 594 details. Calculate Q-value up to  $\epsilon_1$  tolerance which requires  $O_p^{\text{sa}}(S^{-\frac{1}{q}} \epsilon_1)$  from lemma F.3. Let  
 595  $d_1$  and  $d_2$  be  $\epsilon_2$  approximation of  $d_{P_0, \mu}^{\pi}$  and  $d_{P_0, k}^{\pi}$  respectively, which is insignificant compared to  
 596  $O_p^{\text{sa}}(S^{-\frac{1}{q}} \epsilon_2)$ . Now let's approximate the gradient with  $d_1, d_2, Q, \nabla \pi$  as in Theorem 5.1, which has a  
 597 complexity of  $O(SA)$ . Since the uncertainty set  $\mathcal{U}$  is compact, all the quantities are bounded. And  
 598 there are  $O(SA)$  operations in the Theorem 5.1, so taking  $\epsilon_1, \epsilon_2 = O(\frac{\epsilon}{SA})$ , we will get the  $O(\epsilon)$  of  
 599 the gradient. Hence, the total complexity is  $O_p^{\text{sa}}(S^{-\frac{1}{q}-1} A^{-1} \epsilon)$  which is  $\tilde{O}(S^2 A \log(\epsilon^{-1}))$ , by hiding  
 600 log factors, see [12]. A similar analysis follows for the  $s$ -rectangular case.

## 601 G Generalization to arbitrary norms

602 Here we focus on the generalization of our result to a general norm from the existing  $\ell_p$  norm. We do  
603 it case by case.

### 604 G.1 sa-rectangular robust MDPs.

605 Lets consider sa-rectangular uncertainty set  $\mathcal{U} = \mathcal{U}_{\|\cdot\|}^{\text{sa}}$  constrained by  $\|\cdot\|$  norm. Precisely, defined as

$$606 \quad \mathcal{U}_{\|\cdot\|}^{\text{sa}} = (P_0 + \mathcal{P}) \times (R_0 + \mathcal{R}), \quad \text{where} \quad (\mathcal{P}, \mathcal{R}) = (\times_{s,a} \mathcal{P}_{sa}, \times_{s,a} \mathcal{R}_{sa}),$$

$$607 \quad \mathcal{R}_{(s,a)} = \{r \in \mathbb{R} \mid \|r\| \leq \alpha_{s,a}\}, \quad \text{and} \quad \mathcal{P}_{(s,a)} = \{p \in \mathbb{R}^S \mid \langle p, \mathbf{1} \rangle_S = 0, \|p\| \leq \beta_{s,a}\}.$$

607 The robust Bellman operator  $\mathcal{T}_{\mathcal{U}}^\pi$  can be evaluated as

$$608 \quad (\mathcal{T}_{\mathcal{U}}^\pi v)(s) = \sum_a \pi_s(a) \left[ R(s, a) - \gamma \beta_{s,a} \kappa_{\|\cdot\|}(v) + \gamma \sum_{s'} P(s'|s, a) v(s') \right],$$

608 where variance function is defined as

$$\kappa_{\|\cdot\|}(v) := \min_{\langle u, \mathbf{1} \rangle_S = 0, \|u\| \leq 1} \langle u, v_{\mathcal{U}}^\pi \rangle.$$

609 This can be used to compute the robust value function. Then the worst values can found using robust  
610 Bellman operator  $\mathcal{T}_{\mathcal{U}}^\pi$  and robust value function  $v_{\mathcal{U}}^\pi$  as

$$(P_{\mathcal{U}}^\pi, R_{\mathcal{U}}^\pi) \in \arg \min_{(P,R) \in \mathcal{U}} \mathcal{T}_{(P,R)}^\pi v_{\mathcal{U}}^\pi, \quad [30].$$

611 It is easy to see that the worst values are given as

$$612 \quad R_{\mathcal{U}}^\pi(s, a) = R_0(s, a) - \alpha_{s,a} \quad \text{and} \quad P_{\mathcal{U}}^\pi(\cdot|s, a) = P_0^\pi(\cdot|s, a) - \beta_{s,a} u_{\mathcal{U}}^\pi,$$

612 where normalized-balanced value function  $u_{\mathcal{U}}^\pi$  is a solution to

$$\min_{\langle u, \mathbf{1} \rangle_S = 0, \|u\| \leq 1} \langle u, v_{\mathcal{U}}^\pi \rangle.$$

