
Twice regularized MDPs and the equivalence between robustness and regularization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Robust Markov decision processes (MDPs) aim to handle changing or partially
2 known system dynamics. To solve them, one typically resorts to robust optimization
3 methods. However, this significantly increases the computational complexity and
4 limits scalability in both learning and planning. On the other hand, regularized
5 MDPs show more stability in policy learning without impairing time complexity.
6 Yet, they generally do not encompass uncertainty in the model dynamics. In
7 this work, we aim to learn robust MDPs using regularization. We first show that
8 regularized MDPs are a particular instance of robust MDPs with uncertain reward.
9 We thus establish that policy iteration on reward-robust MDPs has the same time
10 complexity as on regularized MDPs. We further extend this relationship to MDPs
11 with uncertain transitions: this leads to a regularization term with an additional
12 dependence on the value function. We finally generalize regularized MDPs to twice
13 regularized MDPs (R^2 MDPs), *i.e.*, MDPs with *both* value and policy regularization.
14 The corresponding Bellman operators enable developing policy iteration schemes
15 with convergence and robustness guarantees. It also reduces planning and learning
16 in robust MDPs to regularized MDPs.

17 1 Introduction

18 Markov decision processes (MDPs) provide a practical framework for solving sequential decision
19 problems under uncertainty [30]. However, the chosen strategy can be very sensitive to sampling
20 errors or inaccurate model estimates. This can lead to complete failure in common situations where
21 the model parameters vary adversarially or are simply unknown [21]. Robust MDPs aim to mitigate
22 such sensitivity by assuming that the transition and/or reward function (P, r) varies arbitrarily inside a
23 given *uncertainty set* \mathcal{U} [16, 26]. In this setting, an optimal solution maximizes a performance measure
24 under the worst-case parameters. It can be thought of as a dynamic zero-sum game with an agent
25 choosing the best action while Nature imposes it the most adversarial model. As such, solving robust
26 MDPs involves max-min problems, which can be computationally challenging and limits scalability.

27 In recent years, several methods have been developed to alleviate the computational concerns raised
28 by robust reinforcement learning (RL). Apart from [22, 23] that consider specific types of coupled
29 uncertainty sets, all rely on a rectangularity assumption without which the problem can be NP-hard
30 [1, 40]. This rectangularity assumption is key to deriving tractable solvers of robust MDPs such
31 as robust value iteration [1, 10] or more general robust modified policy iteration (MPI) [17]. Yet,
32 reducing time complexity in robust Bellman updates remains challenging and is still researched today
33 [14, 10].

34 At the same time, the empirical success of regularization in policy search methods has motivated a
35 wide range of algorithms with diverse motivations such as improved exploration [11, 19] or stability
36 [34, 12]. Geist et al. [9] proposed a unified view from which many existing algorithms can be derived.

37 Their regularized MDP formalism enables error propagation analysis in approximate MPI [33] and
 38 leads to the same bounds as for standard MDPs. Nevertheless, as we further show in Sec. 3, policy
 39 regularization accounts for reward uncertainty only: it does not encompass uncertainty in the model
 40 dynamics. Despite a vast literature on *how* regularized policy search works and convergence rates
 41 analysis [36, 4], little attention has been given to understanding *why* it can generate strategies that
 42 are robust to external perturbations [12].

43 To our knowledge, the only works that relate robustness to regularization in RL are [5, 15, 8]. Derman
 44 & Mannor [5] employ a distributionally robust optimization approach to regularize an empirical
 45 value function. Unfortunately, computing this empirical value necessitates several policy evaluation
 46 procedures, which is quickly unpractical. Husain et al. [15] provide a dual relationship with robust
 47 MDPs under uncertain reward. Their duality result applies to general regularization methods and
 48 gives a robust interpretation of soft-actor-critic [12]. Although these two works justify the use
 49 of regularization for ensuring robustness, they do not enclose any algorithmic novelty. Similarly,
 50 Eysenbach & Levine [8] focus specifically on maximum entropy methods and relate them to either
 51 reward or transition robustness. We shall further detail on these most related studies in Sec. 6.

52 The robustness-regularization duality is well established in statistical learning theory [41, 35, 18],
 53 as opposed to RL theory. In fact, standard setups such as classification or regression may be
 54 considered as single-stage decision-making problems, *i.e.*, one-step MDPs, a particular case of RL
 55 setting. Extending this robustness-regularization duality to RL would yield cheaper learning methods
 56 with robustness properties. As such, we propose a regularization function $\Omega_{\mathcal{U}}$ that depends on the
 57 uncertainty set \mathcal{U} and is defined over both policy and value spaces (see Sec. 5), thus inducing a
 58 *twice regularized* Bellman operator. We show that this regularizer yields an equivalence of the form
 59 $v_{\pi, \mathcal{U}} = v_{\pi, \Omega_{\mathcal{U}}}$, where $v_{\pi, \mathcal{U}}$ is the robust value function for policy π and $v_{\pi, \Omega_{\mathcal{U}}}$ the regularized one.
 60 This equivalence is derived through the objective function each value optimizes. More concretely, we
 61 formulate the robust value function $v_{\pi, \mathcal{U}}$ as an optimal solution of the robust optimization problem:

$$\max_{v \in \mathbb{R}^S} \langle v, \mu_0 \rangle \text{ s. t. } v \leq \inf_{(P, r) \in \mathcal{U}} T_{(P, r)}^{\pi} v, \quad (\text{RO})$$

62 where $T_{(P, r)}^{\pi}$ is the evaluation Bellman operator. Then, we show that $v_{\pi, \mathcal{U}}$ is also an optimal solution
 63 of the convex (non-robust) optimization problem:

$$\max_{v \in \mathbb{R}^S} \langle v, \mu_0 \rangle \text{ s. t. } v \leq T_{(P_0, r_0)}^{\pi} v - \Omega_{\mathcal{U}}(\pi, v), \quad (\text{CO})$$

64 where (P_0, r_0) is the *nominal model*, establishing equivalence between the two optimization problems.
 65 Moreover, the inequality constraint of (CO) enables to derive a *twice regularized* (\mathbb{R}^2) Bellman
 66 operator defined according to $\Omega_{\mathcal{U}}$, a policy and value regularizer. For ball-constrained uncertainty
 67 sets, $\Omega_{\mathcal{U}}$ has an explicit form and under mild conditions, the corresponding \mathbb{R}^2 Bellman operators
 68 are contracting. The equivalence between the two problems (RO) and (CO) together with the
 69 contraction properties of \mathbb{R}^2 Bellman operators enable to circumvent robust optimization problems
 70 at each Bellman update. As such, it alleviates robust planning and learning algorithms by reducing
 71 them to regularized ones, which are as complex as classical methods.

72 To summarize, we make the following contributions: (i) We show that regularized MDPs are a specific
 73 instance of robust MDPs with uncertain reward. Besides formalizing a general connection between
 74 the two settings, our result enables to explicit the uncertainty sets induced by standard regularizers.
 75 (ii) We generalize this dual relationship to MDPs with uncertain transition and provide the first
 76 regularizer that recovers robust MDPs with s -rectangular balls and arbitrary norm. (iii) We introduce
 77 twice regularized MDPs that apply both policy and value regularization to retrieve robust MDPs. We
 78 establish contraction properties of the corresponding Bellman operators. This leads us to proposing a
 79 robust MPI algorithm that opens new perspectives for practical and scalable robust RL algorithms.

