
Near-Optimal Distributionally Robust Reinforcement Learning with General L_p Norms

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 To address the challenges of sim-to-real gap and sample efficiency in reinforcement
2 learning (RL), this work studies distributionally robust Markov decision processes
3 (RMDPs) — optimize the worst-case performance when the deployed environment
4 is within an uncertainty set around some nominal MDP. Despite recent efforts,
5 the sample complexity of RMDPs has remained largely undetermined. While the
6 statistical implications of distributional robustness in RL have been explored in
7 some specific cases, the generalizability of the existing findings remains unclear,
8 especially in comparison to standard RL. Assuming access to a generative model
9 that samples from the nominal MDP, we examine the sample complexity of
10 RMDPs using a class of generalized L_p norms as the 'distance' function for the
11 uncertainty set, under two commonly adopted *sa*-rectangular and *s*-rectangular
12 conditions. Our results imply that RMDPs can be more sample-efficient to solve
13 than standard MDPs using generalized L_p norms in both *sa*- and *s*-rectangular
14 cases, potentially inspiring more empirical research. We provide a near-optimal
15 upper bound and a matching minimax lower bound for the *sa*-rectangular scenarios.
16 For *s*-rectangular cases, we improve the state-of-the-art upper bound and also
17 derive a lower bound using L_∞ norm that verifies the tightness.

18 1 Introduction

19 Reinforcement learning (RL) [Sutton, 1988] is a popular paradigm in machine learning, particularly
20 noted for its success in practical applications. The RL framework, usually modeled within the context
21 of a Markov decision process (MDP), focuses on learning effective decision-making strategies based
22 on interactions with an environment. However, the work of Mannor et al. [2004], among others,
23 has highlighted a vulnerability in RL strategies, revealing the sensitivity to estimation errors in the
24 reward and transition probabilities. A specific example of this is when, because of a sim-to-real gap,
25 policies learned in idealized environments catastrophically fail when deployed in settings with slight
26 changes or adversarial perturbations [Klopp et al., 2017, Mahmood et al., 2018].

27 To address this issue, robust MDPs (RMDPs), proposed by Iyengar [2005] and Nilim and El Ghaoui
28 [2005], have attracted considerable attention. RMDPs are formulated as max-min problems,
29 seeking policies that are resilient to model estimation errors within a specified uncertainty set.
30 Despite the robustness benefits, solving RMDPs is NP-hard for general uncertainty sets [Nilim and
31 El Ghaoui, 2005]. To overcome this challenge, the assumption of rectangularity is often adopted,
32 with uncertainty sets structured as products of independent subsets for each state or state-action pair,
33 denoted as *s*-rectangular or *sa*-rectangular assumptions (see Definitions 4 and 5). These assumptions
34 facilitate the use of methods such as robust value iteration and robust policy iteration, preserving
35 many structural properties of MDPs [Ho et al., 2021]. The *s*-rectangular sets, though less restrictive,
36 pose greater challenges, while the *sa*-rectangular sets allow for deterministic optimal policies akin

Result type	Reference	Distance	<i>sa</i> -rectangularity		<i>s</i> -rectangularity	
			$0 < \sigma \lesssim 1 - \gamma$	$1 - \gamma \lesssim \sigma < \sigma_{\max}$	$0 < \bar{\sigma} \lesssim 1 - \gamma$	$1 - \gamma \lesssim \bar{\sigma} < \bar{\sigma}_{\max}$
Upper bound	Yang et al. [2022a]	TV	$\frac{S^2 A(2+\sigma)^2}{\sigma^2(1-\gamma)^2 \varepsilon^2}$	$\frac{S^2 A(2+\sigma)^2}{\sigma^2(1-\gamma)^2 \varepsilon^2}$	$\frac{S^2 A^2(2+\bar{\sigma})^2}{\bar{\sigma}^2(1-\gamma)^2 \varepsilon^2}$	$\frac{S^2 A^2(2+\bar{\sigma})^2}{\bar{\sigma}^2(1-\gamma)^2 \varepsilon^2}$
	Panaganti and Kalathil [2022]	TV	$\frac{S^2 A}{(1-\gamma)^3 \varepsilon^2}$	$\frac{S^2 A}{(1-\gamma)^3 \varepsilon^2}$	×	×
	Shi et al. [2023]	TV	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{\sigma(1-\gamma)^2 \varepsilon^2}$	×	×
	Clavier et al. [2023]	L_p	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$
	This paper	L_p	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{\sigma(1-\gamma)^2 \varepsilon^2}$	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{(1-\gamma)^2 \bar{\sigma} \min_s \ \pi_s\ _s \varepsilon^2}$
	This paper	General L_p [1]	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{\sigma(1-\gamma)^2 \varepsilon^2}$	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{(1-\gamma)^2 \bar{\sigma} C_D \min_s \ \pi_s\ _s \varepsilon^2}$
Lower bound	Yang et al. [2022a]	TV	$\frac{SA(1-\gamma)}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA(1-\gamma)}{\sigma^2 \varepsilon^2}$	×	×
	Shi et al. [2023]	TV	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{\sigma(1-\gamma)^2 \varepsilon^2}$	×	×
	This paper	L_p	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{\sigma(1-\gamma)^2 \varepsilon^2}$	×	×
	This paper	L_∞	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{\sigma(1-\gamma)^2 \varepsilon^2}$	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{\bar{\sigma}(1-\gamma)^2 \varepsilon^2}$

Table 1: Comparisons with prior results (up to log terms) regarding finding an ε -optimal policy for the distributionally RMDP, where σ is the radius of the uncertainty set and σ_{\max} defined in Theorem 1.

37 to non-robust MDPs [Wiesemann et al., 2013]. Note that, while uncertainty in the reward can be
38 easily handled, dealing with uncertainty in the transition kernel is much more difficult [Kumar et al.,
39 2022, Derman et al., 2021].

40 The question of sample efficiency is central in RL problems ranging from practice to theory. Although
41 minimax rates are achieved in [Azar et al., 2013b, Li et al., 2023c] in the context of classical MDPs,
42 this goal remains open, in general, in the context of RMDPs. Specifically, there exists prior work
43 studying the sample complexity of distributionally robust RL for a few specific divergences such
44 as total variation (*TV*), χ^2 , *KL*, and Wasserstein (see a further discussion in Appendix 6) [Yang
45 et al., 2022b, Zhou et al., 2021, Panaganti and Kalathil, 2022], while such results remain unclear
46 for more general classes of L_p norms defined in 1. To this point, to the best of our knowledge, the
47 results of sample complexity that achieve minimax optimality for the full range of uncertainty level
48 are limited to only one case — *TV* distance [Shi et al., 2023].

49 In this work, we focus on understanding the sample complexity of RMDPs with a general smooth
50 L_p that will be defined in Def. 1. This generalization is appealing for both practice and theory. In
51 practice, numerous applications are based on optimizations or learning approaches that involve
52 general norms beyond those that have already been studied. Additionally, optimizing norm weighted
53 ambiguity sets for Robust MDPs has been proposed in the context of RMDPs in Russel et al. [2019],
54 which justifies our formulation. Theoretically, prior work has characterized the sample complexity of
55 RMDPs for some specific norms have suggested intriguing insights about the statistical implications
56 of distributional robustness in RL. It is interesting to further understand the statistical cost of robust
57 RL in more general scenarios. One area of focus is the contrast between the sample efficiency of
58 solving distributionally robust RL and solving standard RL. In particular, for the specific case of
59 *TV* distance, Shi et al. [2023] shows that the sample complexity for solving robust RL is at least
60 the same as and sometimes (when the uncertainty level is relatively large) could be smaller than
61 that of standard RL. This motivates the following open question:

62 *Is distributionally robust RL more sample efficient than standard RL for norms defined in Def. (1) ?*

63 A second question is about the comparisons between the sample complexity of solving *s*-rectangular
64 RMDPs and that of solving *sa*-rectangular RMDPs. Note that *s*-rectangular RMDPs have more
65 complicated optimization formulations with additional variables (uncertainty levels for each action) to
66 optimize. This leads to a richer class of optimal policy candidates—stochastic policies in *s*-rectangular
67 cases, in contrast to the class of deterministic policies for *sa*-rectangular cases. In addition, existing
68 sample complexity upper bounds for solving *s*-rectangular RMDPs are larger than that for solving
69 *sa*-rectangularity [Yang et al., 2022b] for the investigated cases. This motivates the curious question:

70 *Does solving s-rectangular RMDPs require more samples than solving sa-rectangular RMDPs with*
71 *general smooth L_p norms defined in Def. 1?*

72 **Main contributions.** In this paper, we address each of the two questions discussed above. In
73 particular, we provide the first sample complexity analysis for RMDPs with general L_p norms defined
74 in 1 under both the s - and sa -rectangularity conditions. For convenience, we present a detailed
75 comparison between the existing state-of-the-art and our results in Table 1 for quick reference and
76 discuss the contributions and their implications below.

77 • Considering the first question, we illustrate our results in both sa - and s -rectangular case in
78 Figure 2. In the case of sa -rectangularity, we derive a sample complexity upper bound for RMDPs
79 using general smooth L_p norms (cf. Theorem 1) in the order of $\tilde{O}\left(\frac{SA}{(1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} \varepsilon^2}\right)$. with
80 $C_g > 0$ a positive constant related to the geometry of the norm defined in 1. For classical L_P norms,
81 $C_g \geq 1$ so we can directly relax this constant to 1 to obtain the result in table 1. In addition, we
82 provide a matching minimax lower bound (cf. Theorem 2) that confirms the near-optimality of
83 the upper bound for almost full range of the uncertainty level. Our results match the near-optimal
84 sample complexity derived in Shi et al. [2023] for the specific case using TV distance, while holding
85 for broader cases using general L_p norms. The results rely on a new dual optimization form for
86 sa -rectangular RMDPs and reveal the relationship between the sample complexity and this new dual
87 form — the infinite span seminorm (controlled in Lemma 5), which may be of independent interest.

88 In the case of s -rectangularity, we provide a sample complexity upper bound for solving RMDPs
89 with general smooth L_p norms in the order of $\tilde{O}\left(\frac{SA}{(1-\gamma)^2 \max\{1-\gamma, C_g \min_s \|\pi_s\|_* \sigma\} \varepsilon^2}\right)$. This result
90 improves the prior art $\tilde{O}\left(\frac{SA}{(1-\gamma)^4 \varepsilon^2}\right)$ in Clavier et al. [2023] for classical L_p when $\tilde{\sigma} \lesssim 1 - \gamma$ — by
91 at least a factor of $O\left(\frac{1}{1-\gamma}\right)$. Furthermore, we present a lower bound for a representative case with
92 L_∞ norm, which corroborates the tightness of the upper bound. To the best of our knowledge, this
93 is the first lower bound for solving RMDPs with s -rectangularity.

94 • Considering the second question, as illustrated in Figure 2, our results highlight that robust RL is at
95 least the same as and sometimes can be more sample-efficient to solve than standard RL for general
96 smooth L_p norms in 1. This insight is of significant practical importance and serves to provide
97 crucial motivation for the use and study of distributionally robustness in RL. Notably, robust RL
98 does not only reduce the vulnerability of RL policy to estimation errors and sim-to-real gaps, but
99 also leads to better data efficiency. In terms of comparing the statistical implications of sa - and
100 s -rectangularity, our results show that solving s -rectangular RMDPs is not harder than solving
101 sa -rectangular RMDPs in terms of sample requirement (See Theorem 3 and Figure 2, Right).

102 • We highlight the technical contributions as below. For the upper bounds, regarding optimization
103 contribution, we derive new dual optimization problem forms for both sa - and s -rectangular
104 cases (Lemma 3 and 4), which is the foundation of the covering number argument in finite-sample
105 analysis. From a statistical point of view, a new concentration lemma (See Lemma 8 for dual
106 forms and two new lemmas to obtain sample complexity lower than classical RL, controlling the
107 infinite span semi norm of the value function, both for sa - and s -rectangular case are derived
108 (See Lemmas 5 and 6). For the lower bound, the technical contributions are mainly in s -rectangular
109 cases, which involves entire new challenges compared to sa -rectangularity case: the optimal policies
110 can be stochastic and hard to be characterized as a closed form, compared to the deterministic one
111 in sa -rectangular cases. Therefore, we construct new hard instances for s -rectangular cases that
112 is distinct from those used in sa -rectangular cases or standard RL.

113 2 Problem Formulation: Robust Markov Decision Processes

114 In this section, we formulate distributionally robust Markov decision processes (RMDPs) in the
115 discounted infinite-horizon setting, introduce the sampling mechanism, and describe our goal.

116 **Standard Markov decision processes (MDPs).** A discounted infinite-horizon MDP is represented
117 by $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \gamma, P, r)$, where $\mathcal{S} = \{1, \dots, S\}$ and $\mathcal{A} = \{1, \dots, A\}$ are the finite state and action
118 spaces, respectively, $\gamma \in [0, 1]$ is the discounted factor, $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ denotes the probability
119 transition kernel, and $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is the immediate reward function, which is assumed to
120 be deterministic. Moreover, we assume that the reward function is bounded in $(0, 1)$ without loss of
121 generality of the results due to the variance reward invariance. Finally we denote 1_A or 1_S the unitary
122 vector of respectively dimension A or S . Moreover, e_s is the standard unitary vector supported

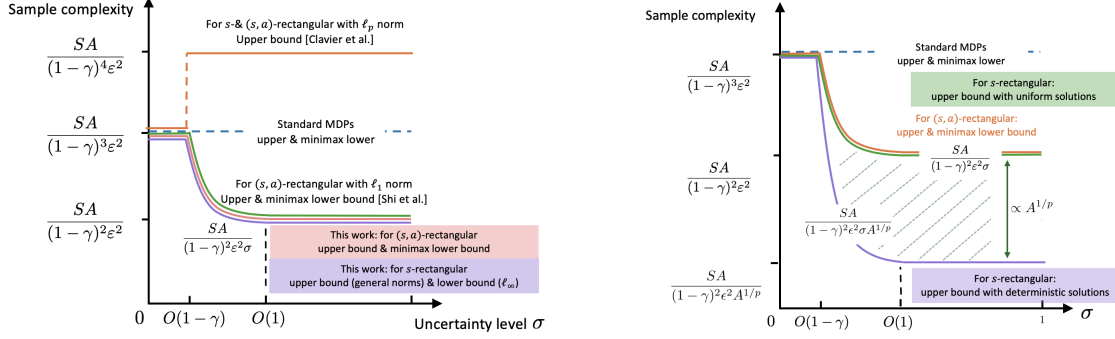


Figure 1: **Left:** Sample complexity results for RMDPs with sa - and s -rectangularity with L_p with comparisons to prior arts [Shi et al., 2023] (for L_1 norm, or called total variation distance) and [Clavier et al., 2023] ; **Right:** The data and instance-dependent sample complexity upper bound of solving s -rectangular dependency RMDPs with L_P norms.

123 on s . The policy we are looking for is denoted by $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$, which specifies the probability
 124 of action selection over the action space in any state. Note that if the policy is deterministic in the
 125 sa -rectangular case, we overload the notation and refer to $\pi(s)$ as the action selected by the policy
 126 π in state s . Finally, to characterize the cumulative reward, the value function $V^{\pi, P}$ for any policy
 127 π under the transition kernel P is defined by $\forall s \in \mathcal{S}$

$$V^{\pi, P}(s) := \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right]. \quad (1)$$

128 The expectation is taken over the randomness of the trajectory $\{s_t, a_t\}_{t=0}^{\infty}$ generated by executing
 129 the policy π under the transition kernel P , such that $a_t \sim \pi(\cdot \mid s_t)$ and $s_{t+1} \sim P(\cdot \mid s_t, a_t)$ for all
 130 $t \geq 0$. In the same way, the Q function $Q^{\pi, P}$ associated with any policy π under the transition kernel
 131 P is defined using expectation taken over the randomness of the trajectory under policy π as

$$Q^{\pi, P}(s, a) := \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0, a_0 = s, a \right], \quad (2)$$

132 **Distributionally robust MDPs.** We consider distributionally robust MDPs (RMDPs) in the
 133 discounted infinite-horizon setting, denoted by $\mathcal{M}_{\text{rob}} = \{\mathcal{S}, \mathcal{A}, \gamma, \mathcal{U}_{\|\cdot\|}^{\sigma}(P^0), r\}$, where $\mathcal{S}, \mathcal{A}, \gamma, r$
 134 are the same sets and parameters as in standard MDPs. The main difference compared to standard
 135 MDPs is that instead of assuming a fixed transition kernel P , it allows the transition kernel to be
 136 arbitrarily chosen from a prescribed uncertainty set $\mathcal{U}_{\|\cdot\|}^{\sigma}(P^0)$ centered around a *nominal* kernel
 137 $P^0 : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$, where the uncertainty set is specified using some called smooth norm denoted
 138 $\|\cdot\|$ defined in of radius $\sigma > 0$ defined in 1.

139 **Definition 1** (General smooth L_p norms and dual norms). A norm $\|\cdot\|$ is said to be a general smooth
 140 L_p norm if

- 141 • for all $x \in \mathbb{R}^n$, $\|x\| = \|x\|_{p, w} = (\sum_{k=1}^n w_k (|x_k|)^p)^{1/p}$, where $w \in \mathbb{R}_+^n$, is an arbitrary
 142 positive vector;
- 143 • it is twice continuously differentiable Rudin et al. [1964] with the supremum of the Hessian
 144 Matrix over the simple $C_S = \sup_{x \in \Delta_S} \|\nabla^2 \|x\|\|_2$, where $\|\cdot\|_2$ here is the spectral norm

145 Finally, we denote the dual norm of $\|\cdot\|$ as $\|\cdot\|_*$ s.t. $\|y\|_* = \max_x x^T y : \|x\| \leq 1$. Moreover, for any
 146 metric $\|\cdot\|$, we define C_g as $C_g = 1 / \min_s \|e_s\|$ where $e_s \in \mathbb{R}^S$ is the standard basis of supported in s .

147 Note the quantity C_S exists as the Hessian is continuous for C^2 functional and the simplex is a com-
 148 pact set, so by Extreme Value Theorem Rudin et al. [1964], C_S is finite. Moreover, to give an example,

149 considering $L_p, p \geq 2$, norms, C_s is bounded by $S^{1/q}$. (See (151)) This definition is general and
 150 includes $L_p, p \geq 2$, all rescaled and weighted norms. Moreover, we could extend our result to a larger
 151 set than the one of the norms defined in Def. 1, this is why a complete discussion about the set of norms
 152 can be found in Appendix 7. However, it does not include divergences such as KL and χ^2 . Not that
 153 the case of TV which is not C^2 smooth is treated independently with different arguments in the proof
 154 but has the same sample complexity. In particular, given the nominal transition kernel P^0 and some un-
 155 certainty level σ , the uncertainty set—with arbitrary smooth norm metric $\|\cdot\| : \mathbb{R}^S \times \mathbb{R}^+ \rightarrow \mathbb{R}^+$ in sa rect-
 156 angular case or from $\mathbb{R}^{S \times \mathcal{A}}$ in the s -rectangular case, is specified as $\mathcal{U}_{\|\cdot\|}^\sigma(P^0) := \otimes_{s,a} \mathcal{U}_{\|\cdot\|}^{sa,\sigma}(P_{s,a}^0)$

$$\mathcal{U}_{\|\cdot\|}^{sa,\sigma}(P_{s,a}^0) := \{P_{s,a} \in \Delta(\mathcal{S}) : \|P_{s,a} - P_{s,a}^0\| \leq \sigma\}, \quad (3)$$

$$P_{s,a} := P(\cdot | s, a) \in \mathbb{R}^{1 \times S}, P_{s,a}^0 := P^0(\cdot | s, a) \in \mathbb{R}^{1 \times S}. \quad (4)$$

157 where we denote a vector of the transition kernel P or P^0 at state-action pair (s, a) . In other
 158 words, the uncertainty is imposed in a decoupled manner for each state-action pair, obeying the
 159 so-called sa -rectangularity [Zhou et al., 2021, Wiesemann et al., 2013]. More generally, we
 160 define s -rectangular MDPs as $\mathcal{U}_{\|\cdot\|}^\sigma(P) = \otimes_s \mathcal{U}_{\|\cdot\|}^{s,\tilde{\sigma}}(P_s)$, for the general smooth L_p norm $\|\cdot\|$. The
 161 uncertainty is imposed in a decoupled manner for each state pair, and a fixed budget given a state
 162 for all action is defined. To get a similar meaning for the radius of the ball between sa -rectangular
 163 and s -rectangular assumptions, we need to rescale the radius depending on the norm like in Yang
 164 et al. [2022b]. The s -uncertainty set is then defined using the rescaled radius $\tilde{\sigma}$ as

$$\mathcal{U}_{\|\cdot\|}^{s,\tilde{\sigma}}(P_s) := \left\{ P'_s \in \Delta(\mathcal{S})^{\mathcal{A}} : \|P'_s - P_s\| \leq \tilde{\sigma} = \sigma \|1_A\| \right\}, \quad (5)$$

$$P_s := P(\cdot, \cdot | s) \in \mathbb{R}^{1 \times SA}, \quad P_s^0 := P^0(\cdot, \cdot | s) \in \mathbb{R}^{1 \times SA}. \quad (6)$$

165 where $1_A \in \mathbb{R}^A$ denotes the unitary vector. For the specific case of respectively L_1, L_p and L_∞ norm,
 166 $\tilde{\sigma}$ is equal to $|\sigma \mathcal{A}|, |\sigma \mathcal{A}|^{1/p}$ and σ . Note that this scaling allows for a fair comparison between sa -
 167 and s -rectangular MDPs. In RMDPs, we are interested in the worst-case performance of a policy
 168 π over all the possible transition kernels in the uncertainty set. This is measured by the *robust value*
 169 *function* $V^{\pi,\sigma}$ and the *robust Q-function* $Q^{\pi,\sigma}$ in \mathcal{M}_{rob} , defined respectively as $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$

$$V^{\pi,\sigma}(s) := \inf_{P \in \mathcal{U}_{\|\cdot\|}^{sa,\sigma}(P^0)} V^{\pi,P}(s), \quad Q^{\pi,\sigma}(s, a) := \inf_{P \in \mathcal{U}_{\|\cdot\|}^{sa,\sigma}(P^0)} Q^{\pi,P}(s, a). \quad (7)$$

170 Similarly for s -rectangularity, the value function is denoted $V_s^{\pi,\sigma}(s) := \inf_{P \in \mathcal{U}_{\|\cdot\|}^{s,\tilde{\sigma}}(P^0)} V^{\pi,P}(s)$.