613 Observe that the worst kernel is still a rank-one perturbation of the nominal kernel. Hence, the robust  
614 occupation measure can be obtained using the Lemma 4.4 as

$$d_{\mathcal{U}, \mu}^\pi = d_{P_0, \mu}^\pi - \gamma \frac{\langle d_{P_0, \mu}^\pi, \beta^\pi \rangle_S}{1 + \gamma \langle d_{P_0, u_{\mathcal{U}}^\pi}^\pi, \beta^\pi \rangle_S} d_{P_0, u_{\mathcal{U}}^\pi}^\pi, \quad (12)$$

615 where  $\beta^\pi(s) = \sum_a \pi_s(a) \beta_{s,a}$ . The last ingredient to compute RPG is robust Q-value which can be  
616 computed using robust value function and worst values. However, it can be computed directly using  
617 the following iterates:

$$Q_{n+1}(s) = \min_{(P,R) \in \mathcal{U}} \left[ R(s, a) + \gamma \sum_{s'} P(s'|s, a) v(s') \right]$$

$$= R(s, a) - \alpha_{s,a} - \gamma \beta_{s,a} \kappa_{\|\cdot\|}(v) + \gamma \sum_{s'} P(s'|s, a) v(s'),$$

618 as  $Q_n$  converges to robust Q-value  $Q_{\mathcal{U}}^\pi$  linearly.

619 The proofs of the above claims are easy or similar to the  $\ell_p$  counterparts. Finally, computation of  
620 the variance function  $\kappa_{\|\cdot\|}$  and normalized-value function  $u_{\mathcal{U}}^\pi$  can be done via numerical convex  
621 optimization methods for general norms. However for the  $\ell_p$  case, they can be obtained in concrete  
622 forms, hence we choose it for the main presentation.

### 623 G.2 s-rectangular Case

624 Generalization of our methods to s-rectangular balls of a general norm is not straightforward and may  
625 not be possible for all kinds of norms. The crucial property of  $\ell_p$  norm that we exploited to prove  
626 rank-one perturbation is 'decoupling', that is, for  $x \in \mathbf{R}^{\mathcal{A} \times \mathcal{S}}$ ,

$$\|x\|_p^y = \sum_{a \in \mathcal{A}} \|x_a\|_w^z,$$

627 for some  $w, y, z$ . This holds for the  $\ell_p$  norm with  $w = y = z = p$ . We leave the further analysis of  
628 this setting for future work.

### 629 G.3 Generalization to non-norms

630 Further, generalization of our results to distance such as KL, can be tricky. The ability of our methods  
631 to compute RPG (particularly robust occupation measure) crucially relies on the rank-one perturbation  
632 result, which might not be the case for distance measures such as KL. We leave this analysis for  
633 future work.

## 634 H Experiments

635 **Parameters.** All the nominal transition kernels and reward functions are generated randomly. The  
636 number of states and the number of actions are varied. Discount factor  $\gamma = 0.9$ , reward noise radius  
637  $\alpha_{s,a}, \alpha_s = 0.1$ , transition noise kernel  $\beta_{s,a}, \beta_s = \frac{0.01}{SA}$ .

638 **Hardware** Experiments are done on the machine with the following configuration: Intel(R) Core(TM)  
639 i7-6700 CPU @3.40GHZ, size:3598MHz, capacity 4GHz, width 64 bits, memory size 64 GiB.

640 **Software and codes** All the experiments were done in Python using numpy, matplotlib. All codes  
641 and results will be made public on GitHub after the publication to preserve anonymity.

642 **Procedure and Results.** All the experiments were repeated 100 times, except for Linear Programming  
643 (LP) cases as LP methods were very time-consuming. In LP methods, experiments were repeated  
644 5 times except for the case ( $S = 500, A = 100$ ) which was done only once. As this case was  
645 prohibitively expensive. Standard deviation in all cases was less than 10%, and typically 1 – 2%.  
646 This conveniently illustrates the superiority of our methods over LP methods.

### 647 Observations

- 648 • **Scalability of our methods.** Note that our methods scale very well with large state-action  
649 space. It takes a (small) constant times the time required by non-robust MDPs. On the other  
650 hand, LP methods explode. Both observations confirm the theoretical time complexity.
- 651 • **sa-case vs s case in LP methods.** We see s-case outperforms sa-case for small state-action  
652 spaces via LP methods. This is opposite to the theoretical time complexity of s-case which  
653 expensive than sa-case. We believe this is due to the internal implementation issues. Note  
654 that computing the robust value function is the most expensive step which requires evaluation  
655 of the robust Bellman operator. In sa case, one evaluation requires solving  $SA$  LP programs  
656 with  $S$  variables each, while for s-case, it is  $S$  LP programs with  $SA$  variables each. To  
657 solve LP, `scipy.linprog` is used, we believe it does some parallelization for large LPs. Hence,  
658 we observe less cost for sa-case. However, we observe that the cost of s-case increases  
659 much faster than s-case, and eventually under-performing than sa-case.