80 **Notations.** We designate the extended reals by $\overline{\mathbb{R}} := \{-\infty, \infty\}$. Given a finite set \mathcal{Z} , the set of
 81 real-valued functions (resp. probability distributions) over \mathcal{Z} is denoted by $\mathbb{R}^{\mathcal{Z}}$ (resp. $\Delta_{\mathcal{Z}}$), while the
 82 constant function equal to 1 over \mathcal{Z} is denoted by $\mathbb{1}_{\mathcal{Z}}$. The inner product of two functions $\mathbf{a}, \mathbf{b} \in \mathbb{R}^{\mathcal{Z}}$
 83 is defined as $\langle \mathbf{a}, \mathbf{b} \rangle := \sum_{z \in \mathcal{Z}} \mathbf{a}(z) \mathbf{b}(z)$, which induces the ℓ_2 -norm $\|\mathbf{a}\| := \sqrt{\langle \mathbf{a}, \mathbf{a} \rangle}$. The ℓ_2 -norm
 84 coincides with its dual norm, *i.e.*, $\|\mathbf{a}\| = \max_{\|\mathbf{b}\| \leq 1} \langle \mathbf{a}, \mathbf{b} \rangle =: \|\mathbf{a}\|_*$. Let a function $f : \mathbb{R}^{\mathcal{Z}} \rightarrow \overline{\mathbb{R}}$.
 85 The Legendre-Fenchel transform (or convex conjugate) of f is $f^*(\mathbf{y}) := \max_{\mathbf{a} \in \mathbb{R}^{\mathcal{Z}}} \{\langle \mathbf{a}, \mathbf{y} \rangle - f(\mathbf{a})\}$.
 86 Given a set $\mathfrak{Z} \subseteq \mathbb{R}^{\mathcal{Z}}$, the characteristic function $\delta_{\mathfrak{Z}} : \mathbb{R}^{\mathcal{Z}} \rightarrow \overline{\mathbb{R}}$ is $\delta_{\mathfrak{Z}}(\mathbf{a}) = 0$ if $\mathbf{a} \in \mathfrak{Z}$; $+\infty$ otherwise.
 87 The Legendre-Fenchel transform of $\delta_{\mathfrak{Z}}$ is the support function $\sigma_{\mathfrak{Z}}(\mathbf{y}) = \max_{\mathbf{a} \in \mathfrak{Z}} \langle \mathbf{a}, \mathbf{y} \rangle$ [2, Ex. 1.6.1].

88 2 Preliminaries

89 This section describes the background material that we use throughout our work. Firstly, we recall
 90 useful properties in convex analysis. Secondly, we address classical discounted MDPs and their
 91 linear program (LP) formulation. Thirdly, we briefly detail on regularized MDPs and the associated
 92 operators and lastly, we focus on the robust MDP setting.

93 2.1 Convex Analysis

94 Let $\Omega : \Delta_{\mathcal{Z}} \rightarrow \mathbb{R}$ be a strongly convex function. Throughout this work, the function Ω plays the
 95 role of a policy and/or value regularization function. Its Legendre-Fenchel transform Ω^* satisfies
 96 several smoothness properties, hence its alternative name "smoothed max operator" [24]. Our work
 97 will make use of the following result [13, 24].

98 **Proposition 2.1.** *Let $\Omega : \Delta_{\mathcal{Z}} \rightarrow \mathbb{R}$ be a strongly convex function. The following properties hold:*

99 (i) $\nabla \Omega^*$ is Lipschitz and satisfies $\nabla \Omega^*(\mathbf{y}) = \arg \max_{\mathbf{a} \in \Delta_{\mathcal{Z}}} \langle \mathbf{a}, \mathbf{y} \rangle - \Omega(\mathbf{a}), \forall \mathbf{y} \in \mathbb{R}^{\mathcal{Z}}$.

100 (ii) For any $c \in \mathbb{R}, \mathbf{y} \in \mathbb{R}^{\mathcal{Z}}, \Omega^*(\mathbf{y} + c\mathbb{1}_{\mathcal{Z}}) = \Omega^*(\mathbf{y}) + c$.

101 (iii) The Legendre-Fenchel transform Ω^* is non-decreasing.

102 2.2 Discounted MDPs and LP formulation

103 Consider an infinite horizon MDP $\mathcal{M} := (\mathcal{S}, \mathcal{A}, \mu_0, \gamma, P, r)$ with \mathcal{S} and \mathcal{A} finite state and action
 104 spaces respectively, $0 < \mu_0 \in \Delta_{\mathcal{S}}$ an initial state distribution and $\gamma \in (0, 1)$ a discount factor.
 105 Denoting $\mathcal{X} := \mathcal{S} \times \mathcal{A}, P \in \Delta_{\mathcal{S}}^{\mathcal{X}}$ is a transition kernel mapping each state-action pair to a probability
 106 distribution over \mathcal{S} and $r \in \mathbb{R}^{\mathcal{X}}$ a reward function. A policy $\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}$ maps any state $s \in \mathcal{S}$ to an
 107 action distribution $\pi_s \in \Delta_{\mathcal{A}}$, and we evaluate its performance through the following measure:

$$\rho(\pi) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \mu_0, \pi, P \right] = \langle v_{(P,r)}^{\pi}, \mu_0 \rangle, \quad (1)$$

108 where the expectation is conditioned on the process distribution determined by μ_0, π and P , and for
 109 all $s \in \mathcal{S}, v_{(P,r)}^{\pi}(s) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, \pi, P]$ is the *value function* at state s . Maximizing
 110 (1) defines the standard RL objective, which can be solved thanks to the Bellman operators:

$$\begin{aligned} T_{(P,r)}^{\pi} v &:= r^{\pi} + \gamma P^{\pi} v \quad \forall v \in \mathbb{R}^{\mathcal{S}}, \pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}, \\ T_{(P,r)} v &:= \max_{\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}} T_{(P,r)}^{\pi} v \quad \forall v \in \mathbb{R}^{\mathcal{S}}, \\ \mathcal{G}(v) &:= \{ \pi \in \Delta_{\mathcal{A}}^{\mathcal{S}} : T_{(P,r)}^{\pi} v = T_{(P,r)} v \} \quad \forall v \in \mathbb{R}^{\mathcal{S}}, \end{aligned}$$

111 where $r^{\pi} := [\langle \pi_s, r(s, \cdot) \rangle]_{s \in \mathcal{S}}$ and $P^{\pi} = [P^{\pi}(s' | s)]_{s', s \in \mathcal{S}}$ with $P^{\pi}(s' | s) := \langle \pi_s, P(s' | s, \cdot) \rangle$. Both
 112 $T_{(P,r)}^{\pi}$ and $T_{(P,r)}$ are γ -contractions with respect to (w.r.t.) the supremum norm, so each admits a
 113 unique fixed point $v_{(P,r)}^{\pi}$ and $v_{(P,r)}^*$, respectively. The set of greedy policies w.r.t. value v defines
 114 $\mathcal{G}_{(P,r)}(v)$, and any policy $\pi \in \mathcal{G}_{(P,r)}(v)$ is optimal [30]. For all $v \in \mathbb{R}^{\mathcal{S}}$, the associated function
 115 $q \in \mathbb{R}^{\mathcal{X}}$ is given by $q(s, a) = r(s, a) + \gamma \langle P(\cdot | s, a), v \rangle, \forall (s, a) \in \mathcal{X}$. In particular, the fixed point
 116 $v_{(P,r)}^{\pi}$ satisfies $v_{(P,r)}^{\pi} = \langle \pi_s, q_{(P,r)}^{\pi} \rangle$ where $q_{(P,r)}^{\pi}$ is its associated q -function.