171 **Optimal robust policy and robust Bellman operator.** As a generalization of properties of standard
 172 MDPs in the sa -rectangular robust case, it is well-known that there exists at least one deterministic
 173 policy that maximizes the robust value function (resp. robust Q-function) simultaneously for all states
 174 (resp. state-action pairs) [Iyengar, 2005, Nilim and El Ghaoui, 2005] but not in the s -rectangular case.
 175 Therefore, we denote the *optimal robust value function* (resp. *optimal robust Q-function*) as $V^{*,\sigma}$
 176 (resp. $Q^{*,\sigma}$), and the optimal robust policy as π^* , which satisfy $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$

$$V^{*,\sigma}(s) := V^{\pi^*,\sigma}(s) = \max_{\pi} V^{\pi,\sigma}(s), \quad Q^{*,\sigma}(s, a) := Q^{\pi^*,\sigma}(s, a) = \max_{\pi} Q^{\pi,\sigma}(s, a). \quad (8a)$$

177 A key concept in RMDPs is a generalization of Bellman's optimality principle, encapsulated in the
 178 following *robust Bellman consistency equation* (resp. *robust Bellman optimality equation*):

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \quad Q^{\pi,\sigma}(s, a) = r(s, a) + \gamma \inf_{P \in \mathcal{U}_{\|\cdot\|}^{sa,\sigma}(P_{s,a}^0)} \mathcal{P}V^{\pi,\sigma}, \quad (9a)$$

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \quad Q^{*,\sigma}(s, a) = r(s, a) + \gamma \inf_{P \in \mathcal{U}_{\|\cdot\|}^{sa,\sigma}(P_{s,a}^0)} \mathcal{P}V^{*,\sigma}. \quad (9b)$$

179 for the sa -rectangular case and same equation replacing $P_{s,a}^0$ by P_s^0 and σ by $\tilde{\sigma}$. The robust Bellman
 180 operator [Iyengar, 2005, Nilim and El Ghaoui, 2005] is denoted by $\mathcal{T}^\sigma(\cdot) : \mathbb{R}^{SA} \rightarrow \mathbb{R}^{SA}$

$$\mathcal{T}^\sigma(Q^\pi)(s, a) := r(s, a) + \gamma \inf_{P \in \mathcal{U}_{\|\cdot\|}^{sa,\sigma}(P_{s,a}^0)} \mathcal{P}V, \quad \text{with } V(s) := \max_{\pi} Q^\pi(s, a). \quad (10)$$

181 for sa -rectangular MDPs. Given that $Q^{*,\sigma}$ is the unique-fixed point of \mathcal{T}^σ one can recover the
 182 optimal robust value function and Q-function using a procedure termed *distributionally robust*
 183 *value iteration (DRV I)*. Generalizing the standard value iteration, *DRV I* starts from some given
 184 initialization and recursively applies the robust Bellman operator until convergence. As has been
 185 shown previously, this procedure converges rapidly due to the γ -contraction property of \mathcal{T}^σ with
 186 respect to the L_∞ norm [Iyengar, 2005, Nilim and El Ghaoui, 2005].

187 3 Distributionally Robust Value Iteration

188 **Generative model-based sampling.** Following Zhou et al. [2021], Panaganti and Kalathil [2022],
 189 we assume access to a generative model or a simulator [Kearns and Singh, 1999], which allows us
 190 to collect N independent samples for each state-action pair generated based on the *nominal* kernel
 191 $P^0: \forall (s, a) \in \mathcal{S} \times \mathcal{A}, s_{i,s,a} \stackrel{i.i.d.}{\sim} P^0(\cdot | s, a), \quad i = 1, 2, \dots, N$. The total sample size is, therefore,
 192 NSA . We consider a model-based approach tailored to RMDPs, which first constructs an empirical
 193 nominal transition kernel based on the collected samples and then applies distributionally robust
 194 value iteration (DRVI) to compute an optimal robust policy. As we decouple the statistical estimation
 195 error and the optimization error, we exhibit an algorithm that can achieve arbitrary small error ϵ_{opt}
 196 in the empirical MDP defined as an empirical nominal transition kernel $\hat{P}^0 \in \mathbb{R}^{\mathcal{S}A \times \mathcal{S}}$ that can be
 197 constructed on the basis of the empirical frequency of state transitions, i.e. $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$

$$\hat{P}^0(s' | s, a) := \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{s_{i,s,a} = s'\}, \quad (11)$$

198 which leads to an empirical RMDP $\widehat{\mathcal{M}}_{\text{rob}} = \{\mathcal{S}, \mathcal{A}, \gamma, \mathcal{U}_{\|\cdot\|}^\sigma(\hat{P}^0), r\}$. Analogously, we can define
 199 the corresponding robust value function (resp. robust Q-function) of policy π in $\widehat{\mathcal{M}}_{\text{rob}}$ as $\widehat{V}^{\pi, \sigma}$
 200 (resp. $\widehat{Q}^{\pi, \sigma}$) (cf. (8)). In addition, we denote the corresponding *optimal robust policy* as $\widehat{\pi}^*$ and the
 201 *optimal robust value function* (resp. *optimal robust Q-function*) as $\widehat{V}^{*, \sigma}$ (resp. $\widehat{Q}^{*, \sigma}$) (cf. (9)), which
 202 satisfies the robust Bellman optimality equation $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$:

$$\widehat{Q}^{*, \sigma}(s, a) = r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}_{\|\cdot\|}^{\text{sa}, \sigma}(\hat{P}_{s,a}^0)} \mathcal{P} \widehat{V}^{*, \sigma}. \quad (12)$$

203 Equipped with \hat{P}^0 , we can define the empirical robust Bellman operator $\widehat{\mathcal{T}}^\sigma$ as $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$

$$\widehat{\mathcal{T}}^\sigma(Q^\pi)(s, a) := r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}_{\|\cdot\|}^{\text{sa}, \sigma}(\hat{P}_{s,a}^0)} \mathcal{P} V, \quad (13)$$

204 with $V(s) := \max_{\pi} Q^\pi(s, a)$. The aim of this work is given the collected samples, to learn
 205 the robust optimal policy for the RMDP w.r.t. some prescribed uncertainty set $\mathcal{U}^\sigma(P^0)$ around the
 206 nominal kernel using as few samples as possible. Specifically, given some target accuracy level $\epsilon > 0$,
 207 the goal is to seek an ϵ -optimal robust policy $\widehat{\pi}$ obeying

$$\forall s \in \mathcal{S}: \quad V^{*, \sigma}(s) - V^{\widehat{\pi}, \sigma}(s) \leq \epsilon. \quad (14)$$

$$\widehat{V}^{\widehat{\pi}^*, \sigma} - \widehat{V}^{\widehat{\pi}, \sigma} \leq \epsilon_{opt}. \quad (15)$$

208 This formulation allows plugging any solver of RMDPs in this bound, for instance, the distributionally
 209 robust value iteration (DRVI) algorithm detailed in Appendix 12.

210 4 Theoretical guarantees

211 In this section, we present our main results characterizing the sample complexity of solving RMDPs
 212 with *sa*- and *s*-rectangularity. Additionally, we discuss the implications of our results for the com-
 213 parisons between standard and robust RL, and for comparisons between *sa*- versus *s*-rectangularity.

214 4.1 *sa*-rectangular uncertainty set with general smooth norms

215 To begin, we consider the RMDPs with *sa*-rectangularity with general norms. We first provide the
 216 following sample complexity upper bound for certain oracle planning algorithms, whose proof is
 217 postponed to Appendix 9.2. Technically, we derive two new dual forms for RMDPs problems using
 218 arbitrary norms in Lemmas 3 and 4 for respectively *sa*- and *s*-rectangular RMDPs. In these dual
 219 forms, a central quantity denoted $\text{sp}(\cdot)_*$, representing the dispersion of the value function, appears
 220 and is the dual span semi-norm associated with the considered general L_p norm $\|\cdot\|$ defined in 1
 221 in the initial primal problem. The main challenge in this analysis is to derive a tight upper bound
 222 on this quantity in Lemmas (5) and (6), leading to the following sample complexity.

223 **Theorem 1** (Upper bound for sa -rectangularity). Consider the uncertainty set $U_{\|\cdot\|}^{sa,\sigma}(\cdot)$ associated
 224 with arbitrary smooth norm $\|\cdot\|$ defined in 1. We denote $\sigma_{\max} := \max_{p_1, p_2 \in \Delta(S)} \|p_1 - p_2\|$ as
 225 the accessible maximal uncertainty level. Consider any $\delta \in (0, 1)$, discount factor $\gamma \in [\frac{1}{4}, 1)$, and
 226 uncertainty level $\sigma \in (0, \sigma_{\max}]$. Let $\hat{\pi}$ be the output policy of some oracle planning algorithm with
 227 optimization error ε_{opt} introduced in (15). With introduced in 1, one has with probability at least $1 - \delta$,

$$\forall s \in \mathcal{S} : \quad V^{*,\sigma}(s) - V^{\hat{\pi},\sigma}(s) \leq \varepsilon + \frac{8\varepsilon_{\text{opt}}}{1 - \gamma} \quad (16)$$

228 for any $\varepsilon \in (0, \sqrt{1/\max\{1 - \gamma, \sigma C_g\}}]$, as long as the total number of samples obeys

$$NSA \gtrsim \frac{c_1 SA}{(1 - \gamma)^2 \max\{1 - \gamma, C_g \sigma\} \varepsilon^2} + \frac{c_2 SAC_S \|1_S\|_*}{(1 - \gamma)^2 \varepsilon} \quad (17)$$

229 with c_1, c_2, c_3 a universal positive constant. For a sufficiently small level of accuracy
 230 $\varepsilon \leq (\max\{1 - \gamma, C_g \sigma\}) / (C_S \|1_S\|)$, the sample complexity is

$$NSA \gtrsim \frac{c_3 SA}{(1 - \gamma)^2 \max\{1 - \gamma, C_g \sigma\} \varepsilon^2} \quad (18)$$

231 Note that this result is also true for TV without the geometric smooth term depending on C_S . Consid-
 232 ering L_p norms, $C_g \geq 1$ and $C_S \leq S^{1/q}$. In Theorem 1, we introduce the following minimax-optimal
 233 lower bound to verify the tightness of the above upper bound; a proof is provided in Appendix 10.

234 **Theorem 2** (Lower bound for sa -rectangularity). Consider the uncertainty set $U_{\|\cdot\|}^{sa,\sigma}(\cdot)$ associated
 235 with arbitrary L_p norm $\|\cdot\|$ defined in 1. We denote $\sigma_{\max} := \max_{p_1, q_1 \in \Delta(S)} \|p_1 - p_2\|$ as
 236 the accessible maximal uncertainty level. Consider any tuple $(S, A, \gamma, \sigma, \varepsilon)$, where $\gamma \in [\frac{1}{2}, 1)$,
 237 $\sigma \in (0, \sigma_{\max}(1 - c_0)]$ with $0 < c_0 \leq \frac{1}{8}$ being any small enough positive constant, and $\varepsilon \in$
 238 $(0, \frac{c_0}{256(1-\gamma)}]$. We can construct two infinite-horizon RMDPs $\mathcal{M}_0, \mathcal{M}_1$ such that giving a dataset
 239 with N independent samples for each state-action pair over the nominal transition kernel (for either
 240 \mathcal{M}_0 or \mathcal{M}_1 respectively), one has

$$\inf_{\hat{\pi}} \max_{\mathcal{M} \in \{\mathcal{M}_0, \mathcal{M}_1\}} \left\{ \mathbb{P}_{\mathcal{M}} \left(\max_{s \in \mathcal{S}} [V^{*,\sigma}(s) - V^{\hat{\pi},\sigma}(s)] > \varepsilon \right) \right\} \geq \frac{1}{8},$$

241 where the infimum is taken over all estimators $\hat{\pi}$, \mathbb{P}_0 (resp. \mathbb{P}_1) are the probability when the RMDP is
 242 \mathcal{M}_0 (resp. \mathcal{M}_1), as long as, for c_7 is a universal positive constant,

$$NSA \leq \frac{c_7 SA}{(1 - \gamma)^2 \max\{1 - \gamma, C_g \sigma\} \varepsilon^2}. \quad (19)$$

243 • **Near minimax-optimal sample complexity with general L_p norms.** Recall that Theorem 1
 244 shows that the sample complexity upper bound of oracle algorithms for RMDPs is in the order of
 245 $\tilde{O}\left(\frac{SA}{(1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} \varepsilon^2}\right)$. Combined with the lower bound in Theorem 2, we observe that the
 246 above sample complexity is near minimax-optimal, in almost the full range of uncertainty.

247 • **Solving RMDPs with general L_p norms can be easier than solving standard RL.** Recall that
 248 the sample complexity of solving standard RL with a generative model [Agarwal et al., 2020, Li
 249 et al., 2024, Azar et al., 2013a] is: $\tilde{O}\left(\frac{SA}{(1-\gamma)^3 \varepsilon^2}\right)$. Comparing this with the sample complexity in
 250 (18), it highlights that solving robust MDPs (cf. (18)) using any norm as the divergence function for
 251 the uncertainty set is not harder than (and is sometimes easier than) solving standard RL (cf. (4.1)).
 252 Specifically, when the uncertainty level is small $\sigma \lesssim 1 - \gamma$, the sample complexity of solving
 253 robust MDPs matches that of standard MDPs. While when the uncertainty level is relatively larger
 254 $1 - \gamma \lesssim \sigma \leq \sigma_{\max}$, the sample complexity of solving robust MDPs is smaller than that of standard
 255 MDPs by a factor or $\frac{\sigma}{1-\gamma}$, which goes to $\frac{1}{1-\gamma}$ when $\sigma = O(1)$.

256 • **Comparisons with prior arts.** In Figure 2, we illustrate the comparisons with two state-of-the-
 257 art [Clavier et al., 2023, Shi et al., 2023] which use some divergence functions belonging to the class
 258 of general norms considered in this work. In particular, Shi et al. [2023] achieved the state-of-the-art

259 minimax-optimal sample complexity $\tilde{O}\left(\frac{SA}{(1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2}\right)$ for specific L_1 norm (or called total
 260 variation distance). In this work, we attain near minimax-optimal sample complexity for any general
 261 norm (including L_1) which matches the one in Shi et al. [2023] when narrowing down to L_1 norm.
 262 Note that in TV case, $C_g = 1$. This reveals that the finding of robust MDPs can be easier than
 263 standard MDPs [Shi et al., 2023] in terms of sample requirement does not only hold for L_1 norm,
 264 but for any general norm. In addition, compared to Clavier et al. [2023] which focuses on L_p norms
 265 for any $1 \leq p \leq \infty$: when $1 - \gamma \lesssim \sigma \leq \sigma_{\max}$, we improve the sample complexity $\tilde{O}\left(\frac{SA}{(1-\gamma)^4 \varepsilon^2}\right)$ to
 266 $\tilde{O}\left(\frac{SA}{(1-\gamma)^2 \sigma \varepsilon^2}\right)$ by at least a factor of $\frac{1}{1-\gamma}$; otherwise, we match the results in Clavier et al. [2023].

267 **Burn-in Condition, C_g factor and TV case :** In Th. 1 and 3 we need a sufficiently small level
 268 of accuracy $\epsilon \leq (\max\{1 - \gamma, C_g \sigma\}) / (C_S \|1_S\|)$, to obtain the sample complexity. This type of
 269 condition is usual in MDPS analysis Shi et al. [2022] and is equivalent to burn in term. Moreover,
 270 the quantity C_S exists (see 1) and for example, considering L_p norms, C_S is bounded by $S^{1/q}$. (See
 271 (151)) and the product $C_S \|1_S\|$ is upper bounded by S for L_2 norm. Moreover, note that our theorem
 272 for the smooth norm is also true for TV which is not C^2 and has the same complexity as (Shi et al.
 273 [2023]). In this case, the burn-in condition is not needed. (See Lemma 9.3.3). Finally, the factor
 274 $C_g = 1 / \min_s \|e_s\|$ is norm dependent and depends on how big the vector e_{s_0} is in the considered
 275 norm. Note for classical L_p this quantity is bigger than 1, which reduces the sample complexity.

276 4.2 s -rectangular uncertainty set with general norms

277 To continue, we move on to the case when the uncertainty set is constructed under s -rectangularity
 278 smooth norm. The following theorem presents the sample complexity upper bound for learning an
 279 ϵ -optimal policy for RMDPs with s -rectangularity. A proof is shown in Appendix 9.2.

280 **Theorem 3** (Upper bound for s -rectangularity). *Consider the uncertainty set $\mathcal{U}_{\|\cdot\|}^{s, \tilde{\sigma}}(\cdot)$ with
 281 s -rectangularity. Consider any discount factor $\gamma \in [\frac{1}{4}, 1)$, the rescaled uncertainty level $\tilde{\sigma} = \sigma \|1_A\|$,
 282 and denote $\tilde{\sigma}_{\max} := \|1_A\| \max_{p_1, p_2 \in \Delta(S)} \|p_1 - p_2\|$ and $\delta \in (0, 1)$. Let $\hat{\pi}$ be the output policy of
 283 an arbitrary optimization algorithm with error ε_{opt} , with probability at least $1 - \delta$, one has for any
 284 $\varepsilon \in (0, \sqrt{1 / \max\{1 - \gamma, C_g \min_s \|\pi_s\|_* \sigma\}}]$, $\forall s \in \mathcal{S}$: $V^{*, \tilde{\sigma}}(s) - V^{\hat{\pi}, \tilde{\sigma}}(s) \leq \varepsilon + \frac{\delta \varepsilon_{\text{opt}}}{1-\gamma}$ as long
 285 as the total number of samples obeys*

$$NSA \gtrsim \frac{c_4 SA}{(1-\gamma)^2 \varepsilon^2} \min \left\{ \frac{1}{\max\{1-\gamma, C_g \sigma\}}, \frac{1}{\sigma C_g \min_{s \in \mathcal{S}} \{ \|\pi_s^*\|_* \|1_A\|, \|\hat{\pi}_s\|_* \|1_A\| \}} \right\} + \frac{c_5 SAC_S \|1_S\|_*}{(1-\gamma)^2 \varepsilon} \quad (20)$$

286 *For a sufficiently small accuracy, $\epsilon \leq (\max\{1 - \gamma, C_g \tilde{\sigma}\}) / (C_S \|1_S\|)$ the sample complexity is*

$$NSA \gtrsim \frac{c_6 SA}{(1-\gamma)^2 \varepsilon^2} \min \left\{ \frac{1}{\max\{1-\gamma, C_g \sigma\}}, \frac{1}{\sigma C_g \min_{s \in \mathcal{S}} \{ \|\pi_s^*\|_* \|1_A\|, \|\hat{\pi}_s\|_* \|1_A\| \}} \right\} \quad (21)$$

287 where $\hat{\pi}_s \in \Delta_A$ denote the policy of the empirical RMPDs at state s , $\pi_s^* \in \Delta_A$ the optimal policy
 288 given s of the true RMPDs, $\|\cdot\|_*$ the dual norm and c_4, c_5, c_6 are universal constant. Note that this
 289 result is also true for TV without the term depending on smoothness C_S . In addition, we provide the
 290 lower bounds for a representative divergence function — L_∞ norm in the following. Note that for
 291 classical L_p , $C_S = S^{1/q}$ and C_g can be lower bounded by 1. A proof is provided in Appendix 11.

292 **Theorem 4** (Lower bound for s -rectangularity). *Consider the uncertainty set $\mathcal{U}_{L_\infty}^{s, \tilde{\sigma}}(\cdot)$ associated with
 293 the L_∞ norm. Consider any tuple $(S, A, \gamma, \sigma, \varepsilon)$ and $0 < c_0 \leq \frac{1}{8}$ being any small enough positive
 294 constant, where $\gamma \in [\frac{1}{2}, 1)$, and $\varepsilon \in (0, \frac{c_0}{256(1-\gamma)}]$. Correspondingly, we denote the accessible
 295 maximal uncertainty level for $\mathcal{U}_{L_\infty}^{s, \tilde{\sigma}}(\cdot)$ as $\sigma_{\max}^\infty := \max_{p_1, p_2 \in \Delta(S)^A} \|p_1 - p_2\|_\infty = 1$. Then we can
 296 construct a collection of infinite-horizon RMDPs \mathcal{M}_{L_∞} defined by the uncertainty set with $\mathcal{U}_{L_\infty}^{s, \tilde{\sigma}}(\cdot)$
 297 so that for any $\sigma \in (0, \sigma_{\max}^\infty (1 - c_0)]$, and any dataset with in total N_{all} independent samples for all
 298 state-action pairs over the nominal transition kernel (for any RMDP inside \mathcal{M}_{L_∞}), one has*

$$\inf_{\hat{\pi}} \max_{\mathcal{M} \in \mathcal{M}_{L_\infty}} \left\{ \mathbb{P}_{\mathcal{M}} \left(\max_{s \in \mathcal{S}} [V^{*, \sigma}(s) - V^{\hat{\pi}, \sigma}(s)] > \varepsilon \right) \right\} \geq \frac{1}{8}, \quad (22)$$

299 provided that for c_8 is a universal positive constant,

$$N_{\text{all}} \leq \frac{c_8 SA}{(1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2}. \quad (23)$$

300 with $\mathbb{P}_{\mathcal{M}}$ the probability when the RMDP is \mathcal{M} , and the infimum is taken over all estimators $\hat{\pi}$.