### 660 H.1 RPG by LP

661 We compute RPG using LP in the following steps:

- 662 1. **(Robust Value Iteration)** Approximately compute the robust value function  $v_{\mathcal{U}}^{\pi}$  using the  
663 iterates  $v_{n+1} := \mathcal{T}_{\mathcal{U}}^{\pi} v_n$ . Evaluation of robust Bellman operator  $\mathcal{T}_{\mathcal{U}}^{\pi}$  is done via LPs as  
664 described below. This is the most expensive step as it requires evaluating robust Bellman  
665 operators  $O(\log(\epsilon^{-1}))$  times, and each evaluation requires many LPs.
- 666 2. **(Adversarial Values)** Compute the worst values  $(P_{\mathcal{U}}^{\pi}, R_{\mathcal{U}}^{\pi})$  using the robust value function  
667 from the following relation:

$$(P_{\mathcal{U}}^{\pi}, R_{\mathcal{U}}^{\pi}) \in \arg \min_{(P,R) \in \mathcal{U}} \mathcal{T}_{\mathcal{U}}^{\pi} v_{\mathcal{U}}^{\pi}.$$

668 This is also solved by LP.

- 669 3. **(Policy Gradient Theorem)** We now compute the RPG using policy gradient Theorem [23]  
670 w.r.t. the adversarial values computed above, as

$$\partial \rho_{\mathcal{U}}^{\pi} = \sum_{s,a} d_{P_{\mathcal{U}}^{\pi}}^{\pi}(s) Q_{P_{\mathcal{U}}^{\pi}, R_{\mathcal{U}}^{\pi}}^{\pi}(s, a) \nabla \pi_s(a).$$

671 Observe that  $d_P^\pi$  can be approximated as  $\sum_{n=0}^n (\gamma P^\pi)^n$  for large  $n$  enough, and  $Q_{P,R}^\pi$  can  
 672 be approximated by dynamic programming [22]. Notably, this step and the second step are  
 673 negligible as compared to the first step.

674 **Robust Value Iteration by LP**

675 **sa-rectangular robust MDPs.** We first consider sa-rectangular  $L_1$  constrained uncertainty set  
 676  $\mathcal{U}_p^{\text{sa}} = \mathcal{P} \times \mathcal{R}$ . Robust Bellman operator is given by

$$(\mathcal{T}_{\mathcal{U}_p^{\text{sa}}}^\pi v)(s) = \max_a \underbrace{\min_{p \in \mathcal{P}_{sa}, r \in \mathcal{R}_{sa}} \left[ r + \gamma \sum_{s'} p(s') v(s') \right]}_{\text{LP with } S \text{ variable}}.$$

677 Note that the above can be solved by  $A$  LPs as uncertainty set  $\mathcal{U}_p^{\text{sa}} = \mathcal{P} \times \mathcal{R}$  induces linear constraint  
 678 and the objective is also linear with  $S$  variables. Hence, evaluation of  $\mathcal{T}_{\mathcal{U}_p^{\text{sa}}}^\pi v$  requires solving  $SA$  LPs  
 679 with  $S$  variable each.

680  **$s$ -rectangular robust MDPs.** We now consider  $s$ -rectangular  $L_1$  constrained uncertainty set  $\mathcal{U}_p^s =$   
 681  $\mathcal{P} \times \mathcal{R}$ . Robust Bellman operator is given by

$$(\mathcal{T}_{\mathcal{U}_p^s}^\pi v)(s) = \underbrace{\min_{p \in \mathcal{P}_s, r \in \mathcal{R}_s} \sum_a \pi_s(a) \left[ r(a) + \gamma \sum_{s'} p(s'|a) v(s') \right]}_{\text{LP with } SA \text{ variable}}.$$

682 Note that the above can be solved by one LP as uncertainty set  $\mathcal{U}_p^s = \mathcal{P} \times \mathcal{R}$  induces linear constraint  
 683 and the objective is also linear with  $SA$  variables. Hence, evaluation of  $\mathcal{T}_{\mathcal{U}_p^s}^\pi v$  requires solving  $S$  LPs  
 684 with  $SA$  variable each.