117 The problem in (1) can also be formulated as an LP [30]. Given a policy $\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}$, we characterize its
 118 performance $\rho(\pi)$ by the following v -LP [30, 25]:

$$\min_{v \in \mathbb{R}^{\mathcal{S}}} \langle v, \mu_0 \rangle \text{ subject to (s.t.) } v \geq r^{\pi} + \gamma P^{\pi} v. \quad (\mathbf{P}^{\pi})$$

119 This primal objective provides a policy view of the problem. Alternatively, one may take a state
 120 visitation perspective by studying the dual objective instead:

$$\max_{\mu \in \mathbb{R}^{\mathcal{S}}} \langle r^{\pi}, \mu \rangle \text{ s. t. } \mu \geq 0 \text{ and } (\mathbf{Id}_{\mathbb{R}^{\mathcal{S}}} - \gamma P^{\pi}) \mu = \mu_0, \quad (\mathbf{D}^{\pi})$$

where P_*^{π} is the *adjoint policy transition operator* [1]: $[P_*^{\pi} \mu](s) := \sum_{\bar{s} \in \mathcal{S}} P^{\pi}(s | \bar{s}) \mu(\bar{s}), \forall \mu \in \mathbb{R}^{\mathcal{S}}$, and
 $\mathbf{Id}_{\mathcal{S}}$ is the identity function in $\mathbb{R}^{\mathcal{S}}$. Let $\mathbf{I}(s' | s, a) := \delta_{s'=s}, \forall (s, a) \in \mathcal{X}, s' \in \mathcal{S}$ the trivial transition

¹It is the adjoint operator of P^{π} in the sense that $\langle P^{\pi} v, v' \rangle = \langle v, P_*^{\pi} v' \rangle \quad \forall v, v' \in \mathbb{R}^{\mathcal{S}}$.

matrix, and define its *adjoint transition operator* as $\mathbf{I}_* \mu(s) := \sum_{(\bar{s}, \bar{a}) \in \mathcal{X}} \mathbf{I}(s|\bar{s}, \bar{a}) \mu(\bar{s}, \bar{a}), \forall s \in \mathcal{S}$. The correspondence between either the occupancy measure or the policy view lies in the one-to-one mapping $\mu \mapsto \frac{\mu(\cdot, \cdot)}{\mathbf{I}_* \mu(\cdot)} =: \pi_\mu$ and its inverse $\pi \mapsto \mu_\pi$ given by

$$\mu_\pi(s, a) := \sum_{t=0}^{\infty} \gamma^t \mathbb{P} \left(s_t = s, a_t = a \mid \mu_0, \pi, P \right), \forall (s, a) \in \mathcal{X}.$$

121 As such, one can interchangeably work with the primal LP (\mathbf{P}^π) or the dual (\mathbf{D}^π) .

122 2.3 Regularized MDPs

123 A regularized MDP is a tuple $\mathcal{M}_\Omega := (\mathcal{S}, \mathcal{A}, \mu_0, \gamma, P, r, \Omega)$ with $(\mathcal{S}, \mathcal{A}, \mu_0, \gamma, P, r)$ an infinite
124 horizon MDP as defined above, and $\Omega := (\Omega_s)_{s \in \mathcal{S}}$ a finite set of functions such that for all $s \in \mathcal{S}$,
125 $\Omega_s : \Delta_{\mathcal{A}} \rightarrow \mathbb{R}$ is strongly convex. Each function Ω_s plays the role of a policy regularizer $\Omega_s(\pi_s)$.
126 With a slight abuse of notation, we shall denote by $\Omega(\pi) := (\Omega_s(\pi_s))_{s \in \mathcal{S}}$ the family of state-
127 dependent regularizers². The regularized Bellman evaluation operator is given by

$$[T_{(P,r)}^{\pi, \Omega} v](s) := T_{(P,r)}^\pi v(s) - \Omega_s(\pi_s), \quad \forall v \in \mathbb{R}^{\mathcal{S}}, s \in \mathcal{S},$$

128 and the regularized Bellman optimality operator is $T_{(P,r)}^{*, \Omega} v := \max_{\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}} T_{(P,r)}^{\pi, \Omega} v, \forall v \in \mathbb{R}^{\mathcal{S}}$ [9].

129 The unique fixed point of $T_{(P,r)}^{\pi, \Omega}$ (respectively $T_{(P,r)}^{*, \Omega}$) is denoted by $v_{(P,r)}^{\pi, \Omega}$ (resp. $v_{(P,r)}^{*, \Omega}$) and called
130 *regularized value function* (resp. *regularized optimal value function*). Although the regularized
131 MDP formalism starts from the aforementioned Bellman operators in [9], it turns out that regularized
132 MDPs are MDPs with modified reward function. Indeed, for any policy $\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}$, the regularized
133 value function is $v_{(P,r)}^{\pi, \Omega} = (\mathbf{I}_{\mathcal{S}} - \gamma P^\pi)^{-1} (r^\pi - \Omega(\pi))$, which corresponds to a non-regularized value
134 function with expected reward $\tilde{r}^\pi := r^\pi - \Omega(\pi)$. Note that the modified reward $\tilde{r}^\pi(s)$ is no longer
135 linear in π_s because of the strong convexity of Ω_s .

136 2.4 Robust MDPs

137 In general, the MDP model is not explicitly known but rather estimated from sampled trajectories.
138 As this may result in over-sensitive outcome [21], robust MDPs aim to reduce such performance
139 variation. Formally, a robust MDP $(\mathcal{S}, \mathcal{A}, \mu_0, \gamma, \mathcal{U})$ is an MDP with uncertain model belonging to
140 $\mathcal{U} := \mathcal{P} \times \mathcal{R}$, *i.e.*, uncertain transition kernel $P \in \mathcal{P} \subseteq \Delta_{\mathcal{S}}^{\mathcal{X}}$ and reward function $r \in \mathcal{R} \subseteq \mathbb{R}^{\mathcal{X}}$
141 [16, 40]. The uncertainty set \mathcal{U} typically controls the confidence level for the model estimate, which
142 in turn determines the agent's level of robustness. It is given to the agent which seeks to maximize
143 performance under the worst-case model. Although untractable in general, this problem can be
144 solved in polynomial time for *rectangular* uncertainty sets, *i.e.*, for $\mathcal{U} = \times_{s \in \mathcal{S}} \mathcal{U}_s = \times_{s \in \mathcal{S}} (\mathcal{P}_s \times \mathcal{R}_s)$
145 [40, 22]. For any policy $\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}$ and state $s \in \mathcal{S}$, the *robust value function* at s is $v_{(P,r)}^{\pi, \mathcal{U}}(s) :=$
146 $\min_{(P,r) \in \mathcal{U}} v_{(P,r)}^\pi(s)$ and the *robust optimal value function* $v_{(P,r)}^{*, \mathcal{U}}(s) := \max_{\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}} v_{(P,r)}^{\pi, \mathcal{U}}(s)$. Both
147 objectives can be solved using the robust Bellman operators:

$$\begin{aligned} [T^{\pi, \mathcal{U}} v](s) &:= \min_{(P,r) \in \mathcal{U}} T_{(P,r)}^\pi v(s) \quad \forall v \in \mathbb{R}^{\mathcal{S}}, s \in \mathcal{S}, \pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}, \\ [T^{*, \mathcal{U}} v](s) &:= \max_{\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}} [T^{\pi, \mathcal{U}} v](s), \quad \forall v \in \mathbb{R}^{\mathcal{S}}, \end{aligned}$$

148 which are γ -contractions, so each admits a unique fixed point $v_{(P,r)}^{\pi, \mathcal{U}}$ and $v_{(P,r)}^{*, \mathcal{U}}$, respectively. For
149 all $v \in \mathbb{R}^{\mathcal{S}}$, the associated robust q -function is given by $q(s, a) = \min_{(P,r) \in \mathcal{U}} \{r(s, a) +$
150 $\gamma \langle P(\cdot|s, a), v \rangle\}, \forall (s, a) \in \mathcal{X}$, so that $v_{\pi, \mathcal{U}} = \langle \pi_s, q_{\pi, \mathcal{U}} \rangle$ where $q_{\pi, \mathcal{U}}$ is the robust q -function as-
151 sociated to $v_{\pi, \mathcal{U}}$.