301 Now we can present some implications of Theorem 3 and Theorem 4.

302 **• Robust MDPs with s -rectangularity are at least as easy as sa -rectangularity.** Theorem 3
 303 shows that the sample complexity of solving RMDPs with s -rectangularity does not exceed the
 304 order of $\tilde{O}\left(\frac{SA}{(1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} \varepsilon^2}\right)$. This matches the sample complexity for sa -rectangularity
 305 (cf. (18)) and indicates that although s -rectangular RMDPs are of a more complicated formulation,
 306 solving s -rectangular RMDPs is at least as easy as solving sa -rectangular RMDPs in terms of the
 307 sample complexity. In addition to the worst-case sample complexity upper bound, Theorem 3 also
 308 provides a data and instance-dependent sample complexity upper bound for s -rectangular RMDPs
 309 (cf. in (20)). Taking the divergence function $\|\cdot\| = L_p$ for instance, the data and instance-dependent
 310 sample complexity upper bound is

$$\begin{cases} \tilde{O}\left(\frac{SA}{(1-\gamma)^2 \varepsilon^2} \frac{1}{\max\{1-\gamma, \sigma\}}\right) & \text{if } \hat{\pi}_s(a|s) = \pi_s^*(a|s) = \frac{1}{A}, \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A} \\ \tilde{O}\left(\frac{SA}{(1-\gamma)^2 \varepsilon^2} \frac{1}{\max\{1-\gamma, \sigma A^{1/p}\}}\right) & \text{if } \|\hat{\pi}_s(\cdot|s)\|_0 = \|\pi_s^*(\cdot|s)\|_0 = 1, \quad \forall s \in \mathcal{S}. \end{cases}$$

311 where $\|\cdot\|_0$ corresponds to the total number of nonzero elements in a vector. The intuition beyond
 312 this theorem is that when the policy becomes proportional to uniform, the uncertainty budget of
 313 the s -rectangular MDPs is equally spread into all actions, and we retrieve the sa -rectangular case.
 314 When the policy becomes deterministic, all the uncertainty budget concentrates on one action. In
 315 this case, most of the actions are not robust except one, and the problem is simpler than classical
 316 MDP for this only specific action. An illustration of this result can be found in Fig. 2.

317 **• Comparisons with prior arts.** In Figure 2, we illustrate the comparisons with Clavier et al.
 318 [2023] which use L_p norms functions belonging to the class of general norms considered in this
 319 work. We do not compare in this section to Yang et al. [2022a] as it is not anymore state-of-the-art
 320 with regard to the work of Clavier et al. [2023]. In particular, the latest achieves in the s -rectangular
 321 case at sample complexity of $\tilde{O}\left(\frac{SA}{(1-\gamma)^3 \varepsilon^2}\right)$ in the regime where $\tilde{\sigma} \lesssim 1-\gamma$. In this regime, our result
 322 is the same but more general but in the regime where $\tilde{\sigma} \gtrsim 1-\gamma$, they achieve sample complexity
 323 of $\tilde{O}\left(\frac{SA}{(1-\gamma)^4 \varepsilon^2}\right)$ which is bigger than our result $\tilde{O}\left(\frac{SA}{(1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2}\right)$ by a factor at least $\frac{1}{1-\gamma}$.

324 5 Conclusion

325 This work refined sample complexity bounds to learn robust Markov decision processes when the
 326 uncertainty set is characterized by an general L_p metric, assuming the presence of a generative model.
 327 Our findings not only strengthen the current knowledge by improving both the upper and lower bounds,
 328 but also highlight that learning s -rectangular MDPs is less challenging in terms of sample complexity
 329 compared to classical sa -rectangular MDPs. This work is the first to provide results with a minimax
 330 bound, as prior results concerning s -rectangular cases were not minimax optimal. Additionally, we
 331 have established the minimax sample complexity for RMDPs using a general L_p norm, demonstrating
 332 that it is never larger than that required for learning standard MDPs. Our research identifies potential
 333 avenues for future work, such as exploring the characterization of tight sample complexity for RMDPs
 334 under a broader family of uncertainty sets, such as those defined by f -divergence. It would be highly
 335 desirable for a more unified theoretical foundation, as the distance between probability measures
 336 is more natural to define using divergence. Moreover, it would be interesting to focus on the finite-
 337 horizon Setting and linear setting, as our current analytical framework opens the door for potential ex-
 338 tensions to address finite-horizon RMDPs. Such an extension would contribute to a more comprehen-
 339 sive understanding of tabular cases. Finally, the case of linear MDPs would be interesting to explore.

340 References

- 341 Alekh Agarwal, Sham Kakade, and Lin F Yang. Model-based reinforcement learning with a generative
342 model is minimax optimal. In *Conference on Learning Theory*, pages 67–83. PMLR, 2020.
- 343 Mohammad Azar, Rémi Munos, and Hilbert J Kappen. Minimax pac bounds on the sample complexity
344 of reinforcement learning with a generative model. *Machine learning*, 91:325–349, 2013a.
- 345 Mohammad Gheshlaghi Azar, Rémi Munos, and Hilbert J Kappen. Minimax PAC bounds on the
346 sample complexity of reinforcement learning with a generative model. *Machine learning*, 91(3):
347 325–349, 2013b.
- 348 Kishan Panaganti Badrinath and Dileep Kalathil. Robust reinforcement learning using least squares
349 policy iteration with provable performance guarantees. In *International Conference on Machine
350 Learning*, pages 511–520. PMLR, 2021.
- 351 Yu Bai, Tengyang Xie, Nan Jiang, and Yu-Xiang Wang. Provably efficient Q-learning with low
352 switching cost. *arXiv preprint arXiv:1905.12849*, 2019.
- 353 Carolyn L Beck and Rayadurgam Srikant. Error bounds for constant step-size Q-learning. *Systems &
354 control letters*, 61(12):1203–1208, 2012.
- 355 Jose Blanchet and Karthyek Murthy. Quantifying distributional model risk via optimal transport.
356 *Mathematics of Operations Research*, 44(2):565–600, 2019.
- 357 Jose Blanchet, Miao Lu, Tong Zhang, and Han Zhong. Double pessimism is provably efficient
358 for distributionally robust offline reinforcement learning: Generic algorithm and robust partial
359 coverage. *arXiv preprint arXiv:2305.09659*, 2023.
- 360 Zaiwei Chen, Siva Theja Maguluri, Sanjay Shakkottai, and Karthikeyan Shanmugam. Finite-
361 sample analysis of stochastic approximation using smooth convex envelopes. *arXiv preprint
362 arXiv:2002.00874*, 2020.
- 363 Pierre Clavier, Stéphanie Allasonnière, and Erwan Le Pennec. Robust reinforcement learning with
364 distributional risk-averse formulation. *arXiv preprint arXiv:2206.06841*, 2022.
- 365 Pierre Clavier, Erwan Le Pennec, and Matthieu Geist. Towards minimax optimality of model-based
366 robust reinforcement learning. *arXiv preprint arXiv:2302.05372*, 2023.
- 367 Esther Derman and Shie Mannor. Distributional robustness and regularization in reinforcement
368 learning. *arXiv preprint arXiv:2003.02894*, 2020.
- 369 Esther Derman, Matthieu Geist, and Shie Mannor. Twice regularized MDPs and the equivalence
370 between robustness and regularization. *Advances in Neural Information Processing Systems*, 34,
371 2021.
- 372 Jing Dong, Jingwei Li, Baoxiang Wang, and Jingzhao Zhang. Online policy optimization for robust
373 MDP. *arXiv preprint arXiv:2209.13841*, 2022.
- 374 Kefan Dong, Yuanhao Wang, Xiaoyu Chen, and Liwei Wang. Q-learning with UCB exploration is
375 sample efficient for infinite-horizon MDP. *arXiv preprint arXiv:1901.09311*, 2019.
- 376 John Duchi and Hongseok Namkoong. Learning models with uniform performance via distributionally
377 robust optimization. *arXiv preprint arXiv:1810.08750*, 2018.
- 378 Rui Gao. Finite-sample guarantees for wasserstein distributionally robust optimization: Breaking the
379 curse of dimensionality. *arXiv preprint arXiv:2009.04382*, 2020.
- 380 Vineet Goyal and Julien Grand-Clement. Robust markov decision processes: Beyond rectangularity.
381 *Mathematics of Operations Research*, 2022.
- 382 Songyang Han, Sanbao Su, Sihong He, Shuo Han, Haizhao Yang, and Fei Miao. What is the solution
383 for state adversarial multi-agent reinforcement learning? *arXiv preprint arXiv:2212.02705*, 2022.
- 384 Chin Pang Ho, Marek Petrik, and Wolfram Wiesemann. Fast bellman updates for robust MDPs. In
385 *International Conference on Machine Learning*, pages 1979–1988. PMLR, 2018.

- 386 Chin Pang Ho, Marek Petrik, and Wolfram Wiesemann. Partial policy iteration for 11-robust markov
387 decision processes. *Journal of Machine Learning Research*, 22(275):1–46, 2021.
- 388 Garud N Iyengar. Robust dynamic programming. *Mathematics of Operations Research*, 30(2):
389 257–280, 2005.
- 390 Mehdi Jafarnia-Jahromi, Chen-Yu Wei, Rahul Jain, and Haipeng Luo. A model-free learning
391 algorithm for infinite-horizon average-reward MDPs with near-optimal regret. *arXiv preprint*
392 *arXiv:2006.04354*, 2020.
- 393 Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is Q-learning provably efficient?
394 In *Advances in Neural Information Processing Systems*, pages 4863–4873, 2018.
- 395 Chi Jin, Akshay Krishnamurthy, Max Simchowitz, and Tiancheng Yu. Reward-free exploration for
396 reinforcement learning. In *International Conference on Machine Learning*, pages 4870–4879.
397 PMLR, 2020.
- 398 Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline RL? In
399 *International Conference on Machine Learning*, pages 5084–5096, 2021.
- 400 William Karush. Minima of functions of several variables with inequalities as side conditions. In
401 *Traces and emergence of nonlinear programming*, pages 217–245. Springer, 2013.
- 402 David L Kaufman and Andrew J Schaefer. Robust modified policy iteration. *INFORMS Journal on*
403 *Computing*, 25(3):396–410, 2013.
- 404 Michael J Kearns and Satinder P Singh. Finite-sample convergence rates for Q-learning and indirect
405 algorithms. In *Advances in neural information processing systems*, pages 996–1002, 1999.
- 406 Olga Klopp, Karim Lounici, and Alexandre B Tsybakov. Robust matrix completion. *Probability*
407 *Theory and Related Fields*, 169(1-2):523–564, 2017.
- 408 Aounon Kumar, Alexander Levine, Tom Goldstein, and Soheil Feizi. Certifying model accuracy
409 under distribution shifts. *arXiv preprint arXiv:2201.12440*, 2022.
- 410 Navdeep Kumar, Esther Derman, Matthieu Geist, Kfir Levy, and Shie Mannor. Policy gradient for
411 s-rectangular robust markov decision processes. *arXiv preprint arXiv:2301.13589*, 2023.
- 412 Gen Li, Laixi Shi, Yuxin Chen, Yuantao Gu, and Yuejie Chi. Breaking the sample complexity barrier
413 to regret-optimal model-free reinforcement learning. *Advances in Neural Information Processing*
414 *Systems*, 34, 2021.
- 415 Gen Li, Yuejie Chi, Yuting Wei, and Yuxin Chen. Minimax-optimal multi-agent RL in Markov games
416 with a generative model. *Neural Information Processing Systems*, 2022a.
- 417 Gen Li, Laixi Shi, Yuxin Chen, Yuejie Chi, and Yuting Wei. Settling the sample complexity of
418 model-based offline reinforcement learning. *arXiv preprint arXiv:2204.05275*, 2022b.
- 419 Gen Li, Changxiao Cai, Yuxin Chen, Yuting Wei, and Yuejie Chi. Is Q-learning minimax optimal? a
420 tight sample complexity analysis. *Operations Research*, 2023a.
- 421 Gen Li, Yuting Wei, Yuejie Chi, and Yuxin Chen. Breaking the sample size barrier in model-based
422 reinforcement learning with a generative model. *accepted to Operations Research*, 2023b.
- 423 Gen Li, Yuling Yan, Yuxin Chen, and Jianqing Fan. Minimax-optimal reward-agnostic exploration in
424 reinforcement learning. *arXiv preprint arXiv:2304.07278*, 2023c.
- 425 Gen Li, Yuting Wei, Yuejie Chi, and Yuxin Chen. Breaking the sample size barrier in model-based
426 reinforcement learning with a generative model. *Operations Research*, 72(1):203–221, 2024.
- 427 Yan Li, Tuo Zhao, and Guanghui Lan. First-order policy optimization for robust markov decision
428 process. *arXiv preprint arXiv:2209.10579*, 2022c.
- 429 A Rupam Mahmood, Dmytro Korenkevych, Gautham Vasan, William Ma, and James Bergstra.
430 Benchmarking reinforcement learning algorithms on real-world robots. In *Conference on robot*
431 *learning*, pages 561–591. PMLR, 2018.

- 432 Shie Mannor, Duncan Simester, Peng Sun, and John N Tsitsiklis. Bias and variance in value function
433 estimation. In *Proceedings of the twenty-first international conference on Machine learning*,
434 page 72, 2004.
- 435 Colin McDiarmid et al. On the method of bounded differences. *Surveys in combinatorics*, 141(1):
436 148–188, 1989.
- 437 Janosch Moos, Kay Hansel, Hany Abdulsamad, Svenja Stark, Debora Clever, and Jan Peters. Robust
438 reinforcement learning: A review of foundations and recent advances. *Machine Learning and
439 Knowledge Extraction*, 4(1):276–315, 2022.
- 440 Arnab Nilim and Laurent El Ghaoui. Robust control of Markov decision processes with uncertain
441 transition matrices. *Operations Research*, 53(5):780–798, 2005.
- 442 Kishan Panaganti and Dileep Kalathil. Sample complexity of robust reinforcement learning with
443 a generative model. In *International Conference on Artificial Intelligence and Statistics*, pages
444 9582–9602. PMLR, 2022.
- 445 You Qiaoben, Xinning Zhou, Chengyang Ying, and Jun Zhu. Strategically-timed state-observation
446 attacks on deep reinforcement learning agents. In *ICML 2021 Workshop on Adversarial Machine
447 Learning*, 2021.
- 448 Hamed Rahimian and Sanjay Mehrotra. Distributionally robust optimization: A review. *arXiv
449 preprint arXiv:1908.05659*, 2019.
- 450 Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging offline
451 reinforcement learning and imitation learning: A tale of pessimism. *Neural Information Processing
452 Systems (NeurIPS)*, 2021.
- 453 Aurko Roy, Huan Xu, and Sebastian Pokutta. Reinforcement learning under model mismatch.
454 *Advances in neural information processing systems*, 30, 2017.
- 455 Walter Rudin et al. *Principles of mathematical analysis*, volume 3. McGraw-hill New York, 1964.
- 456 Reazul Hasan Russel, Bahram Behzadian, and Marek Petrik. Optimizing norm-bounded weighted
457 ambiguity sets for robust mdps. *arXiv preprint arXiv:1912.02696*, 2019.
- 458 Laixi Shi and Yuejie Chi. Distributionally robust model-based offline reinforcement learning with
459 near-optimal sample complexity. *arXiv preprint arXiv:2208.05767*, 2022.
- 460 Laixi Shi, Gen Li, Yuting Wei, Yuxin Chen, and Yuejie Chi. Pessimistic Q-learning for offline rein-
461 forcement learning: Towards optimal sample complexity. In *Proceedings of the 39th International
462 Conference on Machine Learning*, volume 162, pages 19967–20025. PMLR, 2022.
- 463 Laixi Shi, Gen Li, Yuting Wei, Yuxin Chen, Matthieu Geist, and Yuejie Chi. The curious price
464 of distributional robustness in reinforcement learning with a generative model. *arXiv preprint
465 arXiv:2305.16589*, 2023.
- 466 Aaron Sidford, Mengdi Wang, Xian Wu, Lin Yang, and Yinyu Ye. Near-optimal time and sample
467 complexities for solving Markov decision processes with a generative model. In *Advances in
468 Neural Information Processing Systems*, pages 5186–5196, 2018.
- 469 Elena Smirnova, Elvis Dohmatob, and Jérémie Mary. Distributionally robust reinforcement learning.
470 *arXiv preprint arXiv:1902.08708*, 2019.
- 471 Richard S Sutton. Learning to predict by the methods of temporal differences. *Machine learning*, 3
472 (1):9–44, 1988.
- 473 Aviv Tamar, Shie Mannor, and Huan Xu. Scaling up robust MDPs using function approximation. In
474 *International conference on machine learning*, pages 181–189. PMLR, 2014.
- 475 Kai Liang Tan, Yasaman Esfandiari, Xian Yeow Lee, and Soumik Sarkar. Robustifying reinforcement
476 learning agents via action space adversarial training. In *2020 American control conference (ACC)*,
477 pages 3959–3964. IEEE, 2020.

- 478 Chen Tessler, Yonathan Efroni, and Shie Mannor. Action robust reinforcement learning and applica-
479 tions in continuous control. In *International Conference on Machine Learning*, pages 6215–6224.
480 PMLR, 2019.
- 481 A. B. Tsybakov. *Introduction to nonparametric estimation*, volume 11. Springer, 2009.
- 482 J v. Neumann. Zur theorie der gesellschaftsspiele. *Mathematische annalen*, 100(1):295–320, 1928.
- 483 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*,
484 volume 47. Cambridge university press, 2018.
- 485 Martin J Wainwright. Stochastic approximation with cone-contractive operators: Sharp ℓ_∞ -bounds
486 for Q-learning. *arXiv preprint arXiv:1905.06265*, 2019.
- 487 Shengbo Wang, Nian Si, Jose Blanchet, and Zhengyuan Zhou. A finite sample complexity bound for
488 distributionally robust q-learning. *arXiv preprint arXiv:2302.13203*, 2023.
- 489 Yue Wang and Shaofeng Zou. Online robust reinforcement learning with model uncertainty. *Advances*
490 *in Neural Information Processing Systems*, 34, 2021.
- 491 Wolfram Wiesemann, Daniel Kuhn, and Berç Rustem. Robust markov decision processes. *Mathe-*
492 *matics of Operations Research*, 38(1):153–183, 2013.
- 493 Eric M Wolff, Ufuk Topcu, and Richard M Murray. Robust control of uncertain markov decision
494 processes with temporal logic specifications. In *2012 IEEE 51st IEEE Conference on Decision*
495 *and Control (CDC)*, pages 3372–3379. IEEE, 2012.
- 496 Tengyang Xie, Nan Jiang, Huan Wang, Caiming Xiong, and Yu Bai. Policy finetuning: Bridg-
497 ing sample-efficient offline and online reinforcement learning. *Advances in neural information*
498 *processing systems*, 34, 2021.
- 499 Huan Xu and Shie Mannor. Distributionally robust Markov decision processes. *Mathematics of*
500 *Operations Research*, 37(2):288–300, 2012.
- 501 Zaiyan Xu, Kishan Panaganti, and Dileep Kalathil. Improved sample complexity bounds for distribu-
502 tionally robust reinforcement learning. *arXiv preprint arXiv:2303.02783*, 2023.
- 503 Yuling Yan, Gen Li, Yuxin Chen, and Jianqing Fan. The efficacy of pessimism in asynchronous
504 Q-learning. *arXiv preprint arXiv:2203.07368*, 2022.
- 505 Yuling Yan, Gen Li, Yuxin Chen, and Jianqing Fan. The efficacy of pessimism in asynchronous
506 q-learning. *IEEE Transactions on Information Theory*, 2023.
- 507 Kunhe Yang, Lin Yang, and Simon Du. Q-learning with logarithmic regret. In *International*
508 *Conference on Artificial Intelligence and Statistics*, pages 1576–1584. PMLR, 2021.
- 509 Wei H Yang. On generalized holder inequality. 1991.
- 510 Wenhao Yang, Liangyu Zhang, and Zhihua Zhang. Toward theoretical understandings of robust
511 Markov decision processes: Sample complexity and asymptotics. *The Annals of Statistics*, 50(6):
512 3223–3248, 2022a.
- 513 Wenhao Yang, Liangyu Zhang, and Zhihua Zhang. Toward theoretical understandings of robust
514 markov decision processes: Sample complexity and asymptotics. *The Annals of Statistics*, 50(6):
515 3223–3248, 2022b.
- 516 Wenhao Yang, Han Wang, Tadashi Kozuno, Scott M Jordan, and Zhihua Zhang. Avoiding
517 model estimation in robust markov decision processes with a generative model. *arXiv preprint*
518 *arXiv:2302.01248*, 2023.
- 519 Ming Yin, Yu Bai, and Yu-Xiang Wang. Near-optimal offline reinforcement learning via double
520 variance reduction. *arXiv preprint arXiv:2102.01748*, 2021.
- 521 Huan Zhang, Hongge Chen, Chaowei Xiao, Bo Li, Mingyan Liu, Duane Boning, and Cho-Jui
522 Hsieh. Robust deep reinforcement learning against adversarial perturbations on state observations.
523 *Advances in Neural Information Processing Systems*, 33:21024–21037, 2020a.

- 524 Huan Zhang, Hongge Chen, Duane Boning, and Cho-Jui Hsieh. Robust reinforcement learning on
525 state observations with learned optimal adversary. *arXiv preprint arXiv:2101.08452*, 2021.
- 526 Zihan Zhang, Yuan Zhou, and Xiangyang Ji. Almost optimal model-free reinforcement learning
527 via reference-advantage decomposition. *Advances in Neural Information Processing Systems*, 33,
528 2020b.
- 529 Zhengqing Zhou, Qinxun Bai, Zhengyuan Zhou, Linhai Qiu, Jose Blanchet, and Peter Glynn.
530 Finite-sample regret bound for distributionally robust offline tabular reinforcement learning. In
531 *International Conference on Artificial Intelligence and Statistics*, pages 3331–3339. PMLR, 2021.

532 6 Other related works

533 Here we provide additional discussion of related work that could not be fit into the main paper due
534 to space considerations. We limit our discussions to the tabular setting with finite state and action
535 spaces provable RL algorithms.