152 3 Reward-robust MDPs

153 This section focuses on reward-robust MDPs, *i.e.*, robust MDPs with uncertain reward function but
154 known transition model. We show that the regularized MDP formalism introduced in [9] represents a

²In the formalism of Geist et al. [9], Ω_s is initially constant over \mathcal{S} . However, later in the paper (Sec. 5), it changes according to policy iterates. Here, we alternatively define a family Ω of state-dependent regularizers, in order to account for state-dependent uncertainty sets (see Sec. 5).

155 particular instance of reward-robust MDP, as both solve the same optimization problem. This equivalence
 156 provides a general robustness motivation for the heuristic success of policy regularization. Then,
 157 we explicit the uncertainty set underlying some standard regularization functions, thus suggesting an
 158 interpretable explanation of their empirical robustness.

159 We first show the following general result that applies to transition and reward-robust MDPs with
 160 random policies. It slightly generalizes [16], as Lemma 3.2 there focuses on uncertain-transition
 161 MDPs with deterministic policies. For completeness, we provide a proof of Prop. 3.1 in Appx. A.1

162 **Proposition 3.1.** *For any policy $\pi \in \Delta_{\mathcal{A}}^S$, the robust value function $v^{\pi, \mathcal{U}}$ is the optimal solution of
 163 the robust optimization problem:*

$$\max_{v \in \mathbb{R}^S} \langle v, \mu_0 \rangle \text{ s. t. } v \leq T_{(P,r)}^{\pi} v \text{ for all } (P, r) \in \mathcal{U}. \quad (\text{P}\mathcal{U})$$

164 In the robust optimization problem (P \mathcal{U}), the inequality constraint must hold over the whole uncertainty
 165 set \mathcal{U} . As such, a function $v \in \mathbb{R}^S$ is said to be *robust feasible* for (P \mathcal{U}) if $v \leq T_{(P,r)}^{\pi} v$ for all
 166 $(P, r) \in \mathcal{U}$ or equivalently, if $\max_{(P,r) \in \mathcal{U}} \{v(s) - T_{(P,r)}^{\pi} v(s)\} \leq 0$ for all $s \in \mathcal{S}$. Therefore, checking
 167 robust feasibility requires to solve a maximization problem. For properly structured uncertainty sets,
 168 a closed form solution can be derived, which we shall see in the sequel. As standard in the robust RL
 169 literature [32, 14, 27], the remaining of this work focuses on uncertainty sets that are centered around
 170 a known *nominal model*. Formally, given P_0 (resp. r_0) the nominal transition kernel (resp. reward
 171 function), we consider uncertainty sets of the form $(P_0 + \mathcal{P}) \times (r_0 + \mathcal{R})$. Here, the size of $\mathcal{P} \times \mathcal{R}$
 172 quantifies the level of uncertainty or alternatively, the degree of robustness.

173 3.1 Reward-robust and regularized MDPs: an equivalence

174 We now focus on reward-robust MDPs, *i.e.*, robust MDPs with $\mathcal{U} = \{P_0\} \times (r_0 + \mathcal{R})$. Thm. 3.1
 175 establishes that reward-robust MDPs are in fact regularized MDPs whose regularizer is given by a
 176 support function. Its proof can be found in Appx. A.2 This result brings two take-home messages: (i)
 177 policy regularization is equivalent to reward uncertainty; (ii) policy iteration on reward-robust MDPs
 178 has the same convergence rate as regularized MDPs, which in turn is the same as standard MDPs [9].

179 **Theorem 3.1** (Reward-robust MDP). *Assume that $\mathcal{U} = \{P_0\} \times (r_0 + \mathcal{R})$. Then, for any policy
 180 $\pi \in \Delta_{\mathcal{A}}^S$, the robust value function $v^{\pi, \mathcal{U}}$ is the optimal solution of the convex optimization problem:*

$$\max_{v \in \mathbb{R}^S} \langle v, \mu_0 \rangle \text{ s. t. } v(s) \leq T_{(P_0, r_0)}^{\pi} v(s) - \sigma_{\mathcal{R}_s}(-\pi_s) \text{ for all } s \in \mathcal{S}.$$

181 Thm. 3.1 clearly highlights a convex regularizer $\Omega_s(\pi_s) := \sigma_{\mathcal{R}_s}(-\pi_s), \forall s \in \mathcal{S}$. We thus recover a
 182 regularized MDP by setting $[T^{\pi, \Omega} v](s) = T_{(P_0, r_0)}^{\pi} v(s) - \sigma_{\mathcal{R}_s}(-\pi_s), \forall s \in \mathcal{S}$. In particular, when
 183 \mathcal{R}_s is a ball of radius α_s^r , the support function (or regularizer) can be written in closed form as
 184 $\Omega_s(\pi_s) := \alpha_s^r \|\pi_s\|$, which is strongly convex. We formalize this below (see proof in Appx. A.3).

185 **Corollary 3.1.** *Let $\pi \in \Delta_{\mathcal{A}}^S$. Further assume that for all $s \in \mathcal{S}$, the reward uncertainty set is
 186 $\mathcal{R}_s := \{r_s \in \mathbb{R}^{\mathcal{A}} : \|r_s\| \leq \alpha_s^r\}$ and $\mathcal{U} = \{P_0\} \times (r_0 + \mathcal{R})$. Then, the robust value function $v^{\pi, \mathcal{U}}$ is
 187 the optimal solution of the convex optimization problem:*

$$\max_{v \in \mathbb{R}^S} \langle v, \mu_0 \rangle \text{ s. t. } v(s) \leq T_{(P_0, r_0)}^{\pi} v(s) - \alpha_s^r \|\pi_s\| \text{ for all } s \in \mathcal{S}.$$

188 While regularization induces reward-robustness, Thm. 3.1 and Cor. 3.1 suggest that, on the other
 189 hand, specific reward-robust MDPs recover well-known policy regularization methods. We explicit
 190 the reward-uncertainty sets underlying some of these regularizers in the following.

191 3.2 Related Algorithms

192 Consider uncertainty sets of the type $\mathcal{R} := \times_{(s,a) \in \mathcal{X}} \mathcal{R}_{s,a}$. This defines an (s, a) -rectangular \mathcal{R}
 193 (a particular type of s -rectangular \mathcal{R}) whose rectangles $\mathcal{R}_{s,a}$ are independently defined for all
 194 state-action pairs. For each of the regularizers below, we derive appropriate $\mathcal{R}_{s,a}$ that recover the
 195 same regularized value function. Note that these reward uncertainty sets depend on the policy.
 196 Detailed proofs are in Appx. A.4 There, we also include a summary table that reviews properties
 197 of some RL regularizers, as well as our R^2 function which we shall introduce later in Sec. 5

198 **Negative Shannon entropy.** Let $\mathcal{R}_{s,a}^{\text{NS}}(\pi) := (-\infty, \ln(1/\pi_s(a))]$, $\forall (s, a) \in \mathcal{X}$. The associated
 199 support function enables to write:

$$\sigma_{\mathcal{R}_s^{\text{NS}}(\pi)}(-\pi_s) = \max_{r(s,a'): r(s,a') \in \mathcal{R}_{s,a'}^{\text{NS}}(\pi), a' \in \mathcal{A}} \sum_{a \in \mathcal{A}} -r(s, a) \pi_s(a) = \sum_{a \in \mathcal{A}} \pi_s(a) \ln(\pi_s(a)),$$

200 where the last equality results from maximizing $-r(s, a)$ over $(-\infty, \ln(1/\pi_s(a))]$ for each $a \in \mathcal{A}$.
 201 We thus recover the negative Shannon entropy $\Omega(\pi_s) = \sum_{a \in \mathcal{A}} \pi_s(a) \ln(\pi_s(a))$ [12].

202 **KL divergence.** Let $\mathcal{R}_{s,a}^{\text{KL}}(\pi) := \ln(1/|\mathcal{A}|) + \mathcal{R}_{s,a}^{\text{NS}}(\pi)$, $\forall (s, a) \in \mathcal{X}$. It amounts to translating the
 203 interval $\mathcal{R}_{s,a}^{\text{NS}}$ by the given constant. Then, by similarly writing the support function, we recover
 204 $\Omega(\pi_s) = \sum_{a \in \mathcal{A}} \pi_s(a) \ln(\pi_s(a)) + \ln(|\mathcal{A}|)$, which is exactly the Kullback-Leibler divergence [34].

205 **Negative Tsallis entropy.** Letting $\mathcal{R}_{s,a}^{\text{T}}(\pi) := \left[\frac{1-\pi_s(a)}{2}, +\infty \right)$, $\forall (s, a) \in \mathcal{X}$, we recover the
 206 negative Tsallis entropy $\Omega(\pi_s) = \frac{1}{2}(\|\pi_s\|^2 - 1)$ [19].

207 3.3 Policy-gradient for reward-robust MDPs

208 The equivalence between reward-robust and regularized MDPs leads us to ask whether we can employ
 209 policy-gradient methods [37] on reward-robust MDPs using regularization. The following result
 210 establishes that indeed, a policy-gradient theorem can be derived for reward-robust MDPs (see proof
 211 in Appx. A.5).

212 **Proposition 3.2.** Assume that $\mathcal{U} = \{P_0\} \times (r_0 + \mathcal{R})$ with $\mathcal{R}_s = \{r_s \in \mathbb{R}^{\mathcal{A}} : \|r_s\| \leq \alpha_s^r\}$. Then, the
 213 gradient of the reward-robust objective $J_{\mathcal{U}}(\pi) := \langle v_{\pi, \mathcal{U}}, \mu_0 \rangle$ is given by

$$\nabla J_{\mathcal{R}^2}(\pi) = \mathbb{E}_{(s,a) \sim \mu_{\pi}} \left[\nabla \ln \pi_s(a) \left(q_{\pi, \mathcal{U}}(s, a) - \alpha_s^r \frac{\pi_s(a)}{\|\pi_s\|} \right) \right].$$

214 Although Prop. 3.2 is a particular case of [9, Appx. D.3] for regularized MDPs, its application to
 215 reward-robust MDPs is novel and suggests another simplification of robust methods. Indeed, previous
 216 works that derive policy-gradient for robust MDPs involve the occupancy measure of the worst-case
 217 model, whereas our result sticks to the nominal. In practice, Prop. 3.2 enables to learn a robust policy
 218 by sampling transitions from the nominal model instead of all uncertain models. This has a twofold
 219 advantage: (i) it avoids an additional computation of the minimum as done in [29, 20, 6]; there, the
 220 authors sample next-state transitions and rewards based on all parameters from the uncertainty set,
 221 then update the policy based on the worst outcome; (ii) it releases from restricting to finite uncertainty
 222 sets. In fact, our regularizer accounts for robustness regardless of the sampling procedure, whereas
 223 the parallel simulations in [29, 20, 6] require the uncertainty set to be finite. Technical difficulties are
 224 yet to be addressed for generalizing this to transition-uncertain MDPs, because of the interdependence
 225 between the regularizer and the value function (see Secs. 4.5). We detail more on this in Appx. A.5.

226 4 General robust MDPs

227 Now that we have established policy regularization as a reward-robust problem, we would like to
 228 study the opposite question: can any robust MDP, with both uncertain reward and transition, be
 229 solved using regularization instead of robust optimization? If so, is the regularization function easy to
 230 determine? In this section, we answer positively to both questions for properly defined robust MDPs.
 231 This greatly facilitates robust reinforcement learning, as it avoids the increased complexity of robust
 232 planning algorithms while still reaching robust performance.

233 The following theorem establishes that similarly to reward-robust MDPs, any robust MDP can be
 234 formulated through regularization (see proof in Appx. B.1). Although the regularizer is also a support
 235 function in that case, it depends on both the policy and the value objective, which explains the added
 236 difficulty of dealing with robust MDPs.

237 **Theorem 4.1** (Transition-robust MDP). Assume that $\mathcal{U} = (P_0 + \mathcal{P}) \times (r_0 + \mathcal{R})$. Then, for all policy
 238 $\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}$, the robust value function $v^{\pi, \mathcal{U}}$ is the optimal solution of the convex optimization problem:

$$\max_{v \in \mathbb{R}^{\mathcal{S}}} \langle v, \mu_0 \rangle \text{ s. t. } v(s) \leq T_{(P_0, r_0)}^{\pi} v(s) - \sigma_{\mathcal{R}_s}(-\pi_s) - \sigma_{\mathcal{P}_s}(-\gamma v \cdot \pi_s) \text{ for all } s \in \mathcal{S}, \quad (2)$$

239 where $[v \cdot \pi_s](s', a) := v(s') \pi_s(a)$, $\forall (s', a) \in \mathcal{X}$.

240 The upper-bound in the inequality constraint (2) is of the same spirit as the regularized Bellman
 241 operator: the first term is a standard, non-regularized Bellman operator on the nominal model (P_0, r_0)
 242 to which we subtract a policy-dependent function playing the role of regularization. That support
 243 function reminds that of [5, Thm. 3.1], also coming from conjugacy. This is the only similarity
 244 between both regularizers: in [5], the Legendre-Fenchel transform is applied on a different type of
 245 function, which results in a regularization term that has no closed form but can only be bounded from
 246 above. Moreover, the setup considered in [5] is different from our robust MDP setting since it studies
 247 a Wasserstein distributionally robust MDP. Therefore, it involves a general convex optimization
 248 problem, whereas we focus on the robust formulation of an LP.

249 When the uncertainty set is a ball, the support function simplifies further, as we show in the following
 250 Cor. 4.1. Yet, the dependence of the regularization term in the value function prevents us from
 251 readily applying the tool-set of regularized MDPs. We shall later study the properties of this new
 252 regularization function in Sec. 5.

253 **Corollary 4.1.** *Assume that $\mathcal{U} = (P_0 + \mathcal{P}) \times (r_0 + \mathcal{R})$ with $\mathcal{P}_s := \{P_s \in \mathbb{R}^{\mathcal{X}} : \|P_s\| \leq \alpha_s^P\}$ and
 254 $\mathcal{R}_s := \{r_s \in \mathbb{R}^{\mathcal{A}} : \|r_s\| \leq \alpha_s^r\}$ for all $s \in \mathcal{S}$. Then, the robust value function $v^{\pi, \mathcal{U}}$ is the optimal
 255 solution of the convex optimization problem:*

$$\max_{v \in \mathbb{R}^{\mathcal{S}}} \langle v, \mu_0 \rangle \text{ s. t. } v(s) \leq T_{(P_0, r_0)}^{\pi} v(s) - \alpha_s^r \|\pi_s\| - \alpha_s^P \|\pi_s\| \cdot \gamma \|v\| \text{ for all } s \in \mathcal{S}. \quad (3)$$

256 In fact, Cor. 4.1 can be reformulated for any arbitrary norm by replacing each $\|\cdot\|$ by its dual norm
 257 $\|\cdot\|_*$ in Eq. (3) (see proof in Appx. B.2): we state Cor. 4.1 with the ℓ_2 -norm for notation convenience
 258 only, the dual norm of ℓ_2 being norm- ℓ_2 itself. Thus, our regularization function recovers a robust
 259 value function independently of the chosen norm, which extends previous results in [14, 10]. Indeed,
 260 Ho et al. [14] lighten complexity of robust planning for ℓ_1 -norm only, while Grand-Clément & Kroer
 261 [10] focus on KL and ℓ_2 ball-constrained uncertainty sets. Both works rely on the specific structure
 262 induced by the divergence they consider to derive more efficient robust Bellman updates. Differently,
 263 we circumvent these updates using a generic regularization method that is problem-independent while
 264 still applying to s -rectangular uncertainty sets as [14, 10] do.