536 **Classical reinforcement learning with finite-sample guarantees.** A recent surge in attention
537 for RL has leveraged the methodologies derived from high-dimensional probability and statistics
538 to analyze RL algorithms in non-asymptotic scenarios. Substantial efforts have been devoted to
539 conducting non-asymptotic sample analyses of standard RL in many settings. Illustrative instances
540 encompass investigations employing Probably Approximately Correct (PAC) bounds in the context
541 of *generative model* settings [Kearns and Singh, 1999, Beck and Srikant, 2012, Li et al., 2022a, Chen
542 et al., 2020, Azar et al., 2013b, Sidford et al., 2018, Agarwal et al., 2020, Li et al., 2023a,b, Wainwright,
543 2019] and the *online setting* via both in PAC-base or regret-based analyses [Jin et al., 2018, Bai
544 et al., 2019, Li et al., 2021, Zhang et al., 2020b, Dong et al., 2019, Jin et al., 2020, Li et al., 2023c,
545 Jafarnia-Jahromi et al., 2020, Yang et al., 2021] and finally *offline setting* [Rashidinejad et al., 2021,
546 Xie et al., 2021, Yin et al., 2021, Shi et al., 2022, Li et al., 2022b, Jin et al., 2021, Yan et al., 2022].

547 **Robustness in reinforcement learning.** Reinforcement learning has had notable achievements
548 but has also exhibited significant limitations, particularly when the learned policy is susceptible
549 to deviations in the deployed environment due to perturbations, model discrepancies, or structural
550 modifications. To address these challenges, the idea of robustness in RL algorithms has been studied.
551 Robustness could concern uncertainty or perturbations across different Markov Decision Processes
552 (MDPs) components, encompassing reward, state, action, and the transition kernel. Moos et al. [2022]
553 gives a recent overview of the different work in this field.

554 The distributionally robust MDP (RMDP) framework has been proposed [Iyengar, 2005] to enhance
555 the robustness of RL has been proposed. In addition to this work, various other research efforts,
556 including, but not limited to, Zhang et al. [2020a, 2021], Han et al. [2022], Clavier et al. [2022],
557 Qiaoben et al. [2021], explore robustness regarding state uncertainty. In these scenarios, the agent’s
558 policy is determined on the basis of perturbed observations generated from the state, introducing
559 restricted noise, or undergoing adversarial attacks. Finally, robustness considerations extend to
560 uncertainty in the action domain. Works such as Tessler et al. [2019], Tan et al. [2020] consider
561 the robustness of actions, acknowledging potential distortions introduced by an adversarial agent.

562 Given the focus of our work, we provide a more detailed background on progress related to distribu-
563 tionally robust RL. The idea of distributionally robust optimization has been explored within the con-
564 text of supervised learning [Rahimian and Mehrotra, 2019, Gao, 2020, Duchi and Namkoong, 2018,
565 Blanchet and Murthy, 2019] and has also been extended to distributionally robust dynamic program-
566 ming and Distributionally Robust Markov Decision Processes (DRMDPs) such as in [Iyengar, 2005,
567 Xu and Mannor, 2012, Wolff et al., 2012, Kaufman and Schaefer, 2013, Ho et al., 2018, Smirnova et al.,
568 2019, Ho et al., 2021, Goyal and Grand-Clement, 2022, Derman and Mannor, 2020, Tamar et al., 2014,
569 Badrinath and Kalathil, 2021]. Despite the considerable attention received, both empirically and theo-
570 retically, most previous theoretical analyses in the context of RMDPs adopt an asymptotic perspective
571 [Roy et al., 2017] or focus on planning with exact knowledge of the uncertainty set [Iyengar, 2005, Xu
572 and Mannor, 2012, Tamar et al., 2014]. Many works have focused on the finite-sample performance
573 of verifiable robust Reinforcement Learning (RL) algorithms. These investigations encompass various
574 data generation mechanisms and uncertainty set formulations over the transition kernel. Closely
575 related to our work, various forms of uncertainty sets have been explored, showcasing the versatility
576 of approaches. Divergence such as Kullback-Leibler (KL) divergence is another prevalent choice,
577 extensively studied by Yang et al. [2022a], Panaganti and Kalathil [2022], Zhou et al. [2021], Shi and
578 Chi [2022], Xu et al. [2023], Wang et al. [2023], Blanchet et al. [2023], who investigated the sample
579 complexity of both model-based and model-free algorithms in simulator or offline settings. Xu et al.
580 [2023] considered various uncertainty sets, including those associated with the Wasserstein distance.
581 The introduction of an R-contamination uncertainty set Wang and Zou [2021], has been proposed to
582 tackle a robust Q-learning algorithm for the online setting, with guarantees analogous to standard RL.
583 Finally, the finite-horizon scenario has been studied by Xu et al. [2023], Dong et al. [2022] with finite-
584 sample complexity bounds for (RMDPs) using TV and χ^2 divergence. More broadly, other related
585 topics have been explored, such as the iteration complexity of policy-based methods [Li et al., 2022c,
586 Kumar et al., 2023], and regularization-based robust RL [Yang et al., 2023]. Finally, Badrinath and

587 Kalathil [2021] examined a general sa -rectangular form of the uncertainty set, proposing a model-free
 588 algorithm for the online setting with linear function approximation to address large state spaces.

589 7 Discussion on hypothesis of Theorems 1 and 3.

- 590 • *What norms are included in the Definition 1?* In our upper bound result Theorems
 591 3 and 1, we upper bound the sample complexity for C^2 norms and TV. The set of C^2
 592 smooth norm is very large as it includes all, L_p norm, weighted, rescaled L_p norms for
 593 $p \geq 2$. Weighted norms can be useful in practice, to get more weights on dangerous
 594 specific states in Robust MDPs formulation such as in Russel et al. [2019]. Moreover, note
 595 that our result can generalize to metric or pseudo metric (which are not homogeneous ie
 596 $\|\lambda\| = |\lambda| \|x\| \forall x \in \mathbb{R}^n, \lambda \in \mathbb{R}$) with norms of the form $x \mapsto \phi^{-1}(\sum_{k=1}^n \phi(|x_k|))$ with
 597 ϕ a convex incising function such as the norm is still positive, definite positive. Choosing
 598 $\phi(x) = x^p$ leads to the L_p norms.
- 599 • *Assumptions on γ in Theorems 1 and 3, and Assumptions on γ for lower bound.* When
 600 γ is small (e.g., $\gamma \in (0, \frac{1}{2}]$) leads to the effective horizon length is at most 2), the sequential
 601 structure almost disappears and is much less of interest for RL community. So people Li
 602 et al. [2023b] Yan et al. [2023] usually focus on reasonable range $\gamma \in (c, 1)$ for some small
 603 positive constant c , such as $\gamma \in [\frac{1}{2}, 1)$. However, the theorems can be directly extended
 604 to a broader range of $\gamma \in (c, 1)$ along with c as small as desired so that almost cover the
 605 full range $(0, 1)$.
- 606 • *Why final results on s depend on $\hat{\pi}$*
 607 Theorem 3 is $\hat{\pi}$ data dependent which is randomness-dependent measure. However, taking
 608 the minimum of this quantity leads to the same bound as is sa -rectangular, so to illustrate
 609 that it is possible to get tighter bounds for s -rectangular with instance-dependent RMDPs,
 610 we decide to write also randomness-dependent quantity, while the less tight upper bound
 611 is written also in the theorem, taking the first term in the “min” in (21).
- 612 • *Why our results are still true for TV?* Theorems 1 and 3 are stated for C^2 smooth norms,
 613 however, our result is still true for TV which is not C^2 as in this specific case, the dual
 614 of the optimization problem becomes a 1-dimensional problem. In this case in the main
 615 concentration lemma 8, the additional term involving smoothness term denoted C_S is not
 616 present and the bound is simpler as is not required this additional term.
- 617 • *Why burn-in or sufficiently small ϵ condition is not too restrictive?* The burn-in term in
 618 Th. 1 and 3 is proportional to $1/\epsilon$ where the "sample complexity" term is proportional to
 619 $1/\epsilon^2$. The smooth term depending on C_S or burn-in is then not too large for sufficiently
 620 small ϵ compared to the other term, which will give final the sample complexity.
- 621 • *Why this is not extendable to f -divergence currently?* The f -divergence as a distinct family of
 622 divergence is beyond the scope of this paper. Current proof for arbitrary norms cannot be di-
 623 rectly extended since the key phenomenon of shrinking range of the robust value function has
 624 not been verified for f -divergence yet, while it is promising as an interesting future direction.

625 8 Preliminaries

626 These quantities appear in the dual formulation of the robust optimization problem and more pre-
 627 cisely the dual span semi norm $\text{sp}(\cdot)_*$ note that for L_2 , we retrieve the classical mean with the
 628 definition of ω) With slight abuse of notation, we denote 0 (resp. 1) as the all-zero (resp. all-one)
 629 vector. We then introduce the notation $[T] := \{1, \dots, T\}$ for any positive integer $T > 0$. Then, for
 630 two vectors $x = [x_i]_{1 \leq i \leq n}$ and $y = [y_i]_{1 \leq i \leq n}$, the notation $x \leq y$ (resp. $x \geq y$) means $x_i \leq y_i$
 631 (resp. $x_i \geq y_i$) for all $1 \leq i \leq n$. Finally, for any vector x , we overload the notation by letting
 632 $x^{\circ 2} = [x(s, a)^2]_{(s, a) \in \mathcal{S} \times \mathcal{A}}$ (resp. $x^{\circ 2} = [x(s)^2]_{s \in \mathcal{S}}$), Finally, we drop the subscript $\|\cdot\|$ to write
 633 $\mathcal{U}_{\|\cdot\|}^\sigma(\cdot) = \mathcal{U}^\sigma(\cdot)$ for both sa - and s - rectangular assumptions.

634 **Matrix and Vector Notations.** Throughout the analysis, we need to introduce or recall some matrix
 635 and vector notations in the following.

- 636 • $r \in \mathbb{R}^{\mathcal{S}\mathcal{A}}$: the reward function vector r (so that $r_{(s, a)} = r(s, a)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$).

637
638
639
640
641

- $P^0 \in \mathbb{R}^{SA \times S}$: the nominal transition kernel matrix with $P_{s,a}^0$ as the (s, a) -th row.
- $\hat{P}^0 \in \mathbb{R}^{SA \times S}$: the estimated nominal transition kernel matrix with $\hat{P}_{s,a}^0$ as the (s, a) -th row.
- $\Pi^\pi \in \{0, 1\}^{S \times SA}$: a projection matrix associated with a given policy π taking the following form:

$$\Pi^\pi = \begin{pmatrix} 1_{\pi(1)}^\top & 0^\top & \cdots & 0^\top \\ 0^\top & 1_{\pi(2)}^\top & \cdots & 0^\top \\ \vdots & \vdots & \ddots & \vdots \\ 0^\top & 0^\top & \cdots & 1_{\pi(S)}^\top \end{pmatrix}, \quad (24)$$

642

where $1_{\pi(1)}^\top, 1_{\pi(2)}^\top, \dots, 1_{\pi(S)}^\top \in \mathbb{R}^A$ are simplex vector such as

$$1_{\pi(1)}^\top = (\pi(a_1|s_1), \pi(a_2|s_1), \dots, \pi(a_A|s_1)).$$

643
644
645
646

- $P^V \in \mathbb{R}^{SA \times S}$, $\hat{P}^V \in \mathbb{R}^{SA \times S}$ are the matrices representing the probability transition kernel in the uncertainty set that leads to the worst-case value for any vector $V \in \mathbb{R}^S$. We denote $P_{s,a}^V$ (resp. $\hat{P}_{s,a}^V$) as the (s, a) -th row of the transition matrix P^V (resp. \hat{P}^V). The (s, a) -th rows of these transition matrices are defined for sa -rectangular assumptions as

$$P_{s,a}^V = \operatorname{argmin}_{\mathcal{P} \in \mathcal{U}^{sa, \sigma}(P_{s,a}^0)} \mathcal{P}V, \quad \text{and} \quad \hat{P}_{s,a}^V = \operatorname{argmin}_{\mathcal{P} \in \mathcal{U}^{sa, \sigma}(\hat{P}_{s,a}^0)} \mathcal{P}V. \quad (25a)$$

647

Moreover, we will use of the following shorthand notation:

$$P_{s,a}^{\pi, V} := P_{s,a}^{V, \pi, \sigma} = \operatorname{argmin}_{\mathcal{P} \in \mathcal{U}^{sa, \sigma}(P_{s,a}^0)} \mathcal{P}V^{\pi, \sigma}, \quad P_{s,a}^{\pi, \hat{V}} := P_{s,a}^{\hat{V}, \pi, \sigma} = \operatorname{argmin}_{\mathcal{P} \in \mathcal{U}^{sa, \sigma}(P_{s,a}^0)} \mathcal{P}\hat{V}^{\pi, \sigma}, \quad (25b)$$

$$\hat{P}_{s,a}^{\pi, V} := \hat{P}_{s,a}^{V, \pi, \sigma} = \operatorname{argmin}_{\mathcal{P} \in \mathcal{U}^{sa, \sigma}(\hat{P}_{s,a}^0)} \mathcal{P}V^{\pi, \sigma}, \quad \hat{P}_{s,a}^{\pi, \hat{V}} := \hat{P}_{s,a}^{\hat{V}, \pi, \sigma} = \operatorname{argmin}_{\mathcal{P} \in \mathcal{U}^{sa, \sigma}(\hat{P}_{s,a}^0)} \mathcal{P}\hat{V}^{\pi, \sigma}. \quad (25c)$$

648
649

The corresponding probability transition matrices are denoted by $P^{\pi, V} \in \mathbb{R}^{SA \times S}$, $P^{\pi, \hat{V}} \in \mathbb{R}^{SA \times S}$, $\hat{P}^{\pi, V} \in \mathbb{R}^{SA \times S}$ and $\hat{P}^{\pi, \hat{V}} \in \mathbb{R}^{SA \times S}$, respectively.

650
651

- $P^\pi \in \mathbb{R}^{S \times S}$, $\hat{P}^\pi \in \mathbb{R}^{S \times S}$, $\underline{P}^{\pi, V} \in \mathbb{R}^{S \times S}$, $\underline{P}^{\pi, \hat{V}} \in \mathbb{R}^{S \times S}$, $\hat{\underline{P}}^{\pi, V} \in \mathbb{R}^{S \times S}$ and $\hat{\underline{P}}^{\pi, \hat{V}} \in \mathbb{R}^{S \times S}$: six *square* probability transition matrices w.r.t. policy π over the states, namely

$$\begin{aligned} P^\pi &:= \Pi^\pi P^0, & \hat{P}^\pi &:= \Pi^\pi \hat{P}^0, & \underline{P}^{\pi, V} &:= \Pi^\pi P^{\pi, V}, & \underline{P}^{\pi, \hat{V}} &:= \Pi^\pi P^{\pi, \hat{V}}, \\ \hat{\underline{P}}^{\pi, V} &:= \Pi^\pi \hat{P}^{\pi, V}, & \text{and} & & \hat{\underline{P}}^{\pi, \hat{V}} &:= \Pi^\pi \hat{P}^{\pi, \hat{V}}. \end{aligned} \quad (26)$$

652
653
654

For s -rectangular, we will use the same notation for these transition matrices, removing a subscript for s -rectangular assumptions. Finally, we denote P_s^π as the s -th row of the transition matrix P^π .

655
656
657

- $r_\pi \in \mathbb{R}^S$: a reward restricted to the actions chosen by the policy vector π , $r_\pi = \Pi^\pi r$.
- $\operatorname{Var}_P(V) \in \mathbb{R}^{SA}$: for a given transition kernel $P \in \mathbb{R}^{SA \times S}$ and vector $V \in \mathbb{R}^S$, we denote the (s, a) -th row of $\operatorname{Var}_P(V)$ as

$$\operatorname{Var}_P(s, a) := \operatorname{Var}_{P_{s,a}}(V). \quad (27)$$

658

8.1 Additional definitions and basic facts

659

For any norm smooth $\|\cdot\|$ introduced in 1, we define the span semi norm as

660
661

Definition 2 (Span semi norm). *Given any norm $\|\cdot\|$, we define the span semi norm as: $\operatorname{sp}(x) = \min_{\omega \in \mathbb{R}} \|v - \omega \mathbf{1}\|$ and the generalized mean as $\omega(x) := \arg \min_{\omega \in \mathbb{R}} \|x - \omega \mathbf{1}\|$.*

662

Let vector $P \in \mathbb{R}^{1 \times S}$ and vector $V \in \mathbb{R}^S$, we define the variance

$$\operatorname{Var}_P(V) := P(V \circ V) - (PV) \circ (PV). \quad (28)$$

663

The following lemma bounds the Lipschitz constant of the variance function.

664 **Lemma 1.** (Shi et al. [2023], Lemma 2) Assuming $0 \leq V_1, V_2 \leq \frac{1}{1-\gamma}$ which obey $\|V_1 - V_2\|_\infty \leq x$
665 , then for $P \in \Delta(S)$, one has

$$|\text{Var}_P(V_1) - \text{Var}_P(V_2)| \leq \frac{2x}{(1-\gamma)}. \quad (29)$$

666 **Lemma 2.** [Panaganti and Kalathil, 2022, Lemma 6] Consider any $\delta \in (0, 1)$. For any fixed policy
667 π and fixed value vector $V \in \mathbb{R}^S$, one has with probability at least $1 - \delta$,

$$\left| \sqrt{\text{Var}_{\hat{P}^\pi}(V)} - \sqrt{\text{Var}_{P^\pi}(V)} \right| \leq \sqrt{\frac{2\|V\|_\infty^2 \log\left(\frac{2SA}{\delta}\right)}{N}} 1.$$

668 8.2 Empirical robust MDP $\widehat{\mathcal{M}}_{\text{rob}}$ Bellman equations

669 We define the robust MDP $\widehat{\mathcal{M}}_{\text{rob}} = \{\mathcal{S}, \mathcal{A}, \gamma, \mathcal{U}^\sigma(\widehat{P}^0), r\}$ based on the estimated nominal distribution
670 \widehat{P}^0 in (11). Then, we denote the associated robust value function (resp. robust Q-function) are $\widehat{V}^{\pi, \sigma}$
671 (resp. $\widehat{Q}^{\pi, \sigma}$). We can notice that that $\widehat{Q}^{*, \sigma}$ is the unique-fixed point of $\widehat{\mathcal{T}}^\sigma(\cdot)$ (see Lemma 8.3), the
672 empirical robust Bellman operator constructed using \widehat{P}^0 . Finally, similarly to (9), for $\widehat{\mathcal{M}}_{\text{rob}}$, the
673 Bellman's optimality principle gives the following *robust Bellman consistency equation* (resp. *robust*
674 *Bellman optimality equation*) for *sa*-rectangular assumptions:

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \widehat{Q}^{\pi, \sigma}(s, a) = r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(\widehat{P}_{s, a}^0)} \mathcal{P} \widehat{V}^{\pi, \sigma}, \quad (30a)$$

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \widehat{Q}^{*, \sigma}(s, a) = r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(\widehat{P}_{s, a}^0)} \mathcal{P} \widehat{V}^{*, \sigma}. \quad (30b)$$

675 Using matrix notation, we can write the robust Bellman consistency equations as

$$Q^{\pi, \sigma} = r + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(P^0)} \mathcal{P} V^{\pi, \sigma} \quad \text{and} \quad \widehat{Q}^{\pi, \sigma} = r + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(\widehat{P}^0)} \mathcal{P} \widehat{V}^{\pi, \sigma}, \quad (31)$$

676 which imply

$$\begin{aligned} V^{\pi, \sigma} &= r_\pi + \gamma \Pi^\pi \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(P^0)} \mathcal{P} V^{\pi, \sigma} \stackrel{(i)}{=} r_\pi + \gamma \underline{P}^{\pi, V} V^{\pi, \sigma}, \\ \widehat{V}^{\pi, \sigma} &= r_\pi + \gamma \Pi^\pi \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(\widehat{P}^0)} \mathcal{P} \widehat{V}^{\pi, \sigma} \stackrel{(ii)}{=} r_\pi + \gamma \underline{\widehat{P}}^{\pi, \widehat{V}} \widehat{V}^{\pi, \sigma}, \end{aligned} \quad (32)$$

677 where (i) and (ii) hold by the definitions in (24), (25) and (26). For *s*-rectangular, we can define the
678 same notation, removing a subscript:

$$\begin{aligned} V^{\pi, \sigma} &= r_\pi + \gamma \Pi^\pi \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(P^0)} \mathcal{P} V^{\pi, \sigma} \stackrel{(i)}{=} r_\pi + \gamma \underline{P}^{\pi, V} V^{\pi, \sigma}, \\ \widehat{V}^{\pi, \sigma} &= r_\pi + \gamma \Pi^\pi \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(\widehat{P}^0)} \mathcal{P} \widehat{V}^{\pi, \sigma} \stackrel{(ii)}{=} r_\pi + \gamma \underline{\widehat{P}}^{\pi, \widehat{V}} \widehat{V}^{\pi, \sigma}, \end{aligned} \quad (33)$$

679 8.3 Properties of the robust Bellman operator and dual representation

680 The robust Bellman operator (cf. (10)) shares the γ -contraction property of the standard Bellman
681 operator as:

682 **[Iyengar, 2005, Theorem 3.2]** Given $\gamma \in [0, 1)$, the robust Bellman operator $\mathcal{T}^\sigma(\cdot)$ (cf. (10)) is a
683 γ -contraction w.r.t. $\|\cdot\|_\infty$. More formally, for any $Q_1, Q_2 \in \mathbb{R}^{SA}$ s.t. $Q_1(s, a), Q_2(s, a) \in [0, \frac{1}{1-\gamma}]$
684 for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, one has

$$\|\mathcal{T}^\sigma(Q_1) - \mathcal{T}^\sigma(Q_2)\|_\infty \leq \gamma \|Q_1 - Q_2\|_\infty. \quad (34)$$

685 It can be also shown that, $Q^{*,\sigma}$ is the unique fixed point of $\mathcal{T}^\sigma(\cdot)$ obeying $0 \leq Q^{*,\sigma}(s, a) \leq \frac{1}{1-\gamma}$ for
686 all $(s, a) \in \mathcal{S} \times \mathcal{A}$.

687 One of the main contributions is to derive the dual form of optimization problem using arbitrary
688 norms. These lemma take ideas from Iyengar [2005] and are adapted to arbitrary norms and not only
689 TV distance.