265 5 \mathbb{R}^2 MDPs

266 In Sec. 4 we showed that for transition-robust MDPs, the optimization constraint involves a regu-
 267 larization term that depends on the value function itself. This adds a difficulty to the reward-robust
 268 case where the regularization only depends on the policy. On the other hand, we provided an explicit
 269 regularizer for either reward or transition robust MDPs that are ball-constrained. In this section, we
 270 introduce \mathbb{R}^2 MDPs, a generalization of regularized MDPs that combines policy and value regulariza-
 271 tion. The core idea is to further regularize the Bellman operators with a value-dependent term that
 272 recovers the support functions we derived from the robust optimization problems of Secs. 3-4.

273 **Definition 5.1.** *For all $v \in \mathbb{R}^{\mathcal{S}}$, let $\Omega_{v, \mathbb{R}^2} : \Delta_{\mathcal{A}} \rightarrow \mathbb{R}$ defined as $\Omega_{v, \mathbb{R}^2}(\pi_s) := \|\pi_s\|(\alpha_s^r + \alpha_s^P \gamma \|v\|)$.
 274 The \mathbb{R}^2 Bellman evaluation operator is defined as*

$$[T^{\pi, \mathbb{R}^2} v](s) := T_{(P_0, r_0)}^{\pi} v(s) - \Omega_{v, \mathbb{R}^2}(\pi_s), \quad \forall s \in \mathcal{S}.$$

275 *The \mathbb{R}^2 Bellman optimal operator is defined as*

$$[T^{*, \mathbb{R}^2} v](s) := \max_{\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}} [T^{\pi, \mathbb{R}^2} v](s) = \Omega_{v, \mathbb{R}^2}^*(q_s), \quad \forall s \in \mathcal{S}.$$

For any function $v \in \mathbb{R}^{\mathcal{S}}$, the associated unique greedy policy is defined as

$$\pi_s = \arg \max_{\pi_s \in \Delta_{\mathcal{A}}} T^{\pi, \mathbb{R}^2} v(s) = \nabla \Omega_{v, \mathbb{R}^2}^*(q_s), \quad \forall s \in \mathcal{S},$$

276 *that is, in vector form, $\pi = \nabla \Omega_{v, \mathbb{R}^2}^*(q) =: \mathcal{G}_{\Omega_{v, \mathbb{R}^2}^*}(v) \iff T^{\pi, \mathbb{R}^2} v = T^{*, \mathbb{R}^2} v$.*

277 The \mathbb{R}^2 Bellman evaluation operator is not linear because of the functional norm appearing in the
 278 regularization function. Yet, under mild conditions, it is contracting and we can apply Banach's fixed
 279 point theorem to define the \mathbb{R}^2 value function.

Assumption 5.1 (Bounded radius). *For all $s \in \mathcal{S}$, there exists $\epsilon_s > 0$ such that*

$$\alpha_s^P \leq \min \left(\frac{1 - \gamma - \epsilon_s}{\gamma}; \min_{\substack{\mathbf{u}_A \in \mathbb{R}_+^A, \|\mathbf{u}_A\|=1 \\ \mathbf{v}_S \in \mathbb{R}_+^S, \|\mathbf{v}_S\|=1}} \mathbf{u}_A^\top P_0(\cdot|s, \cdot) \mathbf{v}_S \right).$$

280

281

282 **Algorithm 1:** \mathbb{R}^2 MPI

283 **Result:** π_{k+1}, v_{k+1}

284 Initialize $v_k \in \mathbb{R}^S$;

285 **while not converged do**

286 $\pi_{k+1} \leftarrow \mathcal{G}_{\Omega_{R^2}}(v_k)$;

287 $v_{k+1} \leftarrow (T^{\pi_{k+1}, R^2})^m v_k$;

288 **end**

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

Asm. 5.1 requires an upper-bound on ball radius for transition uncertainty sets. The first term in the minimum is needed for establishing the contraction property of \mathbb{R}^2 Bellman operators, while the second one is used for monotonicity. We remark that the former depends on the original discount factor γ : radius α_s^P must be smaller as γ tends to 1, but as γ decreases to 0, it can grow arbitrarily without altering contraction. Indeed, larger γ implies longer time horizon and higher stochasticity. At the same time, Asm. 5.1 requires tighter level of uncertainty for solving robust MDPs via regularization. Otherwise, both value and policy regularization seem unable to handle the mixed effects of parameter and stochastic uncertainties. Although we recognize

a generalized Rayleigh quotient-type problem in the second minimum [28], its interpretation in our context remains unclear. Asm. 5.1 enables the \mathbb{R}^2 Bellman operators to admit a unique fixed point, among other nice properties. We formalize this in the following proposition (see proof in Appx. C.1).

Proposition 5.1. *Suppose that Asm. 5.1 holds. Then, we have the following properties:*

(i) *Monotonicity:* For all $v_1, v_2 \in \mathbb{R}^S$ such that $v_1 \leq v_2$, we have $T^{\pi, R^2} v_1 \leq T^{\pi, R^2} v_2$ and $T^{*, R^2} v_1 \leq T^{*, R^2} v_2$.

(ii) *Sub-distributivity:* For all $c \in \mathbb{R}$, we have $T^{\pi, R^2}(v_1 + c\mathbb{1}_S) \leq T^{\pi, R^2} v_1 + \gamma c\mathbb{1}_S$ and $T^{*, R^2}(v_1 + c\mathbb{1}_S) \leq T^{*, R^2} v_1 + \gamma c\mathbb{1}_S, \forall c \in \mathbb{R}$.

(iii) *Contraction:* Let $\epsilon_* := \min_{s \in \mathcal{S}} \epsilon_s > 0$. Then, for all $v_1, v_2 \in \mathbb{R}^S$, $\|T^{\pi, R^2} v_1 - T^{\pi, R^2} v_2\|_\infty \leq (1 - \epsilon_*) \|v_1 - v_2\|_\infty$ and $\|T^{*, R^2} v_1 - T^{*, R^2} v_2\|_\infty \leq (1 - \epsilon_*) \|v_1 - v_2\|_\infty$.

The fixed point property of \mathbb{R}^2 Bellman operators leads us to define \mathbb{R}^2 value functions as follows.

Definition 5.2 (\mathbb{R}^2 value functions). (i) *The \mathbb{R}^2 value function v^{π, R^2} is defined as the unique fixed point of the \mathbb{R}^2 Bellman evaluation operator: $v^{\pi, R^2} = T^{\pi, R^2} v^{\pi, R^2}$. The associated q -function is $q^{\pi, R^2}(s, a) = r(s, a) + \gamma \langle P_0(\cdot|s, a), v^{\pi, R^2} \rangle$. (ii) *The \mathbb{R}^2 optimal value function v^{*, R^2} is defined as the unique fixed point of the \mathbb{R}^2 Bellman optimal operator: $v^{*, R^2} = T^{*, R^2} v^{*, R^2}$. The associated q -function is $q^{*, R^2}(s, a) = r(s, a) + \gamma \langle P_0(\cdot|s, a), v^{*, R^2} \rangle$.**

Monotonicity of \mathbb{R}^2 Bellman operators plays a key role for reaching an optimal \mathbb{R}^2 policy, as we show in the following. A proof can be found in Appx. C.2.

Theorem 5.1 (\mathbb{R}^2 optimal policy). *The policy $\pi^{*, R^2} = \mathcal{G}_{\Omega_{R^2}}(v^{*, R^2})$ is the unique optimal \mathbb{R}^2 policy, i.e., for all $\pi \in \Delta_{\mathcal{A}}, v^{\pi, R^2} = v^{*, R^2} \geq v^{\pi, R^2}$.*

All of the results above ensure convergence of MPI in \mathbb{R}^2 MDPs. We call that method \mathbb{R}^2 MPI and provide its pseudo-code in Alg. 1. The convergence proof follows the same lines as in [30]. Regarding time complexity of the greedy step, computing the optimal \mathbb{R}^2 operator amounts to projecting onto the simplex, which can be efficiently performed in linear time [7].

6 Related Work

Connections between regularization and robustness have already been established in standard statistical learning settings such as support vector machines [41], logistic regression [35] or maximum likelihood estimation [18]. As stated in Sec. 1, these single-stage decision-making problems are a particular instance of RL problems. In that regard, the generalization of robustness-regularization duality to sequential decision-making is seldom studied in the RL literature.

Two works that interpret policy regularization from a robustness perspective are [15] and [8]. In [15], regularization is applied on the dual objective instead of the primal, which has two shortcomings:

324 (i) It prevents from deriving regularized Bellman operators and dynamic programming methods; (ii)
 325 The feasible set is that of occupancy measures, so the connection with standard policy regularization
 326 remains unclear. Furthermore, their work focus on reward robustness only. Differently, Eysenbach &
 327 Levine [8] address both reward and transition uncertainty by showing that regularized policies with
 328 maximum entropy solve a particular type of robust MDPs. Yet, their analysis treats each uncertainty
 329 separately, which questions the robustness of the resulting policy when the whole model (P, r) is
 330 adversarial. Moreover, their dual relation between entropy regularization and transition-robust MDPs
 331 is weak besides applying to specific uncertainty sets. Both of these works treat robustness as an
 332 side-effect of regularization more than an objective on its own, whereas we aim to do the opposite,
 333 namely, use regularization to solve robust RL problems.

334 Derman & Mannor [5] similarly view regularization as a tool for achieving robust policies. Through
 335 their distributionally robust MDP framework, they show upper and lower bounds between transition-
 336 robustness and regularization. There again, duality is weak while reward uncertainty is not considered.
 337 Moreover, since the exact regularization term has no explicit form, it is usable through its upper bound
 338 only. Finally, regularization is applied on the mean of several value functions $v_{(\hat{P}_i, r)}^\pi$, where each \hat{P}_i
 339 is a transition model estimated from an episode run. Computing this quantity would require as many
 340 policy evaluations as available model estimates available, yielding at least a linear complexity blowup.

341 Previous studies analyze robust planning algorithms and provide convergence guarantees: in [1, 26]
 342 a value iteration algorithm is proposed, the works [16, 40] introduce robust policy iteration and
 343 Kaufman & Schaefer [17] generalize both schemes by studying the conditions under which robust
 344 MPI converges. All of them guarantee a robust solution in polynomial time which is often insufficient,
 345 as the complexity of each Bellman update grows at least cubically in the number of states [14].

346 In order to reduce time complexity of robust planning algorithms, Ho et al. [14] propose two algo-
 347 rithms that compute robust Bellman updates in $\mathcal{O}(|\mathcal{S}||\mathcal{A}| \log(|\mathcal{S}||\mathcal{A}|))$ operations for ℓ_1 -constrained
 348 uncertainty sets. Advantageously, our regularization approach reduces each robust Bellman update
 349 to its standard, non-robust version of $\mathcal{O}(|\mathcal{S}||\mathcal{A}|)$ computations. Moreover, although (s, a) and s -
 350 rectangular uncertainty sets are considered in [14], they concern transition-uncertainty only, while
 351 our model includes both uncertainties in the general s -rectangular case. Also, their contribution
 352 relies on LP formulations that necessitate restricting to ℓ_1 -norm. Differently, our formulation applies
 353 to any norm, which may come from the fact that our main Theorems (Thms. 3.1, 4.1) use Fenchel-
 354 Rockafellar duality [31, 3], a generalization of LP duality (see Appx. A.2, B.1). More recently, in
 355 [10], Grand-Clément & Kroer propose a first-order method to accelerate robust value iteration under
 356 s -rectangular uncertainty sets that are ellipsoidal or KL balls. To our knowledge, our study is the first
 357 that reduces time complexity for uncertainty sets as general as s -rectangular sets with arbitrary norm.

358 7 Conclusion and future work

359 In this work, we established strong duality between robust and regularized MDPs with a properly
 360 chosen regularization function, thus making regularized MDPs a particular case of robust MDPs
 361 with uncertain reward. This enabled us to derive a policy-gradient theorem for reward-robust MDPs.
 362 When extending this robustness-regularization duality to robust MDPs with uncertain reward and
 363 transition, we found that the regularizer depends on the value function besides the policy. We thus
 364 introduced R^2 MDPs, a generalization of regularized MDPs with both policy and value regularization.
 365 Their related R^2 Bellman operators enable to derive a converging MPI algorithm that achieves the
 366 optimal robust value function.

367 This study provides the theoretical foundations of scalable robust RL. Although our theory focused
 368 on planning, the R^2 Bellman operators and their contracting properties open the path to learning
 369 algorithms that can (i) use existing RL algorithms and robustify them by simply changing the
 370 regularizer; (ii) scale to deep learning settings. Apart from its practical effect, we believe our work
 371 opens the path to more theoretical contributions in robust RL. For example, it would be interesting
 372 to extend R^2 MPI to the approximate case [33], as non-linearity of the R^2 evaluation operator
 373 complicates the analysis. Other possible directions are sample complexity analysis for R^2 MDPs
 374 and its comparison with that of robust MDPs, as done recently in [42]; approximate dynamic
 375 programming for R^2 MDPs in light of their robust analog [38, 27]. Another line of research is to
 376 extend policy-gradient to twice regularized MDPs, as it would avoid parallel learning of adversarial
 377 models [6, 39] and can be very useful for continuous control.

378 **References**

- 379 [1] Bagnell, J. A., Ng, A. Y., and Schneider, J. G. Solving uncertain Markov decision processes.
380 2001.
- 381 [2] Bertsekas, D. P. *Convex optimization theory*. Athena Scientific Belmont, 2009.
- 382 [3] Borwein, J. and Lewis, A. S. *Convex analysis and nonlinear optimization: theory and examples*.
383 Springer Science & Business Media, 2010.
- 384 [4] Cen, S., Cheng, C., Chen, Y., Wei, Y., and Chi, Y. Fast global convergence of natural policy
385 gradient methods with entropy regularization. *arXiv preprint arXiv:2007.06558*, 2020.
- 386 [5] Derman, E. and Mannor, S. Distributional robustness and regularization in reinforcement
387 learning. *ICML Workshop*, 2020.
- 388 [6] Derman, E., Mankowitz, D., Mann, T., and Mannor, S. Soft-robust actor-critic policy-gradient.
389 *AUAI press for Association for Uncertainty in Artificial Intelligence*, pp. 208–218, 2018.
- 390 [7] Duchi, J., Shalev-Shwartz, S., Singer, Y., and Chandra, T. Efficient projections onto the l_1
391 ball for learning in high dimensions. In *Proceedings of the 25th international conference on*
392 *Machine learning*, pp. 272–279, 2008.
- 393 [8] Eysenbach, B. and Levine, S. Maximum entropy RL (provably) solves some robust RL problems.
394 *arXiv preprint arXiv:2103.06257*, 2021.
- 395 [9] Geist, M., Scherrer, B., and Pietquin, O. A theory of regularized Markov decision processes. In
396 *International Conference on Machine Learning*, pp. 2160–2169. PMLR, 2019.
- 397 [10] Grand-Clément, J. and Kroer, C. Scalable first-order methods for robust MDPs. *arXiv preprint*
398 *arXiv:2005.05434*, 2020.
- 399 [11] Haarnoja, T., Tang, H., Abbeel, P., and Levine, S. Reinforcement learning with deep energy-
400 based policies. In *International Conference on Machine Learning*, pp. 1352–1361. PMLR,
401 2017.
- 402 [12] Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum
403 entropy deep reinforcement learning with a stochastic actor. In *International Conference on*
404 *Machine Learning*, pp. 1861–1870. PMLR, 2018.
- 405 [13] Hiriart-Urruty, J.-B. and Lemaréchal, C. *Fundamentals of convex analysis*. Springer Science &
406 Business Media, 2004.
- 407 [14] Ho, C. P., Petrik, M., and Wiesemann, W. Fast bellman updates for robust MDPs. In *Interna-*
408 *tional Conference on Machine Learning*, pp. 1979–1988. PMLR, 2018.
- 409 [15] Husain, H., Ciosek, K., and Tomioka, R. Regularized policies are reward robust. In *International*
410 *Conference on Artificial Intelligence and Statistics*, pp. 64–72. PMLR, 2021.
- 411 [16] Iyengar, G. N. Robust dynamic programming. *Mathematics of Operations Research*, 30(2):
412 257–280, 2005.
- 413 [17] Kaufman, D. L. and Schaefer, A. J. Robust modified policy iteration. *INFORMS Journal on*
414 *Computing*, 25(3):396–410, 2013.
- 415 [18] Kuhn, D., Esfahani, P. M., Nguyen, V. A., and Shafieezadeh-Abadeh, S. Wasserstein distri-
416 butionally robust optimization: Theory and applications in machine learning. In *Operations*
417 *Research & Management Science in the Age of Analytics*, pp. 130–166. INFORMS, 2019.
- 418 [19] Lee, K., Choi, S., and Oh, S. Sparse Markov decision processes with causal sparse tsallis
419 entropy regularization for reinforcement learning. *IEEE Robotics and Automation Letters*, 3(3):
420 1466–1473, 2018.
- 421 [20] Mankowitz, D., Mann, T., Bacon, P.-L., Precup, D., and Mannor, S. Learning robust options. In
422 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.