690 **Dual equivalence of the robust Bellman operator.** Fortunately, the robust Bellman operator can
691 be evaluated efficiently by resorting to its dual formulation, and this idea is central in all proofs for
692 RMPDs. Dual formulation of RMPDs have been introduced in [Iyengar, 2005] but the proof was
693 done uniquely for the TV and the χ^2 case. Before continuing, for any $V \in \mathbb{R}^S$, we denote $[V]_\alpha$ as
694 its clipped version by some non-negative vector α , namely,

$$[V]_\alpha(s) := \begin{cases} \alpha, & \text{if } V(s) > \alpha(s), \\ V(s), & \text{otherwise.} \end{cases} \quad (35)$$

695 Defining the gradient of $P \mapsto \|P\|$ as $\nabla \|P\|$, $\lambda > 0$, a positive scalar and ω is the generalized mean
696 defined as the argmin in the definition of the span semi norm in Def.2, we derive two optimization
697 lemmas.

698 **Lemma 3** (Strong duality using norm $\|\cdot\|$ in the sa -rectangular case.). *Consider any probability*
699 *vector $P \in \Delta(\mathcal{S})$ and any fixed uncertainty level σ , we abbreviate the notation of the uncertainty set*
700 *$\mathcal{U}_{\|\cdot\|}^{sa,\sigma}(P)$ (cf. (3)) as $\mathcal{U}^{sa,\sigma}(P)$. For any vector $V \in \mathbb{R}^S$ obeying $V \geq 0$, recalling the definition of*
701 *$[V]_\alpha$ in (35), one has*

$$\inf_{P \in \mathcal{U}^{sa,\sigma}(P)} \mathcal{P}V = \max_{\mu_P^{\lambda,\omega} \in \mathcal{M}_P^{\lambda,\omega}} \left\{ P(V - \mu_P^{\lambda,\omega}) - \sigma \left(\text{sp}((V - \mu_P^{\lambda,\omega})_*) \right) \right\}. \quad (36)$$

$$= \max_{\alpha_P^{\lambda,\omega} \in \mathcal{A}_P^{\lambda,\omega}} \left\{ P[V]_{\alpha_P^{\lambda,\omega}} - \sigma \left(\text{sp}([V]_{\alpha_P^{\lambda,\omega}})_* \right) \right\} \quad (37)$$

702 where $\text{sp}(\cdot)_*$ is defined in Def.2. Here, the two auxiliary variational family $\mathcal{A}_P^{\lambda,\omega}$, $\mathcal{M}_P^{\lambda,\omega}$ are defined
703 as below:

$$\mathcal{A}_P^{\lambda,\omega} = \left\{ \alpha_P^{\lambda,\omega} : \alpha_P^{\lambda,\omega}(s) = \omega + \lambda |\nabla \|P\|(s) : \lambda > 0, \omega > 0, P \in \Delta(\mathcal{S}), \alpha_P^{\lambda,\omega} \in \left[0, \frac{1}{1-\gamma} \right]^S \right\} \quad (38)$$

$$\mathcal{M}_P^{\lambda,\omega} = \left\{ \mu_P^{\lambda,\omega} = V - \alpha_P^{\lambda,\omega}, \lambda, \omega \in \mathbb{R}^+, P \in \Delta(\mathcal{S}), \mu \in \mathbb{R}_+^S, \mu_P^{\lambda,\omega} = \left[0, \frac{1}{1-\gamma} \right]^S \right\} \quad (39)$$

$$(40)$$

704 For L_1 or TV , case, the vector $\alpha_P^{\lambda,\omega}$ reduces to a 1 dimensional scalar such as $\alpha \in [0, 1/(1-\gamma)]$.

Proof.

$$\begin{aligned} \inf_{P \in \mathcal{U}^{sa,\sigma}(P)} \mathcal{P}V &= \inf_{\{P: P \in \Delta_s, \|P-P\| \leq \sigma\}} \sum_{s'} \mathcal{P}(s')V(s') \\ &= PV + \inf_{\{y: \|y\| \leq \sigma, 1y=0, y \geq -P\}} \sum_{s'} y(s')V(s') \end{aligned}$$

705 where we use the change of variable $y(s') = \mathcal{P}(s') - P(s')$ for all $s' \in \mathcal{S}$. Then the Lagrangian
706 function of the above optimization problem can be written as follows:

$$\inf_{P \in \mathcal{U}_{s,a}^\sigma(P)} \mathcal{P}V = PV + \sup_{\mu \geq 0, \nu \in \mathbb{R}} \inf_{\{y: \|y\| \leq \sigma\}} - \sum_{s'} \mu(s)P(s') + \sum_{s'} (y(s')(V(s') - \mu(s') - \nu)) \quad (41)$$

$$\stackrel{(a)}{=} PV + \sup_{\mu \geq 0, \nu \in \mathbb{R}} - \sum_{s'} \mu(s')P(s') - \sigma \|(V(s') - \mu(s') - \nu \mathbf{1})\|_* \quad (42)$$

$$\stackrel{(b)}{=} \sup_{\mu \geq 0} P(V - \mu) - \sigma \text{sp}(V - \mu)_* \quad (43)$$

707 where $\mu \in \mathbb{R}_+^S$, $\nu \in \mathbb{R}$ are Lagrangian variables, (a) is true using the equality case of Cauchy-Swartz
708 inequality for dual norm Yang [1991], and (b) is due to is the definition of the span semi-norm (see
709 (8)). The value that maximizes the inner maximization problem in (42) in $\omega(V, \mu)$ is the generalized-
710 mean by definition denoted with abbreviate notation ω . If the norm is differentiable, then we have
711 that the equality (a) comes from the generalized Holder's inequality for arbitrary norms Yang [1991],
712 namely, defining $z = (V - \mu - \omega)$, it satisfies

$$z = \|z\|_* \nabla \|y\| \quad (44)$$

713 The quantity ν is replaced by the generalized mean for equality in (b) while (44) comes from Yang
714 [1991]. Using complementary slackness Karush [2013]stackness let $\mathcal{B} = \{s \in \mathcal{S} : \mu(s) > 0\}$

$$\forall s \in \mathcal{B} : \quad y^*(s) = -P(s), \quad (45)$$

715 which leads to the following equality by plugging the previous (45) in (44) and defining $z^* =$
716 $V - \mu^* - \omega$:

$$\forall s \in \mathcal{B}, \quad z^*(s) = \|z^*\|_* \nabla \|P\| (s) \quad (46)$$

717 or

$$\forall s \in \mathcal{B}, \quad V(s) - \mu^*(s) = \omega + \lambda \nabla \|P\| (s) \hat{=} \alpha_P^{\lambda, \omega} \quad (47)$$

718 by letting $\lambda = \|z^*\|_* \in \mathbb{R}^+$. Note that here the hypothesis of 1 are use and especially separability is
719 needed to ensure that for $s \in \mathcal{B}$, $\nabla \|y\| = \nabla \|P\|$ only depend on $P(s)$ and not on other coordinates,
720 which is true form generalized L_p norms. We can remark that $v - \mu^*$ is P dependent, but if P is
721 known, the best μ^* is only determined by one 2 dimensional parameters $\lambda = \|v - \mu^* - \nu\|_*$ and
722 $\omega \in \mathbb{R}^+$. Moreover, when P is fixed, the scalar ω is a constant is fully determined by P , v and μ^* .
723 This is why the quantity defined α_P^λ varies through 2 parameter λ and ω . Given this observation, we
724 can rewrite the optimization problem as :

$$\sup_{\mu \geq 0} P(V - \mu) - \sigma \text{sp}(V - \mu)_* = \sup_{\mu_P^{\lambda, \omega} \in \mathcal{M}_P^{\lambda, \omega}} P(V - \mu_P^{\lambda, \omega}) - \sigma \text{sp}((V - \mu_P^{\lambda, \omega}))_* \quad (48)$$

$$= \sup_{\alpha_P^{\lambda, \omega} \in \mathcal{A}_P^{\lambda, \omega}} P[V]_{\alpha_P^{\lambda, \omega}} - \sigma \text{sp}([V]_{\alpha_P^{\lambda, \omega}})_* \quad (49)$$

where we defined the maximization problem on μ not in \mathbb{R}^S but at the optimal in the variational
family denote $\mathcal{M}_P^{\lambda, \omega} = \{v - \alpha_P^{\lambda, \omega}, (\lambda, \omega) \in \mathbb{R}_+^2, P \in \Delta(S)\}$. We can rewrite the optimization
problem in terms of α_P with

$$[V]_{\alpha_P^{\lambda, \omega}}(s) := \begin{cases} \alpha_P^{\lambda, \omega}, \\ V(s), \quad \text{otherwise.} \end{cases}$$

725 Contrary to the TV case, α is not a scalar but $\alpha_P^{\lambda, \omega}$ belongs to a variational family only determined
726 by two parameter. Note that this lemma is still true writing subgradient and not gradient of P . As
727 we assume C^2 -regularity on norms, the subgradient space of the norm reduce to the singleton of the
728 gradient in our case. C^2 smoothness will be needed in concentration part while it is possible to be
729 more general in optimization lemmas. Note that for TV or L_1 , this lemma holds, but the vector $\alpha_P^{\lambda, \omega}$
730 reduces to a positive scalar denoted α which is equal to $\|v - \mu^*\|_\infty$ according to Iyengar [2005]

731

□

732 **Lemma 4** (Strong duality for the distance induced by the norm $\|\cdot\|$ in the s -rectangular case.).
733 Consider any probability vector $P^\pi := \Pi^\pi P \in \Delta_s$ for $P \in \Delta(S)^{\mathcal{A}}$, any fixed uncertainty level $\tilde{\sigma}$
734 and the uncertainty set $\mathcal{U}_{\|\cdot\|}^{s, \tilde{\sigma}}(P)$, we abbreviate the subscript to use $\mathcal{U}^{s, \tilde{\sigma}}(P) := \mathcal{U}_{\|\cdot\|}^{s, \tilde{\sigma}}(P)$. Then for
735 any vector $V \in \mathbb{R}^S$ obeying $V \geq 0$, recalling the definition of $[V]_\alpha$ in (35), one has

$$\inf_{P \in \mathcal{U}^{s, \tilde{\sigma}}(P)} P^\pi V = \sum_a \pi(a|s) \left(\max_{\alpha_{P_{sa}}^{\lambda, \omega} \in \mathcal{A}_{P_{sa}}^{\lambda, \omega}} P_{sa} [V]_{\alpha_{P_{sa}}^{\lambda, \omega}} - \tilde{\sigma} \|\pi_s\|_* \text{sp}([V]_{\alpha_{P_{sa}}^{\lambda, \omega}})_* \right). \quad (50)$$

736 with the definition of $\text{sp}(\cdot)_*$ in 8 and where the variational family $A_P^{\lambda,\omega}$ is defined as :

$$A_P^{\lambda,\omega} = \{\alpha \in [0, 1/(1-\gamma)]^S, \alpha = \omega + \lambda|\nabla \|P\| \mid := \alpha_P^{\lambda,\omega}\} \quad (51)$$

$$(52)$$

737 with ω is the generalized mean defined as the argmin in the definition of the span semi norm in 2 and
 738 λ, ω a positive scalar. Moreover, for L_1 or TV , case, the vector $\alpha_P^{\lambda,\omega}$ reduces to a 1 dimensional
 739 scalar such as $\alpha \in [0, 1/(1-\gamma)]$.

740 In the proof of the previous lemma, we decompose this problem s -rectangular radius $\tilde{\sigma}$ into sa -
 741 rectangular sub-problem with respectively radius σ_{sa} .

Proof.

$$\begin{aligned} \inf_{\mathcal{P}^\pi \in \mathcal{U}^{s,\tilde{\sigma}}(\mathcal{P}^\pi)} \mathcal{P}^\pi V &= \inf_{\{\sigma_{sa}: \|\sigma_{sa}\| \leq \tilde{\sigma}\}} \inf_{\mathcal{P}' \in \mathcal{U}^{sa,\sigma}(P_{sa})} \sum_a \pi(a|s) \mathcal{P}' V \\ &\stackrel{(a)}{=} \sum_a \pi(a|s) P_{sa} V + \min_{\{\sigma_{sa}: \|\sigma_{sa}\| \leq \tilde{\sigma}\}} \sum_a \pi(a|s) \min_{\{y: \|y\| \leq \sigma_{sa}, 1y=0, y \geq -P_{sa}\}} \sum_{s'} y(s') V \end{aligned}$$

742 where we use the change of variable $y(s') = \mathcal{P}_{sa}(s') - P_{sa}(s')$ in (a). Then we case use the previous
 743 lemma for sa rectangular assumption, Lemma 3. Then,

$$\begin{aligned} &\min_{\{\sigma_{sa}: \|\sigma_{sa}\| \leq \tilde{\sigma}\}} \sum_a \pi(a|s) \min_{\{y: \|y\| \leq \sigma_{sa}, 1y=0, y \geq -P_{sa}\}} \sum_{s'} y(s') V \\ &= \min_{\{\sigma_{sa}: \|\sigma_{sa}\| \leq \tilde{\sigma}\}} \sum_a \pi(a|s) \max_{\mu \geq 0} \left(-P_{sa}\mu - \sigma_{sa} \text{sp}(V - \mu)_* \right) \\ &= \left(\sum_a \pi(a|s) \max_{\mu \geq 0} \left\{ (-P_{sa}\mu) - \max_{\{\sigma_{sa}: \|\sigma_{sa}\| \leq \tilde{\sigma}\}} \sum_a \pi(a|s) \sigma_{sa} \text{sp}(V - \mu)_* \right\} \right) \\ &= \sum_a \pi(a|s) \max_{\mu \geq 0} \left\{ (-P_{sa}\mu) - \tilde{\sigma} \|\pi_s\|_* \text{sp}(V - \mu)_* \right\}. \end{aligned}$$

744 We can exchange the min and the max as we get concave-convex problems in σ and μ in the second
 745 line according to minimax theorem [v. Neumann, 1928] and using Cauchy Swartz inequality which is
 746 attained in the last equality. Finally, we obtain:

$$\begin{aligned} \inf_{\mathcal{P} \in \mathcal{U}^{s,\tilde{\sigma}}(P)} \mathcal{P}^\pi V &= \sum_a \pi(a|s) \left(\max_{\mu \geq 0} P_{sa}(V - \mu) - \tilde{\sigma} \|\pi_s\|_* \text{sp}(V - \mu)_* \right) \\ &\stackrel{(a)}{=} \sum_a \pi(a|s) \left(\max_{\alpha_P^{\lambda,\omega} \in A_{P_{sa}}^{\lambda,\omega}} P_{sa}[V]_{\alpha_P^{\lambda,\omega}} - \tilde{\sigma} \|\pi_s\|_* \text{sp}([V]_{\alpha_P^{\lambda,\omega}})_* \right) \end{aligned}$$

747 where in (a) we use the previous lemma for sa - rectangular case. Note that as we are using sa -
 748 rectangular case, for TV or L_1 , this lemma holds, but the vector α_P^λ reduces to a positive scalar
 749 denoted α which is equal to $\|v - \mu^*\|_\infty$. (See also Iyengar [2005]).

750

□

751 9 Proof of the upper bound : Theorem 1 and 3

752 9.1 Technical lemmas

753 We begin with a key lemma concerning the dynamic range of the robust value function $V^{\pi,\sigma}$ (cf. (7)),
 754 which produces tighter control when σ is large; the proof is deferred to Appendix 9.3.1. This lemma
 755 allows tighter control compared to Clavier et al. [2023].

756 **Lemma 5.** In *sa*-rectangular case (see (3), for any nominal transition kernel $P \in \mathbb{R}^{SA \times S}$, any
 757 fixed uncertainty level σ , and any policy π , its corresponding robust value function $V^{\pi, \sigma}$ (cf. (7))
 758 satisfies

$$\text{sp}(V^{\pi, \sigma})_{\infty} \leq \frac{1}{\gamma \max\{1 - \gamma, C_g \sigma\}} \quad (53)$$

759 where $C_g = 1/(\min_s \|e_s\|)$ is a geometric constant depending on the geometry of the norm. For
 760 example, for L_p , norms $p \geq 1$, $C_g \geq 1$ which reduce the sample complexity. In *s*-rectangular case,
 761 we obtain a slightly different lemma because of the dependency on π .

762 **Lemma 6.** The infinite span semi norm can be controlled as follows for every s in *s*-rectangular case
 763 (See (5)):
 764

$$\text{sp}(V^{\pi, \sigma})_{\infty} \leq \frac{1}{\gamma \max\{1 - \gamma, \|\pi_s\|_* C_g \tilde{\sigma}\}} \leq \frac{1}{\gamma \max\{1 - \gamma, \min_s \|\pi_s\|_* C_g \tilde{\sigma}\}} \quad (54)$$

765 where $C_g = \frac{1}{\min_s \|e_s\|}$ is a geometric constant depending on the geometry of the norm. These lemmas
 766 are required to get tight bounds for the sample complexity. The main difference between *sa*- and *s*-
 767 rectangular case is that we have an extra dependency on $\|\pi_s\|_*$, which represents how stochastic the
 768 policy can be in *s* rectangular MDPs.

769 **Lemma 7.** Consider an MDP with transition kernel matrix P and reward function $0 \leq r \leq 1$. For any
 770 policy π and its associated state transition matrix $P_{\pi} := \Pi^{\pi} P$ and value function $0 \leq V^{\pi, P} \leq \frac{1}{1-\gamma}$
 771 (cf. (1)), one has for *sa*- and *s*- rectangular assumptions.

$$(I - \gamma P_{\pi})^{-1} \sqrt{\text{Var}_{P_{\pi}}(V^{\pi, P})} \leq \sqrt{\frac{8}{\gamma^2(1-\gamma)^2} \text{sp}(V^{\pi, P})_{\infty}} 1.$$

772 See 9.3.7 for the proof

773 9.2 Proof of Theorem 1 and Theorem 3

774 The first decomposition of the proof of Theorem 1 and Theorem 3 Agarwal et al. [2020] while
 775 the argument needs essential adjustments in order to adapt to the robustness setting. One has by
 776 assumptions using any planner in empirical RMDPs :

$$\|\widehat{V}^{*, \sigma} - \widehat{V}^{\widehat{\pi}^*, \sigma}\|_{\infty} \leq \varepsilon_{\text{opt}}, \quad (55)$$

777 using previous inequality, performance gap $\|V^{*, \sigma} - V^{\widehat{\pi}^*, \sigma}\|_{\infty}$, can be upper bounded using 3 steps.

778 **First step: subdivide the performance gap in 3 terms.** We recall the definition of the optimal
 779 robust policy π^* with regard to \mathcal{M}_{rob} and the optimal robust policy $\widehat{\pi}^*$, the optimal robust value
 780 function $\widehat{V}^{*, \sigma}$ (resp. robust value function $\widehat{Q}^{\pi^*, \sigma}$) w.r.t. $\widehat{\mathcal{M}}_{\text{rob}}$. Then, the performance gap $V^{*, \sigma} - V^{\widehat{\pi}^*, \sigma}$
 781 can be decomposed in one optimization term and two statistical error terms

$$\begin{aligned} V^{*, \sigma} - V^{\widehat{\pi}^*, \sigma} &= (V^{\pi^*, \sigma} - \widehat{V}^{\pi^*, \sigma}) + (\widehat{V}^{\pi^*, \sigma} - \widehat{V}^{\widehat{\pi}^*, \sigma}) + (\widehat{V}^{\widehat{\pi}^*, \sigma} - \widehat{V}^{\widehat{\pi}^*, \sigma}) + (\widehat{V}^{\widehat{\pi}^*, \sigma} - V^{\widehat{\pi}^*, \sigma}) \\ &\stackrel{(i)}{\leq} (V^{\pi^*, \sigma} - \widehat{V}^{\pi^*, \sigma}) + (\widehat{V}^{\widehat{\pi}^*, \sigma} - \widehat{V}^{\widehat{\pi}^*, \sigma}) + (\widehat{V}^{\widehat{\pi}^*, \sigma} - V^{\widehat{\pi}^*, \sigma}) \\ &\stackrel{(ii)}{\leq} (V^{\pi^*, \sigma} - \widehat{V}^{\pi^*, \sigma}) + \varepsilon_{\text{opt}} + (\widehat{V}^{\widehat{\pi}^*, \sigma} - V^{\widehat{\pi}^*, \sigma}) \end{aligned} \quad (56)$$

782 where (i) holds by $\widehat{V}^{\pi^*, \sigma} - \widehat{V}^{\widehat{\pi}^*, \sigma} \leq 0$ since $\widehat{\pi}^*$ is the robust optimal policy for $\widehat{\mathcal{M}}_{\text{rob}}$, and (ii) comes
 783 from (55) and definition of optimization error. The proof aims to control the last remaining terms in
 784 (56) using concentration theory and sufficiently big number of step N . To do so, we will consider a
 785 more general term $\widehat{V}^{\pi, \sigma} - V^{\pi, \sigma}$ for any policy π even if control of these two terms slightly differ at

786 the end. Using (32), it holds that for both *sa*- and *s*-rectangular assumptions:

$$\begin{aligned}
\widehat{V}^{\pi,\sigma} - V^{\pi,\sigma} &= r_\pi + \gamma \widehat{\underline{P}}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - (r_\pi + \gamma \underline{P}^{\pi,V} V^{\pi,\sigma}) \\
&= \left(\gamma \widehat{\underline{P}}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - \gamma \underline{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} \right) + \left(\gamma \underline{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - \gamma \underline{P}^{\pi,V} V^{\pi,\sigma} \right) \\
&\stackrel{(i)}{\leq} \gamma \left(\underline{P}^{\pi,V} \widehat{V}^{\pi,\sigma} - \underline{P}^{\pi,V} V^{\pi,\sigma} \right) + \left(\gamma \widehat{\underline{P}}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - \gamma \underline{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} \right),
\end{aligned}$$

787 where (i) holds because $\underline{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} \leq \underline{P}^{\pi,V} \widehat{V}^{\pi,\sigma}$ because of the optimality of $\underline{P}^{\pi,\widehat{V}}$ (see. (25)).
788 Factorizing terms leads to the following equation

$$\widehat{V}^{\pi,\sigma} - V^{\pi,\sigma} \leq \gamma (I - \gamma \underline{P}^{\pi,V})^{-1} \left(\widehat{\underline{P}}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - \underline{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} \right). \quad (57)$$

789 In the same manner, we can also obtain a lower bound of this quantity:

$$\begin{aligned}
\widehat{V}^{\pi,\sigma} - V^{\pi,\sigma} &= r_\pi + \gamma \widehat{\underline{P}}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - (r_\pi + \gamma \underline{P}^{\pi,V} V^{\pi,\sigma}) \\
&= \left(\gamma \widehat{\underline{P}}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - \gamma \underline{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} \right) + \left(\gamma \underline{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - \gamma \underline{P}^{\pi,V} V^{\pi,\sigma} \right) \\
&\geq \gamma \left(\underline{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - \underline{P}^{\pi,\widehat{V}} V^{\pi,\sigma} \right) + \left(\gamma \widehat{\underline{P}}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - \gamma \underline{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} \right) \\
&\geq \gamma (I - \gamma \underline{P}^{\pi,\widehat{V}})^{-1} \left(\widehat{\underline{P}}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - \underline{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} \right). \quad (58)
\end{aligned}$$

790 Using both (57) and (58), we obtain infinite norm control:

$$\begin{aligned}
\|\widehat{V}^{\pi,\sigma} - V^{\pi,\sigma}\|_\infty &\leq \gamma \max \left\{ \left\| (I - \gamma \underline{P}^{\pi,V})^{-1} \left(\widehat{\underline{P}}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - \underline{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} \right) \right\|_\infty, \right. \\
&\quad \left. \left\| (I - \gamma \underline{P}^{\pi,\widehat{V}})^{-1} \left(\widehat{\underline{P}}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - \underline{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} \right) \right\|_\infty \right\}. \quad (59)
\end{aligned}$$

791 By decomposing the error in a symmetric way, he have

$$\begin{aligned}
\|\widehat{V}^{\pi,\sigma} - V^{\pi,\sigma}\|_\infty &\leq \gamma \max \left\{ \left\| (I - \gamma \widehat{\underline{P}}^{\pi,V})^{-1} \left(\widehat{\underline{P}}^{\pi,V} V^{\pi,\sigma} - \underline{P}^{\pi,V} V^{\pi,\sigma} \right) \right\|_\infty, \right. \\
&\quad \left. \left\| (I - \gamma \widehat{\underline{P}}^{\pi,\widehat{V}})^{-1} \left(\widehat{\underline{P}}^{\pi,V} V^{\pi,\sigma} - \underline{P}^{\pi,V} V^{\pi,\sigma} \right) \right\|_\infty \right\}. \quad (60)
\end{aligned}$$

792 Armed with these inequalities, we can use concentration inequalities to upper bound the two remaining
793 terms $\|\widehat{V}^{\pi^*,\sigma} - V^{\pi^*,\sigma}\|_\infty$ and $\|\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma}\|_\infty$ in (56). Taking $\pi = \widehat{\pi}$, applying (59) leads to

$$\begin{aligned}
\|\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma}\|_\infty &\leq \gamma \max \left\{ \left\| (I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}})^{-1} \left(\widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right) \right\|_\infty, \right. \\
&\quad \left. \left\| (I - \gamma \underline{P}^{\widehat{\pi},V})^{-1} \left(\widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right) \right\|_\infty \right\}. \quad (61)
\end{aligned}$$

794 Finally, $\pi = \pi^*$, applying (60) gives us

$$\begin{aligned}
\|\widehat{V}^{\pi^*,\sigma} - V^{\pi^*,\sigma}\|_\infty &\leq \gamma \max \left\{ \left\| (I - \gamma \widehat{\underline{P}}^{\pi^*,V})^{-1} \left(\widehat{\underline{P}}^{\pi^*,V} V^{\pi^*,\sigma} - \underline{P}^{\pi^*,V} V^{\pi^*,\sigma} \right) \right\|_\infty, \right. \\
&\quad \left. \left\| (I - \gamma \widehat{\underline{P}}^{\pi^*,\widehat{V}})^{-1} \left(\widehat{\underline{P}}^{\pi^*,V} V^{\pi^*,\sigma} - \underline{P}^{\pi^*,V} V^{\pi^*,\sigma} \right) \right\|_\infty \right\}. \quad (62)
\end{aligned}$$

795 Note that to control $\|\widehat{V}^{\pi^*,\sigma} - V^{\pi^*,\sigma}\|_\infty$, we use decomposition not depending on $\widehat{\pi}$ for value
796 function as $V^{\pi^*,\sigma}$ is deterministic and fixed, allowing use of classical concentration analysis tools.
797 This decomposition is the same for both *sa*-rectangular and *s*-rectangular case.

798 **Second step: bound first term and second term in (62) to control** $\|\widehat{V}^{\pi^*,\sigma} - V^{\pi^*,\sigma}\|_\infty$ To control
 799 the two terms in (62), we use lemma 8 based Bernstein's concentration argument and whose proof is
 800 in Appendix 9.3.3.

801 **Lemma 8.** For both sa - and s -rectangular setting, consider any $\delta \in (0, 1)$, with probability $1 - \delta$,
 802 it holds:

$$\left| \widehat{\underline{P}}^{\pi^*,V} V^{\pi^*,\sigma} - \underline{P}^{\pi^*,V} V^{\pi^*,\sigma} \right| \leq 2\sqrt{\frac{L}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{*,\sigma})} + \frac{3LC_S \|1\|_*}{N(1-\gamma)} \quad (63)$$

803 with $L = 2 \log(18 \|1\|_* SAN/\delta)$ and where $\text{Var}_{P^{\pi^*}}(V^{*,\sigma})$ is defined in (27). Moreover, for the
 804 specific case of TV, this lemma is true without the smoothness term $\frac{3LC_S \|1\|_*}{N(1-\gamma)}$.

805 Armed with the above lemma, now we control the **first term** on the right-hand side of (62) as follows:

$$\begin{aligned} & \left(I - \gamma \widehat{\underline{P}}^{\pi^*,V} \right)^{-1} \left(\widehat{\underline{P}}^{\pi^*,V} V^{\pi^*,\sigma} - \underline{P}^{\pi^*,V} V^{\pi^*,\sigma} \right) \\ & \stackrel{(a)}{\leq} \left(I - \gamma \widehat{\underline{P}}^{\pi^*,V} \right)^{-1} \left\| \widehat{\underline{P}}^{\pi^*,V} V^{\pi^*,\sigma} - \underline{P}^{\pi^*,V} V^{\pi^*,\sigma} \right\|_\infty \\ & \stackrel{(b)}{\leq} \left(I - \gamma \widehat{\underline{P}}^{\pi^*,V} \right)^{-1} \left(2\sqrt{\frac{L}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{*,\sigma})} + \frac{3LC_S \|1\|_*}{N(1-\gamma)} \right) \\ & \leq \underbrace{\left(I - \gamma \widehat{\underline{P}}^{\pi^*,V} \right)^{-1} \frac{3LC_S \|1\|_*}{N(1-\gamma)} 1}_{=: \mathcal{R}_1} + 2\sqrt{\frac{L}{N}} \underbrace{\left(I - \gamma \widehat{\underline{P}}^{\pi^*,V} \right)^{-1} \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*,V}}(V^{*,\sigma})}}_{=: \mathcal{R}_2} \\ & \quad + 2\sqrt{\frac{L}{N}} \underbrace{\left(I - \gamma \widehat{\underline{P}}^{\pi^*,V} \right)^{-1} \sqrt{\left| \text{Var}_{\widehat{\underline{P}}^{\pi^*}}(V^{*,\sigma}) - \text{Var}_{\widehat{\underline{P}}^{\pi^*,V}}(V^{*,\sigma}) \right|}}_{=: \mathcal{R}_2} \\ & \quad + 2\sqrt{\frac{L}{N}} \underbrace{\left(I - \gamma \widehat{\underline{P}}^{\pi^*,V} \right)^{-1} \left(\sqrt{\text{Var}_{P^{\pi^*}}(V^{*,\sigma})} - \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*}}(V^{*,\sigma})} \right)}_{=: \mathcal{R}_3}, \end{aligned} \quad (64)$$

806 where (a) holds as the matrix $\left(I - \gamma \widehat{\underline{P}}^{\pi^*,V} \right)^{-1}$ is positive definite, (b) holds due to Lemma 8, and
 807 the last point holds from the following decomposition for variance and triangular inequality

$$\begin{aligned} \sqrt{\text{Var}_{P^{\pi^*}}(V^{*,\sigma})} &= \left(\sqrt{\text{Var}_{P^{\pi^*}}(V^{*,\sigma})} - \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*}}(V^{*,\sigma})} \right) + \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*}}(V^{*,\sigma})} \\ &\leq \left(\sqrt{\text{Var}_{P^{\pi^*}}(V^{*,\sigma})} - \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*}}(V^{*,\sigma})} \right) \\ &\quad + \sqrt{\left| \text{Var}_{\widehat{\underline{P}}^{\pi^*}}(V^{*,\sigma}) - \text{Var}_{\widehat{\underline{P}}^{\pi^*,V}}(V^{*,\sigma}) \right|} + \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*,V}}(V^{*,\sigma})}. \end{aligned}$$

808 Finally, the fact that $\widehat{\underline{P}}^{\pi^*,V}$ is a stochastic matrix, so

$$\left(I - \gamma \widehat{\underline{P}}^{\pi^*,V} \right)^{-1} 1 = \left(I + \sum_{t=1}^{\infty} \gamma^t \left(\widehat{\underline{P}}^{\pi^*,V} \right)^t \right) 1 \leq \frac{1}{1-\gamma} 1. \quad (65)$$

809 Armed with these inequalities, the three terms $\mathcal{R}_1, \mathcal{R}_2, \mathcal{R}_3$ in (64) can be controlled separately.

810 • Consider \mathcal{R}_1 . We first introduce the following lemma, whose proof is postponed to Ap-
 811 pendix 9.3.4.

812

Lemma 9. Consider any $\delta \in (0, 1)$. With probability at least $1 - \delta$, one has

$$\begin{aligned} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, V}}(V^{*, \sigma})} &\leq 4 \sqrt{\frac{\left(1 + \left(\sqrt{\frac{L}{(1-\gamma)^2 N}} + \frac{C_S \|1\|_* L}{N(1-\gamma)} \right) \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, C_g \sigma\}}} 1 \\ &\leq 4 \sqrt{\frac{\left(1 + \left(\sqrt{\frac{L}{(1-\gamma)^2 N}} + \frac{C_S \|1\|_* L}{N(1-\gamma)} \right) \right)}{\gamma^3 (1-\gamma)^3}} 1 \end{aligned}$$

813

with $L = 2 \log\left(\frac{18 \|1\|_* S A N}{\delta}\right)$ in the *sa*-rectangular case. In the *s*-rectangular case, it holds:

$$\begin{aligned} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, V}}(V^{*, \sigma})} &\leq 4 \sqrt{\frac{\left(1 + \left(\sqrt{\frac{L}{(1-\gamma)^2 N}} + \frac{C_S \|1\|_* L}{N(1-\gamma)} \right) \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, C_g \tilde{\sigma} \min_s \|\pi_s\|_*\}}} 1 \\ &\leq 4 \sqrt{\frac{\left(1 + \left(\sqrt{\frac{L}{(1-\gamma)^2 N}} + \frac{C_S \|1\|_* L}{N(1-\gamma)} \right) \right)}{\gamma^3 (1-\gamma)^3}} 1 \end{aligned}$$

814

Using Lemma 9 and inserting back to (64) gives in *sa*-rectangular case

$$\begin{aligned} \mathcal{R}_1 &= 2 \sqrt{\frac{L}{N}} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, V}}(V^{*, \sigma})} \\ &\leq 8 \sqrt{\frac{L}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N} \left(1 + \sqrt{\frac{L}{(1-\gamma)^2 N}} + \frac{C_S \|1\|_* L}{N(1-\gamma)} \right)} 1. \end{aligned} \quad (66)$$

815

• Consider \mathcal{R}_2 . First, denote $V' := V^{*, \sigma} - \eta 1$ $\eta \in \mathbb{R}$, by Lemma 5, we have for any π ,

$$0 \leq \min_{\eta} \|V\|_{\infty} - \eta 1 \leq \frac{1}{\gamma \max\{1-\gamma, C_g \sigma\}}. \quad (67)$$

816

for *sa*-rectangular case or in *s*-rectangular we obtain

$$0 \leq \min_{\eta} \|V - \eta 1\|_{\infty} \leq \frac{1}{\gamma \max\{1-\gamma, \tilde{\sigma} C_g \|\pi_s\|_*\}} \quad (68)$$

817

by the definition of the span semi norm. Moreover, we can use Holder with L_1 and L_{∞} we have for both *sa* and *s*-rectangular case to as it holds that:

818

$$\begin{aligned} \left| \text{Var}_{\tilde{P}_{s,a}}(V^{*, \sigma}) - \text{Var}_{P_{s,a}}(V^{*, \sigma}) \right| &= \left| \text{Var}_{\tilde{P}_{s,a}}(V') - \text{Var}_{P_{s,a}}(V') \right| \\ &\leq \|\tilde{P}_{s,a} - P_{s,a}\|_1 \|V'\|_{\infty}^2 \stackrel{a}{\leq} \frac{\sigma_1}{(\gamma^2 (\max(1-\gamma), C_g \sigma))^2} \\ &\leq \frac{1}{\gamma^2 \max\{(1-\gamma), \sigma C_g\}} \end{aligned} \quad (69)$$

819

In the first inequality, we use $\|V'\|_{\infty}^2 = \|V'^2\|_{\infty}$ and we use Lemma 5 in (a) where $C_g \sigma = \sigma_1$.

820

821

With the same arguments for *s*-rectangular, we obtain for $V' := V^{*, \sigma} - \eta 1$ $\eta \in \mathbb{R}$,

$$\begin{aligned} \left| \Pi^{\pi^*} \left(\text{Var}_{\tilde{P}_s}(V^{*, \sigma}) - \text{Var}_{P_s}(V^{*, \sigma}) \right) \right| &= \left| \Pi^{\pi^*} \left(\text{Var}_{\tilde{P}_s}(V') - \text{Var}_{P_s}(V') \right) \right| \\ &\leq \sum_a \pi(a|s) (\tilde{P}_s(s', a) - P_s(s', a)) V(s')^2 \end{aligned} \quad (70)$$

$$\stackrel{a}{\leq} \|V'\|_{\infty}^2 \sum_a \pi(a|s) (\tilde{P}_s(s', a) - P_s(s', a)) \stackrel{b}{\leq} \|V'\|_{\infty}^2 \tilde{\sigma} \|\pi_s\|_* \quad (71)$$

$$\stackrel{c}{\leq} \frac{\tilde{\sigma} C_g \|\pi_s^*\|_* \|V'\|_{\infty}}{\gamma \|\pi_s^*\|_* \tilde{\sigma} C_g} 1 \leq \frac{\|V'\|_{\infty}}{\gamma} 1. \quad (72)$$

822
823

where where (a) and (b) comes Cauchy Swartz inequality, , (c) comes lemma 6. Then, taking the sup over s in the previous equations, it holds

$$|\Pi^{\pi^*} (\text{Var}_{\hat{P}_s}(V^{*,\sigma}) - \text{Var}_{P_s}(V^{*,\sigma}))| \leq \frac{\inf_{\eta \in \mathbb{R}^+} \|V - \eta \mathbf{1}'\|_1}{\gamma} \quad (73)$$

$$\leq \frac{1}{\gamma^2 \tilde{\sigma} \min_s \|\pi_s^*\|_* C_g} \mathbf{1}. \quad (74)$$

824

Applying the previous inequality, it holds in sa -rectangular case:

$$\begin{aligned} \mathcal{R}_2 &= 2\sqrt{\frac{L}{N}} \left(I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \sqrt{|\text{Var}_{\hat{P}^{\pi^*}}(V^{*,\sigma}) - \text{Var}_{\hat{P}^{\pi^*, V}}(V^{*,\sigma})|} \\ &= 2\sqrt{\frac{L}{N}} \left(I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \sqrt{|\Pi^{\pi^*} (\text{Var}_{\hat{P}_0}(V^{*,\sigma}) - \text{Var}_{\hat{P}^{\pi^*, V}}(V^{*,\sigma}))|} \\ &\leq 2\sqrt{\frac{L}{N}} \left(I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \sqrt{\|\text{Var}_{\hat{P}_0}(V^{*,\sigma}) - \text{Var}_{\hat{P}^{\pi^*, V}}(V^{*,\sigma})\|_{\infty} \mathbf{1}} \\ &\leq 2\sqrt{\frac{L}{N}} \left(I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \sqrt{\frac{1}{\gamma^2 \max\{1 - \gamma, C_g \sigma\}}} \mathbf{1} \end{aligned} \quad (75)$$

$$\leq 4\sqrt{\frac{L}{\gamma^2 (1 - \gamma)^2 \max\{1 - \gamma, C_g \sigma\} N}} \mathbf{1}, \quad (76)$$

825

where the last inequality uses $\left(I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \mathbf{1} \leq \frac{1}{1 - \gamma} \mathbf{1}$ (cf. (65)). for sa -rectangular

826

In the s -rectangular case, we obtain a different result as

$$\begin{aligned} \mathcal{R}_2 &= 2\sqrt{\frac{L}{N}} \left(I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \sqrt{|\text{Var}_{\hat{P}^{\pi^*}}(V^{*,\sigma}) - \text{Var}_{\hat{P}^{\pi^*, V}}(V^{*,\sigma})|} \\ &= 2\sqrt{\frac{L}{N}} \left(I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \sqrt{|\Pi^{\pi^*} (\text{Var}_{\hat{P}_0}(V^{*,\sigma}) - \text{Var}_{\hat{P}^{\pi^*, V}}(V^{*,\sigma}))|} \\ &\leq 2\sqrt{\frac{L}{N}} \left(I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \sqrt{\frac{1}{\gamma^2 \max\{1 - \gamma, \min_s \|\pi_s^*\|_{\infty} C_g \tilde{\sigma}\}}} \mathbf{1} \end{aligned} \quad (77)$$

$$\leq 2\sqrt{\frac{L}{\gamma^2 (1 - \gamma)^2 \max\{1 - \gamma, \min_s \|\pi_s^*\|_{\infty} \tilde{\sigma} C_g\} N}} \mathbf{1}, \quad (78)$$

827

• Consider \mathcal{R}_3 . The following lemma plays an important role.

828

Applying Lemma 2 and using $\pi = \pi^*$ and $V = V^{*,\sigma}$, it holds

$$\sqrt{\text{Var}_{P^{\pi^*}}(V^{*,\sigma})} - \sqrt{\text{Var}_{\hat{P}^{\pi^*}}(V^{*,\sigma})} \leq \sqrt{\frac{2\|V^{*,\sigma}\|_{\infty}^2 \log(\frac{2SA}{\delta})}{N}} \mathbf{1},$$

829

which can be inserted in (64) to gives

$$\begin{aligned} \mathcal{R}_3 &= 2\sqrt{\frac{L}{N}} \left(I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \left(\sqrt{\text{Var}_{P^{\pi^*}}(V^{*,\sigma})} - \sqrt{\text{Var}_{\hat{P}^{\pi^*}}(V^{*,\sigma})} \right) \\ &\leq \frac{4}{(1 - \gamma)} \frac{\log(\frac{SA}{\delta}) \|V^{*,\sigma}\|_{\infty} \mathbf{1}}{N} \leq \frac{4L}{(1 - \gamma)^2 N} \mathbf{1}, \end{aligned} \quad (79)$$

830

where the last line uses $\left(I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \mathbf{1} \leq \frac{1}{1 - \gamma} \mathbf{1}$ (cf. (65)).

831 Finally, inserting the results of \mathcal{R}_1 in (66), \mathcal{R}_2 in (78), \mathcal{R}_3 in (79), and (65) back into (64) gives

$$\begin{aligned}
& \left(I - \gamma \widehat{\underline{P}}^{\pi^*, V}\right)^{-1} \left(\widehat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma}\right) \tag{80} \\
& \leq 8 \sqrt{\frac{L}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N} \left(1 + \sqrt{\frac{L}{(1-\gamma)^2 N} + \frac{C_S \|1\|_* L}{N(1-\gamma)}}\right) 1 + \frac{3LC_S \|1\|_*}{N(1-\gamma)^2} 1} \\
& \quad + 2 \sqrt{\frac{2L}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N} 1 + \frac{4L}{(1-\gamma)^2 N} 1} \\
& \leq 10 \sqrt{\frac{2L}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N} \left(1 + \sqrt{\frac{L}{(1-\gamma)^2 N} + \frac{C_S \|1\|_* L}{N(1-\gamma)}}\right) 1 + \frac{4L}{(1-\gamma)^2 N} 1 + \frac{3LC_S \|1\|_*}{N(1-\gamma)^2} 1} \\
& \leq 160 \sqrt{\frac{L(1 + \frac{C_S \|1\|_*}{N(1-\gamma)})}{(1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N} 1 + \frac{7LC_S \|1\|_*}{N(1-\gamma)^2} 1}, \tag{81}
\end{aligned}$$

832 where the last inequality holds by the fact $\gamma \geq \frac{1}{4}$ and letting $N \geq \frac{L}{(1-\gamma)^2}$. We have the same result
833 for s -rectangular, replacing, $\max\{1-\gamma, C_g \sigma\}$ by $\max\{1-\gamma, \min_s \|\pi_s^*\|_* \bar{\sigma} C_g\}$.