- 423 [21] Mannor, S., Simester, D., Sun, P., and Tsitsiklis, J. N. Bias and variance approximation in value
424 function estimates. *Management Science*, 53(2):308–322, 2007.
- 425 [22] Mannor, S., Mebel, O., and Xu, H. Lightning does not strike twice: Robust MDPs with coupled
426 uncertainty. *ICML*, 2012.
- 427 [23] Mannor, S., Mebel, O., and Xu, H. Robust MDPs with k-rectangular uncertainty. *Mathematics
428 of Operations Research*, 41(4):1484–1509, 2016.
- 429 [24] Mensch, A. and Blondel, M. Differentiable dynamic programming for structured prediction and
430 attention. In *International Conference on Machine Learning*, pp. 3462–3471. PMLR, 2018.
- 431 [25] Nachum, O. and Dai, B. Reinforcement learning via fenchel-rockafellar duality. *arXiv preprint
432 arXiv:2001.01866*, 2020.
- 433 [26] Nilim, A. and El Ghaoui, L. Robust control of Markov decision processes with uncertain
434 transition matrices. *Operations Research*, 53(5):780–798, 2005.
- 435 [27] Panaganti, K. and Kalathil, D. Robust reinforcement learning using least squares policy iteration.
436 *arXiv preprint arXiv:2006.11608*, 2020.
- 437 [28] Parlett, B. N. The rayleigh quotient iteration and some generalizations for nonnormal matrices.
438 *Mathematics of Computation*, 28(127):679–693, 1974.
- 439 [29] Pinto, L., Davidson, J., Sukthankar, R., and Gupta, A. Robust adversarial reinforcement learning.
440 In *International Conference on Machine Learning*, pp. 2817–2826. PMLR, 2017.
- 441 [30] Puterman, M. L. *Markov decision processes: discrete stochastic dynamic programming*. John
442 Wiley & Sons, 2014.
- 443 [31] Rockafellar, R. T. *Convex analysis*, volume 36. Princeton university press, 1970.
- 444 [32] Roy, A., Xu, H., and Pokutta, S. Reinforcement learning under model mismatch. *NIPS*, 2017.
- 445 [33] Scherrer, B., Ghavamzadeh, M., Gabillon, V., Lesner, B., and Geist, M. Approximate modified
446 policy iteration and its application to the game of Tetris. *J. Mach. Learn. Res.*, 16:1629–1676,
447 2015.
- 448 [34] Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P. Trust region policy optimization.
449 In *International conference on machine learning*, pp. 1889–1897. PMLR, 2015.
- 450 [35] Shafieezadeh-Abadeh, S., Esfahani, P. M., and Kuhn, D. Distributionally robust logistic
451 regression. *NIPS*, 2015.
- 452 [36] Shani, L., Efroni, Y., and Mannor, S. Adaptive trust region policy optimization: Global
453 convergence and faster rates for regularized MDPs. In *Proceedings of the AAAI Conference on
454 Artificial Intelligence*, volume 34, pp. 5668–5675, 2020.
- 455 [37] Sutton, R. S., McAllester, D. A., Singh, S. P., Mansour, Y., et al. Policy gradient methods
456 for reinforcement learning with function approximation. In *NIPS*, volume 99, pp. 1057–1063.
457 Citeseer, 1999.
- 458 [38] Tamar, A., Mannor, S., and Xu, H. Scaling up robust MDPs using function approximation. In
459 *International Conference on Machine Learning*, pp. 181–189. PMLR, 2014.
- 460 [39] Tessler, C., Efroni, Y., and Mannor, S. Action robust reinforcement learning and applications in
461 continuous control. In *International Conference on Machine Learning*, pp. 6215–6224. PMLR,
462 2019.
- 463 [40] Wiesemann, W., Kuhn, D., and Rustem, B. Robust Markov decision processes. *Mathematics of
464 Operations Research*, 38(1):153–183, 2013.
- 465 [41] Xu, H., Caramanis, C., and Mannor, S. Robustness and regularization of support vector
466 machines. *Journal of machine learning research*, 10(7), 2009.
- 467 [42] Yang, W. and Zhang, Z. Non-asymptotic performances of Robust Markov Decision Processes.
468 *arXiv preprint arXiv:2105.03863*, 2021.

469 **Checklist**

- 470 1. For all authors...
- 471 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
472 contributions and scope? [Yes] The main contributions are consistently listed in the
473 abstract and introduction. These are: Showing that regularized MDPs are a particular
474 instance of robust MDPs with uncertain reward (Sec. 3); Extending this relationship to
475 MDPs with uncertain transitions (Sec. 4); Generalizing regularized MDPs to R^2 MDPs
476 (Sec. 5).
- 477 (b) Did you describe the limitations of your work? [Yes] These are described when
478 discussing Prop. 3.2 and Asm. 5.1 near the corresponding statements.
- 479 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 480 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
481 them? [Yes]
- 482 2. If you are including theoretical results...
- 483 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Assumptions
484 are introduced by "Assume that..." or "If..." in each theoretical result.
- 485 (b) Did you include complete proofs of all theoretical results? [Yes] All proofs can be
486 found in the Appendix. At each theoretical statement in the text body, we precisely
487 refer to the corresponding proof.
- 488 3. If you ran experiments...
- 489 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
490 mental results (either in the supplemental material or as a URL)? [N/A]
- 491 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
492 were chosen)? [N/A]
- 493 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
494 ments multiple times)? [N/A]
- 495 (d) Did you include the total amount of compute and the type of resources used (e.g., type
496 of GPUs, internal cluster, or cloud provider)? [N/A]
- 497 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 498 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 499 (b) Did you mention the license of the assets? [N/A]
- 500 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 501
- 502 (d) Did you discuss whether and how consent was obtained from people whose data you're
503 using/curating? [N/A]
- 504 (e) Did you discuss whether the data you are using/curating contains personally identifiable
505 information or offensive content? [N/A]
- 506 5. If you used crowdsourcing or conducted research with human subjects...
- 507 (a) Did you include the full text of instructions given to participants and screenshots, if
508 applicable? [N/A]
- 509 (b) Did you describe any potential participant risks, with links to Institutional Review
510 Board (IRB) approvals, if applicable? [N/A]
- 511 (c) Did you include the estimated hourly wage paid to participants and the total amount
512 spent on participant compensation? [N/A]