834 Now we are ready to control **second term in** (62) to control $\|\widehat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma}\|_\infty$. To proceed,
835 applying Lemma 8 on the second term of the right-hand side of (62) leads to

$$\begin{aligned}
& \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \left(\widehat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma}\right) \\
& \leq \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \left(2 \sqrt{\frac{L}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{\pi^*, \sigma})} + \frac{3LC_S \|1\|_*}{N(1-\gamma)}\right) \\
& \leq \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \frac{L' C_S \|1\|_*}{N(1-\gamma)} + 2 \underbrace{\sqrt{\frac{L}{N}} \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*, \widehat{V}}}(\widehat{V}^{\pi^*, \sigma})}}_{=: \mathcal{R}_4} \\
& \quad 2 \underbrace{\sqrt{\frac{L}{N}} \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \left(\sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*, \widehat{V}}}(V^{\pi^*, \sigma} - \widehat{V}^{\pi^*, \sigma})}\right)}_{=: \mathcal{R}_5} \\
& \quad + 2 \underbrace{\sqrt{\frac{L}{N}} \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \left(\sqrt{\left|\text{Var}_{\widehat{\underline{P}}^{\pi^*}}(V^{\pi^*, \sigma}) - \text{Var}_{\widehat{\underline{P}}^{\pi^*, \widehat{V}}}(V^{\pi^*, \sigma})\right|}\right)}_{=: \mathcal{R}_6} \\
& \quad + 2 \underbrace{\sqrt{\frac{L}{N}} \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \left(\sqrt{\text{Var}_{P^{\pi^*}}(V^{\pi^*, \sigma})} - \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*}}(V^{\pi^*, \sigma})}\right)}_{=: \mathcal{R}_7}. \tag{82}
\end{aligned}$$

836 We now bound the above four terms $\mathcal{R}_4, \mathcal{R}_5, \mathcal{R}_6, \mathcal{R}_7$ separately.

837 • Using Lemma 7 with $P = \widehat{\underline{P}}^{\pi^*, \widehat{V}}$, $\pi = \pi^*$ and $V = \widehat{V}^{\pi^*, \sigma}$ which follow $\widehat{V}^{\pi^*, \sigma} =$
838 $r_{\pi^*} + \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}} \widehat{V}^{\pi^*, \sigma}$, and in view of (65), the term \mathcal{R}_4 in (82) can be controlled as follows:

$$\begin{aligned}
\mathcal{R}_4 & = 2 \sqrt{\frac{L}{N}} \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*, \widehat{V}}}(\widehat{V}^{\pi^*, \sigma})} \\
& \leq 2 \sqrt{\frac{L}{N}} \sqrt{\frac{8 \min\{\text{sp}(\widehat{V}^{\pi^*, \sigma})_*, 1/(1-\gamma)\}}{\gamma^2(1-\gamma)^2}} 1 \\
& \leq 8 \sqrt{\frac{L}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N}} 1, \tag{83}
\end{aligned}$$

839 where the last inequality is due to Lemma 5 for *sa*-rectangular case and with the same
 840 quantity replacing $\max\{1 - \gamma, \sigma\}$ by $\max\{1 - \gamma, \min_s \|\pi_s^*\|_* \tilde{\sigma}\}$ in the *s*-rectangular
 841 case.

842 • For bounding \mathcal{R}_5 , we can simply use (65) to get

$$\begin{aligned} \mathcal{R}_5 &= 2\sqrt{\frac{L}{N}} \left(I - \gamma \hat{P}^{\pi^*, \hat{V}} \right)^{-1} \sqrt{\text{Var}_{\hat{P}^{\pi^*, \hat{V}}} (V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma})} \\ &\leq 2\sqrt{\frac{L}{(1-\gamma)^2 N}} \|V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma}\|_\infty. \end{aligned} \quad (84)$$

843 moreover,

$$\|V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma}\|_\infty \leq \|V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma}\|_\infty \leq \|V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma}\|_\infty \quad (85)$$

844 as for $a > 0, b > 0$, we have $[a] - [b] < [a - b]$. Finally, we obtain

$$\mathcal{R}_5 \leq 2\sqrt{\frac{L}{(1-\gamma)^2 N}} \|V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma}\|_\infty. \quad (86)$$

845 • The term \mathcal{R}_6 can upper bounded as (78) as follows:

$$\mathcal{R}_6 \leq 2\sqrt{\frac{2L}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, C_g\sigma\} N}}. \quad (87)$$

846 for *sa*-rectangular case and with the same quantity replacing $\max\{1 - \gamma, C_g\sigma\}$ by $\max\{1 -$
 847 $\gamma, \min_s \|\pi_s^*\|_* \tilde{\sigma} C_g\}$ in the *s*-rectangular case.

848 • Finally, \mathcal{R}_7 can be controlled the same as (79) shown below:

$$\mathcal{R}_7 \leq \frac{4L}{(1-\gamma)^2 N}. \quad (88)$$

849 Combining the results in (83), (86), (87), and (88) and inserting back to (82) leads to for $N \geq \frac{L}{(1-\gamma)^2}$

$$\begin{aligned} &\left(I - \gamma \hat{P}^{\pi^*, \hat{V}} \right)^{-1} \left(\hat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \leq 8\sqrt{\frac{L(1 + \frac{C_S \|1\|_*}{N(1-\gamma)})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, C_g\sigma\} N}} \\ &+ 2\sqrt{\frac{L}{(1-\gamma)^2 N}} \|V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma}\|_\infty + 2\sqrt{\frac{2L}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, C_g\sigma\} N}} + \frac{7LC_S \|1\|_*}{N(1-\gamma)^2} \\ &\leq 80\sqrt{\frac{L(1 + \frac{C_S \|1\|_*}{N(1-\gamma)})}{(1-\gamma)^2 \max\{1-\gamma, C_g\sigma\} N}} + 2\sqrt{\frac{L}{(1-\gamma)^2 N}} \|V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma}\|_\infty + \frac{7LC_S \|1\|_*}{N(1-\gamma)^2}, \end{aligned} \quad (89)$$

850 where the last inequality follows from the assumption $\gamma \geq \frac{1}{4}$. Finally, inserting (81) and (89) back to
 851 (62) yields

$$\begin{aligned} &\|\hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma}\|_\infty \leq \max \left\{ 160\sqrt{\frac{L(1 + \frac{C_S \|1\|_*}{N(1-\gamma)})}{(1-\gamma)^2 \max\{1-\gamma, C_g\sigma\} N}} + \frac{7LC_S \|1\|_*}{N(1-\gamma)^2}, \right. \\ &80\sqrt{\frac{L(1 + \frac{C_S \|1\|_*}{N(1-\gamma)})}{(1-\gamma)^2 \max\{1-\gamma, C_g\sigma\} N}} + 2\sqrt{\frac{L}{(1-\gamma)^2 N}} \|V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma}\|_\infty + \left. \frac{7LC_S \|1\|_*}{N(1-\gamma)^2} \right\} \\ &\leq 160\sqrt{\frac{L(1 + \frac{C_S \|1\|_*}{N(1-\gamma)})}{(1-\gamma)^2 \max\{1-\gamma, C_g\sigma\} N}} + \frac{14LC_S \|1\|_*}{N(1-\gamma)^2}, \end{aligned} \quad (90)$$

852 where the last inequality holds by taking $N \geq \frac{16 \log(\frac{SAN}{(1-\gamma)^2})}{(1-\gamma)^2}$ rearranging terms. In *s*-rectangular case,
 853 we obtain the same result, replacing $\max\{1 - \gamma, C_g\sigma\}$ by $\max\{1 - \gamma, \min_s \|\pi_s^*\|_* C_g \tilde{\sigma}\}$.

854 **Third step: controlling $\|\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma}\|_\infty$ or bounding the first and second term in (61).** Unlike
 855 the earlier term, one has to face a more complicated statistical dependency between $\widehat{\pi}$ and the
 856 empirical RMDP. To begin with, we introduce the following lemma which controls the main term on the
 857 right-hand side of (61), which is proved in Appendix 9.3.5.

858 **Lemma 10.** Consider any $\delta \in (0, 1)$. Taking $N \geq L''$ with probability at least $1 - \delta$, one has for sa-
 859 or s-rectangular case :

$$\begin{aligned} \left| \widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right| &\leq 2\sqrt{\frac{L'}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma})} \mathbf{1} + 2\varepsilon_{\text{opt}} \mathbf{1} + \frac{15L''C_S \|\mathbf{1}\|_*}{N(1-\gamma)} \\ &\leq 2\sqrt{\frac{L''}{(1-\gamma)^2 N}} \mathbf{1} + 2\varepsilon_{\text{opt}} \mathbf{1} + \frac{14L''C_S \|\mathbf{1}\|_*}{N(1-\gamma)} \mathbf{1}. \end{aligned} \quad (91)$$

860 with $L'' = 2 \log\left(\frac{54\|\mathbf{1}\|_* S A N^2}{(1-\gamma)\delta}\right)$. Moreover, For TV this lemma holds but without the geometric term
 861 $\frac{14L''C_S \|\mathbf{1}\|_*}{N(1-\gamma)} \mathbf{1}$. Taking the sup over s gives the final result.

862 With Lemma 10 in hand, we have to control **first term** in (61)

$$\begin{aligned} &\left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \left(\widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right) \\ &\stackrel{(i)}{\leq} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \left| \widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right| \\ &\leq 2\sqrt{\frac{L'}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\text{Var}_{P^{\widehat{\pi}}}(\widehat{V}^{*,\sigma})} + \left(I - \gamma \underline{P}^{\widehat{\pi},V^{\widehat{\pi}}} \right)^{-1} \left(2\varepsilon_{\text{opt}} \right) \mathbf{1} \\ &+ \left(I - \gamma \underline{P}^{\widehat{\pi},V^{\widehat{\pi}}} \right)^{-1} \frac{14L''C_S \|\mathbf{1}\|_*}{N(1-\gamma)} \mathbf{1} \\ &\stackrel{(ii)}{\leq} \underbrace{\left(\frac{2\varepsilon_{\text{opt}}}{1-\gamma} \right) \mathbf{1} + 2\sqrt{\frac{L'}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})}}_{=:S_1} \\ &+ \underbrace{2\sqrt{\frac{L'}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\left| \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{*,\sigma}) - \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma}) \right|}}_{=:S_2} \\ &+ \underbrace{2\sqrt{\frac{L'}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\left| \text{Var}_{P^{\widehat{\pi}}}(\widehat{V}^{*,\sigma}) - \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{*,\sigma}) \right|}}_{=:S_3}, \end{aligned} \quad (92)$$

863 where (i) and (ii) hold by the fact that each row of $(1-\gamma) \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1}$ is a probability vector
 864 that falls into $\Delta(\mathcal{S})$. The remainder of the proof will focus on controlling the three terms in (93)
 865 separately.

866 • For S_1 , we introduce the following lemma, whose proof is postponed to 9.3.6.

867 **Lemma 11.** Consider any $\delta \in (0, 1)$. Taking $N \geq \frac{L''}{(1-\gamma)^2}$ one has with probability at least
 868 $1 - \delta$, for sa-rectangular

$$\begin{aligned} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})} &\leq 6\sqrt{\frac{\left(1 + \varepsilon_{\text{opt}} + \frac{L''C_S \|\mathbf{1}\|_*}{N(1-\gamma)} \right)}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}}} \mathbf{1} \\ &\leq 6\sqrt{\frac{\left(1 + \varepsilon_{\text{opt}} + \frac{L''C_S \|\mathbf{1}\|_*}{N(1-\gamma)} \right)}{(1-\gamma)^3 \gamma^3}} \mathbf{1}. \end{aligned}$$

869 and for s-rectangular

$$\begin{aligned}
(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}})^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} &\leq 6 \sqrt{\frac{L'' \left(1 + \varepsilon_{\text{opt}} + \frac{C_S \|1\|_*}{N(1-\gamma)}\right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, C_g \tilde{\sigma} \min_s \|\hat{\pi}_s\|_*\}}} 1 \\
&\leq 6 \sqrt{\frac{L'' \left(1 + \varepsilon_{\text{opt}} + \frac{C_S \|1\|_*}{N(1-\gamma)}\right)}{(1-\gamma)^3 \gamma^2}} 1.
\end{aligned}$$

870 Applying Lemma 11 and (65) to (93) leads to

$$\begin{aligned}
\mathcal{S}_1 &= 2 \sqrt{\frac{L'}{N}} (I - \gamma \underline{P}^{\hat{\pi}, \hat{V}})^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} \\
&\leq 12 \sqrt{\frac{L''}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N}} 1. \tag{94}
\end{aligned}$$

871 for sa -rectangular and the same quantity replacing $\max\{1-\gamma, C_g \sigma\}$ by $\max\{1-$
872 $\gamma, C_g \tilde{\sigma} \min_s \|\hat{\pi}_s\|_*\}$ for s -rectangular case.

873 • Applying Lemma 1 with $\|\hat{V}^{*, \sigma} - \hat{V}^{\hat{\pi}, \sigma}\|_\infty \leq \varepsilon_{\text{opt}}$ and (65), \mathcal{S}_2 can be controlled as

$$\begin{aligned}
\mathcal{S}_2 &= 2 \sqrt{\frac{L''}{N}} (I - \gamma \underline{P}^{\hat{\pi}, \hat{V}})^{-1} \sqrt{\left| \text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{*, \sigma}) - \text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma}) \right|} \\
&\leq 4 \sqrt{\frac{L''}{N}} (I - \gamma \underline{P}^{\hat{\pi}, \hat{V}})^{-1} \sqrt{\varepsilon_{\text{opt}} \frac{1}{1-\gamma}} \leq 8 \sqrt{\frac{\varepsilon_{\text{opt}} L''}{(1-\gamma)^4 N}} 1. \tag{95}
\end{aligned}$$

874 • \mathcal{S}_3 can be controlled similar to \mathcal{R}_2 in (78) as follows:

$$\begin{aligned}
\mathcal{S}_3 &= 2 \sqrt{\frac{L''}{N}} (I - \gamma \underline{P}^{\hat{\pi}, \hat{V}})^{-1} \sqrt{\left| \text{Var}_{\underline{P}^{\hat{\pi}}}(\hat{V}^{*, \sigma}) - \text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{*, \sigma}) \right|} \\
&\leq 4 \sqrt{\frac{L''}{N}} (I - \gamma \underline{P}^{\hat{\pi}, \hat{V}})^{-1} \sqrt{\frac{1}{\gamma^2 \max\{1-\gamma, C_g \sigma\}}} 1 \leq 8 \sqrt{\frac{L''}{\gamma^2 (1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N}} 1 \tag{96}
\end{aligned}$$

875 for sa -rectangular and replacing $\max\{1-\gamma, \sigma\}$ by $\max\{1-\gamma, \tilde{\sigma} \min_s \|\hat{\pi}_s\|_*\}$ for s -
876 rectangular case.

877 Finally, summing up the results in (94), (95), and (96) and inserting them back to (93) yields: taking
878 $N \geq \frac{L''}{(1-\gamma)^2}$, with probability at least $1 - \delta$,

$$\begin{aligned}
&(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}})^{-1} \left(\hat{\underline{P}}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right) \leq \left(\frac{2\varepsilon_{\text{opt}}}{1-\gamma} \right) 1 + \frac{14L'' C_S \|1\|_*}{N(1-\gamma)^2} 1 \\
&+ 12 \sqrt{\frac{L'' \left(1 + \varepsilon_{\text{opt}} + \frac{C_S \|1\|_*}{N(1-\gamma)}\right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N}} 1 + 8 \sqrt{\frac{\varepsilon_{\text{opt}} L'}{(1-\gamma)^4 N}} 1 + 8 \sqrt{\frac{L'}{\gamma^2 (1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N}} 1 \\
&\leq 16 \sqrt{\frac{L'' \left(1 + \varepsilon_{\text{opt}} + \frac{C_S \|1\|_*}{N(1-\gamma)}\right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1 + \left(\frac{2\varepsilon_{\text{opt}} \gamma}{(1-\gamma)} + 8 \sqrt{\frac{\varepsilon_{\text{opt}} \gamma L'}{(1-\gamma)^4 N}} 1 + \frac{15L'' C_S \|1\|_*}{N(1-\gamma)^2} 1 \right) \tag{97}
\end{aligned}$$

(98)

879 for sa -rectangular and the same quantity replacing $\max\{1-\gamma, \sigma\}$ by $\max\{1-\gamma, \tilde{\sigma} \min_s \|\hat{\pi}_s\|_*\}$
880 for s -rectangular case. In this step, it is harder to decouple terms as $\hat{V}^{\hat{\pi}}$ depends on data both in $\hat{\pi}$
881 and \hat{V} .

882 **Step 5: controlling $\|\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma}\|_\infty$: bounding the second term in (61).** Towards this, applying
 883 Lemma 10 leads to in *sa*-rectangular case:

$$\begin{aligned}
 & (I - \gamma \underline{P}^{\widehat{\pi},V})^{-1} \left(\underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right) \leq (I - \gamma \underline{P}^{\widehat{\pi},V})^{-1} \left| \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right| \\
 & \leq 2\sqrt{\frac{L''}{N}} (I - \gamma \underline{P}^{\widehat{\pi},V})^{-1} \sqrt{\text{Var}_{P^{\widehat{\pi}}}(\widehat{V}^{\star,\sigma})} + (I - \gamma \underline{P}^{\widehat{\pi},V})^{-1} \left(2\varepsilon_{\text{opt}} \right) \mathbf{1} \tag{99} \\
 & + (I - \gamma \underline{P}^{\widehat{\pi},V})^{-1} \frac{L'' 14 C_S \|1\|_* \mathbf{1}}{N(1-\gamma)} \\
 & \leq \left(\frac{2\varepsilon_{\text{opt}}}{(1-\gamma)} \right) \mathbf{1} + 2 \underbrace{\sqrt{\frac{L''}{N}} (I - \gamma \underline{P}^{\widehat{\pi},V})^{-1} \sqrt{\text{Var}_{P^{\widehat{\pi},V}}(V^{\widehat{\pi},\sigma})}}_{=: \mathcal{S}_4} + (I - \gamma \underline{P}^{\widehat{\pi},V})^{-1} \frac{L'' C_S \|1\|_* \mathbf{1}}{N(1-\gamma)} \\
 & + 2 \underbrace{\sqrt{\frac{L'}{N}} (I - \gamma \underline{P}^{\widehat{\pi},V})^{-1} \sqrt{\text{Var}_{P^{\widehat{\pi},V}}(\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma})}}_{=: \mathcal{S}_5} \\
 & + 2 \underbrace{\sqrt{\frac{L''}{N}} (I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}})^{-1} \sqrt{\left| \text{Var}_{P^{\widehat{\pi},V}}(\widehat{V}^{\star,\sigma}) - \text{Var}_{P^{\widehat{\pi},V}}([\widehat{V}^{\widehat{\pi},\sigma}] \right|}}_{=: \mathcal{S}_6} \\
 & + 2 \underbrace{\sqrt{\frac{L''}{N}} (I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}})^{-1} \sqrt{\left| \text{Var}_{P^{\widehat{\pi}}}(\widehat{V}^{\star,\sigma}) - \text{Var}_{P^{\widehat{\pi},V}}([\widehat{V}^{\star,\sigma}]) \right|}}_{=: \mathcal{S}_7}. \tag{100}
 \end{aligned}$$

884 We shall bound each of the terms separately.

885 • Applying Lemma 7 with $P = \underline{P}^{\widehat{\pi},V}$, $\pi = \widehat{\pi}$, and taking $V = V^{\widehat{\pi},\sigma}$ which obeys $V^{\widehat{\pi},\sigma} =$
 886 $r_{\widehat{\pi}} + \gamma \underline{P}^{\widehat{\pi},V} V^{\widehat{\pi},\sigma}$, the term \mathcal{S}_4 can be controlled similar to (83) as follows:

$$\mathcal{S}_4 \leq 8 \sqrt{\frac{L'' \left(1 + \varepsilon_{\text{opt}} + \frac{C_S \|1\|_*}{N(1-\gamma)} \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N}} \mathbf{1}. \tag{101}$$

887 for *sa*-rectangular and the same quantity replacing $\max\{1-\gamma, C_g \sigma\}$ by $\max\{1-$
 888 $\gamma, \min_s \|\widehat{\pi}_s\|_* \tilde{\sigma} C_g\}$ for *s*-rectangular case.

889 • For \mathcal{S}_5 , it is observed that

$$\begin{aligned}
 \mathcal{S}_5 & = 2 \sqrt{\frac{L''}{N}} (I - \gamma \underline{P}^{\widehat{\pi},V})^{-1} \sqrt{\text{Var}_{P^{\widehat{\pi},V}}(\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma})} \\
 & \leq 2 \sqrt{\frac{L''}{(1-\gamma)^2 N}} \left\| V^{\widehat{\pi},\sigma} - \widehat{V}^{\widehat{\pi},\sigma} \right\|_\infty \mathbf{1}. \tag{102}
 \end{aligned}$$

890 • Next, observing that \mathcal{S}_6 and \mathcal{S}_7 are almost the same as the terms \mathcal{S}_2 (controlled in (95)) and
 891 \mathcal{S}_3 (controlled in (96)) in (93), it is easily verified that they can be controlled as follows

$$\mathcal{S}_6 \leq 4 \sqrt{\frac{\varepsilon_{\text{opt}} L''}{(1-\gamma)^4 N}} \mathbf{1}, \quad \mathcal{S}_7 \leq 4 \sqrt{\frac{L''}{\gamma^2 (1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N}} \mathbf{1}. \tag{103}$$

892 for sa -rectangular and the same quantity replacing $\max\{1 - \gamma, \sigma\}$ by $\max\{1 - \gamma, \min_s \|\hat{\pi}_s\|_* \tilde{\sigma}\}$
 893 for s -rectangular case. Then inserting the results in (101), (102), and (103) back to (100) leads to

$$\begin{aligned}
 & (I - \gamma \underline{P}^{\hat{\pi}, V})^{-1} \left(\hat{\underline{P}}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right) \tag{104} \\
 & \leq \left(\frac{2\varepsilon_{\text{opt}}}{(1-\gamma)} \right) 1 + 8 \sqrt{\frac{L'' \left(1 + \varepsilon_{\text{opt}} + \frac{C_S \|1\|_*}{N(1-\gamma)} \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1 + \frac{14L'' C_S \|1\|_*}{N(1-\gamma)^2} 1 \\
 & \quad + 2 \sqrt{\frac{L''}{(1-\gamma)^2 N}} \|V^{\hat{\pi}, \sigma} - \hat{V}^{\hat{\pi}, \sigma}\|_{\infty} 1 + 4 \sqrt{\frac{L'' \varepsilon_{\text{opt}}}{(1-\gamma)^4 N}} 1 + 4 \sqrt{\frac{L''}{\gamma^2 (1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N}} 1 \\
 & \leq 12 \sqrt{\frac{L'' \left(1 + \varepsilon_{\text{opt}} + \frac{C_S \|1\|_*}{N(1-\gamma)} \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} + 4 \sqrt{\frac{L''}{(1-\gamma)^2 N}} \|V^{\hat{\pi}, \sigma} - \hat{V}^{\hat{\pi}, \sigma}\|_{\infty} 1 \tag{105} \\
 & \quad + \frac{3\varepsilon_{\text{opt}}}{(1-\gamma)} + \frac{14L'' C_S \|1\|_*}{N(1-\gamma)^2} 1. \tag{106} \\
 & \tag{107}
 \end{aligned}$$

894 Taking $N \geq \frac{16L''}{1-\gamma}$, we obtain $\frac{2\varepsilon_{\text{opt}}}{(1-\gamma)} + 4\varepsilon_{\text{opt}} \sqrt{\frac{L''}{(1-\gamma)^4 N}} 1 \leq \frac{3\varepsilon_{\text{opt}}}{(1-\gamma)}$ with probability at least $1 - \delta$,
 895 inserting (97) and (105) back to (61)

$$\begin{aligned}
 & \|\hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma}\|_{\infty} \leq \max \left\{ 16 \sqrt{\frac{L'' \left(1 + \varepsilon_{\text{opt}} + \frac{C_S \|1\|_*}{N(1-\gamma)} \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1 + \left(\frac{2\varepsilon_{\text{opt}} \gamma}{(1-\gamma)} + \frac{14L'' C_S \|1\|_*}{N(1-\gamma)^2} 1 \right), \right. \\
 & \quad 12 \sqrt{\frac{L'' \left(1 + \varepsilon_{\text{opt}} + \frac{C_S \|1\|_*}{N(1-\gamma)} \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} + 4 \sqrt{\frac{L''}{(1-\gamma)^2 N}} \|V^{\hat{\pi}, \sigma} - \hat{V}^{\hat{\pi}, \sigma}\|_{\infty} 1 \tag{108} \\
 & \quad \left. + \frac{3\varepsilon_{\text{opt}}}{(1-\gamma)} + \frac{14L'' C_S \|1\|_*}{N(1-\gamma)^2} 1. \right\} \\
 & \leq 48 \sqrt{\frac{L'' \left(1 + \varepsilon_{\text{opt}} + \frac{C_S \|1\|_*}{N(1-\gamma)} \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N}} + \frac{6\varepsilon_{\text{opt}}}{(1-\gamma)} + \frac{28L'' C_S \|1\|_*}{N(1-\gamma)^2} 1 \tag{109}
 \end{aligned}$$

896 for sa -rectangular and the same quantity, replacing $\max\{1 - \gamma, C_g \sigma\}$ by $\max\{1 - \gamma, \tilde{\sigma} \min_s \|\hat{\pi}_s\|_*\}$
 897 for s -rectangular case. The proof is similar for TV without the geometric term depending on C_S .

898 **Step 6: summing all the previous inequalities results.** Using all the previous results in (90) and
 899 (109) and inserting back to (56) complete the proof as follows: taking $N \geq \frac{16L''}{(1-\gamma)^2}$, $\gamma > 1/4$, with
 900 probability at least $1 - \delta$, for sa -rectangular

$$\begin{aligned}
 & \|V^{*, \sigma} - V^{\hat{\pi}, \sigma}\|_{\infty} \leq \|V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma}\|_{\infty} + \varepsilon_{\text{opt}} + \|\hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma}\|_{\infty} \\
 & \leq \varepsilon_{\text{opt}} + 48 \sqrt{\frac{L'' \left(1 + \varepsilon_{\text{opt}} + \frac{C_S \|1\|_*}{N(1-\gamma)} \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N}} + \frac{6\varepsilon_{\text{opt}}}{(1-\gamma)} + \frac{28L'' C_S \|1\|_*}{N(1-\gamma)^2} 1 \\
 & \quad + 160 \sqrt{\frac{L \left(1 + \frac{C_S \|1\|_*}{N(1-\gamma)} \right)}{(1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N}} + \frac{14L C_S \|1\|_*}{N(1-\gamma)^2} \\
 & \leq \frac{8\varepsilon_{\text{opt}}}{1-\gamma} + \frac{42L'' C_S \|1\|_*}{N(1-\gamma)^2} + 1508 \sqrt{\frac{L'' \left(1 + \frac{C_S \|1\|_*}{N(1-\gamma)} \right)}{(1-\gamma)^2 \max\{1-\gamma, C_g \sigma\} N}} \tag{110}
 \end{aligned}$$

901 where the last inequality holds by $\gamma \geq \frac{1}{4}$ and $N \geq \frac{16L''}{(1-\gamma)^2}$ for sa -rectangular and the same quantity
 902 replacing $\max\{1 - \gamma, \sigma\}$ by $\max\{1 - \gamma, \tilde{\sigma} \min_s \{\|\pi_s^*\|_*\}\}$ for s -rectangular case. The proof is
 903 similar for TV without the geometric term depending on C_S .

904 9.3 Proof of the auxiliary lemmas

905 9.3.1 Proof of Lemma 5

906 Similarly to Shi et al. [2023], denoting s_0 the argmax of $V^{\pi, \sigma}$ such that $V^{\pi, \sigma}(s_0) = \min_{s \in \mathcal{S}} V^{\pi, \sigma}(s)$
 907 using recursive Bellman's equation

$$\max_{s \in \mathcal{S}} V^{\pi, \sigma}(s) = \max_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a})} \mathcal{P}V^{\pi, \sigma} \right] \quad (111)$$

$$\leq \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} \left(1 + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a})} \mathcal{P}V^{\pi, \sigma} \right) \quad (112)$$

908 where the second line holds since the reward function $r(s, a) \in [0, 1]$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.

909 Then we construct for any $(s, a) \in \mathcal{S} \times \mathcal{A}$ $\tilde{P}_{s,a} \in \mathbb{R}^{\mathcal{S}}$ by reducing the values of some elements of
 910 $P_{s,a}$ such that $P_{s,a} \geq \tilde{P}_{s,a} \geq 0$ and $\sum_{s'} (P_{s,a}(s') - \tilde{P}_{s,a}(s')) = \sigma C_g^{s,a}$. with $C_g^{s,a} = \frac{1}{\|e_{s_0}\|}$ It
 911 lead to $\tilde{P}_{s,a} + \sigma C_g^{s,a} e_{s_0}^\top \in \mathcal{U}_{\|\cdot\|}^\sigma(P_{s,a})$, where e_{s_0} is the standard basis vector supported on s_0 , since

$$\frac{1}{2} \left\| \tilde{P}_{s,a} + \sigma C_g^{s,a} e_{s_0}^\top - P_{s,a} \right\| \leq \frac{1}{2} \left\| \tilde{P}_{s,a} - P_{s,a} \right\| + \frac{C_g^{s,a} \sigma \|e_{s_0}\|}{2} = \sigma/2 + \sigma/2 = \sigma \quad (113)$$

912 Consequently,

$$\inf_{\mathcal{P} \in \mathcal{U}_{\|\cdot\|}^\sigma(P_{s,a})} \mathcal{P}V^{\pi, \sigma} \leq \left(\tilde{P}_{s,a} + \sigma C_g^{s,a} e_{s_0}^\top \right) V^{\pi, \sigma} \leq \left\| \tilde{P}_{s,a} \right\|_1 \|V^{\pi, \sigma}\|_\infty + \sigma V^{\pi, \sigma}(s_0) C_g \quad (114)$$

$$\leq (1 - C_g^{s,a} \sigma) \max_{s \in \mathcal{S}} V^{\pi, \sigma}(s) + \sigma C_g^{s,a} \min_{s \in \mathcal{S}} V^{\pi, \sigma}(s) \quad (115)$$

913 where the second inequality holds by

$$\left\| \tilde{P}_{s,a} \right\|_1 = \sum_{s'} \tilde{P}_{s,a}(s') = - \sum_{s'} (P_{s,a}(s') - \tilde{P}_{s,a}(s')) + \sum_{s'} P_{s,a}(s') = 1 - \sigma C_g^{s,a} \quad (116)$$

914 Plugging this back to the previous relation gives

$$\max_{s \in \mathcal{S}} V^{\pi, \sigma}(s) \leq 1 + \gamma(1 - C_g^{s,a} \sigma) \max_{s \in \mathcal{S}} V^{\pi, \sigma}(s) + \gamma C_g^{s,a} \sigma \min_{s \in \mathcal{S}} V^{\pi, \sigma}(s) \quad (117)$$

915 which, by rearranging terms, yields

$$\max_{s \in \mathcal{S}} V^{\pi, \sigma}(s) \leq \frac{1 + \gamma C_g^{s,a} \sigma \min_{s \in \mathcal{S}} V^{\pi, \sigma}(s)}{1 - \gamma(1 - C_g^{s,a} \sigma)} \quad (118)$$

$$\leq \frac{1}{(1 - \gamma) + \gamma C_g^{s,a} \sigma} + \min_{s \in \mathcal{S}} V^{\pi, \sigma}(s) \leq \frac{1}{\gamma \max\{1 - \gamma, C_g^{s,a} \sigma\}} + \min_{s \in \mathcal{S}} V^{\pi, \sigma}(s) \quad (119)$$

916 So rearranging term it holds :

$$\text{sp}(V^{\pi, \sigma})_\infty \leq \frac{1}{\gamma \max\{1 - \gamma, C_g \sigma\}} \quad (120)$$

917 As we pick the supreme over s on this quantity, $C_g^{s,a}$ is replaced by $C_g = 1/(\min_s \|e_s\|)$ to obtain a
 918 control for every s .

919 **9.3.2 Proof of Lemma 6**

920 Similarly to 5 denoting s_0 the argmax of $V^{\pi,\sigma}$ such that $V^{\pi,\sigma}(s_0) = \min_{s \in \mathcal{S}} V^{\pi,\sigma}(s)$ using recursive
921 Bellman's equation

$$\max_{s \in \mathcal{S}} V^{\pi,\sigma}(s) = \max_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_s)} \mathcal{P} V^{\pi,\tilde{\sigma}} \right] \quad (121)$$

$$\leq \max_{(s) \in \mathcal{S}} \left(1 + \gamma \inf_{\mathcal{P}^\pi \in \mathcal{U}^\sigma(P_s^\pi)} \mathcal{P}^\pi V^{\pi,\tilde{\sigma}} \right) \quad (122)$$

922 where the second line holds since the reward function $r(s, a) \in [0, 1]$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. Then
923 we construct for any $(s) \in \mathcal{S}$ $\tilde{P}_s \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ by reducing the values of some elements of P_s such that
924 $P_s \geq \tilde{P}_s \geq 0$ and

$$\forall a \in \mathcal{A}, \sum_{s'} \left(P_s(s', a) - \tilde{P}_s(s', a) \right) = \sigma_{s,a} C_g^s$$

925 Writing $\|\sigma_{s,a}\| \leq \tilde{\sigma}$ we construction $\sigma_{s,a}$ such that

$$\sum_a \pi(a|s) \sum_{s'} \left(P_s(s', a) - \tilde{P}_s(s', a) \right) = \|\pi_s\|_* \tilde{\sigma} C_g^s \quad (123)$$

926 Not that this construction is possible as it is simply Cauchy Swartz equality case.

927 It leads to $\tilde{P}_s + \sigma e_{s_0,a}^\top \in \mathcal{U}^{\tilde{\sigma}}(P_s)$, where $e_{s_0,a} \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ is the standard basis vector supported on s_0
928 which is equal to 1 at s_0 for every a and otherwise.

$$\frac{1}{2} \left\| \tilde{P}_s + \sigma_{s,a} C_g^s e_{s_0,a}^\top - P_s \right\| \leq \frac{1}{2} \left\| \tilde{P}_s - P_s \right\| + \frac{\tilde{\sigma} \|e_{s_0}\| C_g}{2} = \tilde{\sigma}/2 + \tilde{\sigma}/2 \quad (124)$$

929 as $C_g^s \|\sigma_{s,a} e_{s_0,a}\|$ is equal to $C_g^s \tilde{\sigma} \|e_{s_0}\|$ Consequently,

$$\inf_{\mathcal{P}^\pi \in \mathcal{U}^\sigma(P_s)} \mathcal{P}^\pi V^{\pi,\tilde{\sigma}} \leq \Pi^\pi \left(\tilde{P}_s^\pi + \sigma C_g^s e_{s_0}^\top \right) V^{\pi,\tilde{\sigma}} \quad (125)$$

$$= \sum_a \sum_{s'} \tilde{P}_s(s', a) \pi(a|s) V^{\pi,\tilde{\sigma}}(s') + \sigma e_{s_0,a} C_g^s V^{\pi,\tilde{\sigma}}(s_0) \pi(a|s) \quad (126)$$

$$= \sum_a \sup_{s'} V(s') \left(\sum_{s'} \tilde{P}_s(s', a) \right) \pi(a|s) + V^{\pi,\tilde{\sigma}}(s_0) \pi(a|s) \sigma_{s,a} C_g^s \quad (127)$$

$$\stackrel{(a)}{=} \max_{s \in \mathcal{S}} V^{\pi,\sigma}(s) \sum_a (1 - \sigma C_g^s) \pi(a|s) + \sum_a V^{\pi,\tilde{\sigma}}(s_0) \pi(a|s) \sigma_{s,a} C_g^s \quad (128)$$

$$\stackrel{(b)}{=} \max_{s \in \mathcal{S}} V^{\pi,\sigma}(s) (1 - \tilde{\sigma} C_g^s) \|\pi_s\|_* + \|\pi_s\|_* \tilde{\sigma} C_g^s \min_{s \in \mathcal{S}} V^{\pi,\tilde{\sigma}}(s) \quad (129)$$

$$\leq (1 - C_g^s \tilde{\sigma}) \max_{s \in \mathcal{S}} V^{\pi,\sigma}(s) + \sigma C_g^s \min_{s \in \mathcal{S}} V^{\pi,\tilde{\sigma}}(s) \quad (130)$$

930 where $\|\pi\|_\infty$ is the norm of the vector $\pi(\cdot|s)$ and where (a) holds because

$$\sum_{s'} \tilde{P}_s(s') = - \sum_{s'} \left(P_s(s') - \tilde{P}_s(s') \right) + \sum_{s'} P_s(s') = 1 - \sigma_{s,a} C_g^s \quad (131)$$

931 Finally (b) is due to (123). Plugging this back to the previous relation gives

$$\max_{s \in \mathcal{S}} V^{\pi,\tilde{\sigma}}(s) \leq 1 + \gamma (1 - \tilde{\sigma} C_g^s \|\pi_s\|_*) \max_{s \in \mathcal{S}} V^{\pi,\sigma}(s) + \gamma \|\pi_s\|_* \sigma C_g^s \min_{s \in \mathcal{S}} V^{\pi,\tilde{\sigma}}(s) \quad (132)$$

932 which, by rearranging terms, yields

$$\max_{s \in \mathcal{S}} V^{\pi, \tilde{\sigma}}(s) \leq \frac{1 + \gamma \tilde{\sigma} \|\pi_s\|_* C_g^s \min_{s \in \mathcal{S}} V^{\pi, \tilde{\sigma}}(s)}{1 - \gamma(1 - C_g^s \tilde{\sigma} \|\pi_s\|_*)} \quad (133)$$

$$\leq \frac{1}{(1 - \gamma) + \|\pi_s\|_* \gamma C_g^s \tilde{\sigma}} + \min_{s \in \mathcal{S}} V^{\pi, \tilde{\sigma}}(s) \quad (134)$$

$$\leq \frac{1}{(1 - \gamma) + \gamma \|\pi_s\|_* C_g^s \tilde{\sigma}} + \min_{s \in \mathcal{S}} V^{\pi, \tilde{\sigma}}(s) \quad (135)$$

$$\leq \frac{1}{\gamma \max\{1 - \gamma, C_g^s \|\pi_s\|_* \tilde{\sigma}\}} + \min_{s \in \mathcal{S}} V^{\pi, \tilde{\sigma}}(s) \quad (136)$$

933 So rearranging and taking the sumpremum over all stern it holds :

$$\text{sp}(V^{\pi, \tilde{\sigma}})_\infty \leq \frac{1}{\gamma \max\{1 - \gamma, \min_s \|\pi_s\|_* C_g \tilde{\sigma}\}} \quad (137)$$

934 As we pick the supreme over s ovf this quantity, C_g^s is replaced by $C_g = 1/\min_s \|\pi_s\|_*$

935 9.3.3 Proof of Lemma 8

936 *Proof.* Concentration of the robust values function. with probability $1 - \delta$, it holds:

$$\left| P_{s,a}^{\pi, V} V - \widehat{P}_{s,a}^{\pi, V} V \right| \leq 2\sqrt{\frac{L}{N}} \sqrt{\text{Var}[V]_{\alpha^{**}}(V)} + \frac{3LC_S \|1\|_*}{N(1 - \gamma)}$$

937 with $L = 2 \log(18 \|1\|_* SAN/\delta)$ and First we can use optimization duality such as in (50):

$$\left| P_{s,a}^{\pi, V} V - \widehat{P}_{s,a}^{\pi, V} V \right| \quad (138)$$

$$\begin{aligned} &= \left| \max_{\mu_{P_{s,a}^0}^{\lambda, \omega} \in \mathcal{M}_{P_{s,a}^0}^{\lambda, \omega}} \left\{ P_{s,a}^0(V - \mu) - \sigma(\text{sp}((V - \mu)_*)) \right\} \right. \\ &\quad \left. - \max_{\mu_{\widehat{P}_{s,a}^0}^{\lambda, \omega} \in \mathcal{M}_{\widehat{P}_{s,a}^0}^{\lambda, \omega}} \left\{ \widehat{P}_{s,a}^0(V - \mu_{\widehat{P}_{s,a}^0}^{\lambda, \omega}) - \sigma(\text{sp}((V - \mu_{\widehat{P}_{s,a}^0}^{\lambda, \omega}))_*) \right\} \right| \\ &\leq \max \left\{ \left| \max_{\mu_{P_{s,a}^0}^{\lambda, \omega} \in \mathcal{M}_{P_{s,a}^0}^{\lambda, \omega}} \left\{ P_{s,a}^0(V - \mu_{P_{s,a}^0}^{\lambda, \omega}) - \sigma(\text{sp}((V - \mu_{P_{s,a}^0}^{\lambda, \omega}))_*) \right\} \right. \right. \\ &\quad \left. \left. - \max_{\mu_{\widehat{P}_{s,a}^0}^{\lambda, \omega} \in \mathcal{M}_{\widehat{P}_{s,a}^0}^{\lambda, \omega}} \left\{ \widehat{P}_{s,a}^0(V - \mu_{\widehat{P}_{s,a}^0}^{\lambda, \omega}) - \sigma(\text{sp}((V - \mu_{\widehat{P}_{s,a}^0}^{\lambda, \omega}))_*) \right\} \right| \right\}; \quad (139) \end{aligned}$$

$$\left| \max_{\mu_{\widehat{P}_{s,a}^0}^{\lambda, \omega} \in \mathcal{M}_{\widehat{P}_{s,a}^0}^{\lambda, \omega}} \left\{ \widehat{P}_{s,a}^0(V - \mu_{\widehat{P}_{s,a}^0}^{\lambda, \omega}) - \sigma(\text{sp}((V - \mu_{\widehat{P}_{s,a}^0}^{\lambda, \omega}))_*) \right\} \right| \quad (140)$$

$$\begin{aligned} &- \max_{\mu_{\widehat{P}_{s,a}^0}^{\lambda, \omega} \in \mathcal{M}_{\widehat{P}_{s,a}^0}^{\lambda, \omega}} \left\{ P_{s,a}^0(V - \mu_{\widehat{P}_{s,a}^0}^{\lambda, \omega}) - \sigma(\text{sp}((V - \mu_{\widehat{P}_{s,a}^0}^{\lambda, \omega}))_*) \right\} \quad \left| \right\} \\ &\leq \max \left\{ \underbrace{\left| \max_{\mu \in \mathcal{M}_{P_{s,a}^0}^{\lambda, \omega}} (P_{s,a}^0 - \widehat{P}_{s,a}^0)(V - \mu_{P_{s,a}^0}^{\lambda, \omega}) \right|}_{=: g_{s,a}(\alpha_{P^0}^{\lambda, \omega}, V)}, \underbrace{\left| \max_{\mu_{\widehat{P}_{s,a}^0}^{\lambda, \omega} \in \mathcal{M}_{\widehat{P}_{s,a}^0}^{\lambda, \omega}} (P_{s,a}^0 - \widehat{P}_{s,a}^0)(V - \mu_{\widehat{P}_{s,a}^0}^{\lambda, \omega}) \right|}_{=: g_{s,a}(\alpha_{\widehat{P}^0}^{\lambda, \omega}, V)} \right\} \quad (141) \end{aligned}$$

938 where in the first equality we use Lemma 3. The final inequality is a consequence of the 1-
939 Lipschitzness of the max operator. First, we control $g_{s,a}(\alpha_{P^0}^{\lambda, \omega}, V)$. To do so, we use for a fixed $\alpha_P^{\lambda, \omega}$

940 and any vector V that is independent with \widehat{P}^0 , the Bernstein's inequality, one has with probability at
 941 least $1 - \delta$ with sa -rectangular notations,

$$g_{s,a}(\alpha_P^{\lambda,\omega}, V) = \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_{\alpha_P^{\lambda,\omega}} \right| \leq \sqrt{\frac{2 \log(\frac{2}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{2 \log(\frac{2}{\delta})}{3N(1-\gamma)}. \quad (142)$$

942 Once pointwise concentration derived, we will use uniform concentration to yield this lemma. First,
 943 union bound, is obtained noticing that $g_{s,a}(\alpha_P^{\lambda,\omega}, V)$ is 1-Lipschitz w.r.t. λ and ω as it is linear in
 944 λ and ω . Moreover, $\lambda^* = \|V - \mu^* - \omega\|_*$ obeying $\lambda^* \leq \frac{\|1\|_*}{1-\gamma}$. The quantity $\omega \in [0, 1/(1-\gamma)]$
 945 as it is always smaller than V by definition. We construct then a 2-dimensional ε_1 -net N_{ε_1} over
 946 $\lambda^* \in [0, \frac{\|1\|_*}{1-\gamma}]$ and $\omega \in [0, 1/(1-\gamma)]$ whose size satisfies $|N_{\varepsilon_1}| \leq \left(\frac{3\|1\|_*}{\varepsilon_1(1-\gamma)} \right)^2$ [Vershynin, 2018].
 947 Using union bound and (142), it holds with probability at least $1 - \frac{\delta}{SA}$ that for all $\lambda \in N_{\varepsilon_1}$,

$$g_{s,a}(\alpha_P^\lambda, V) \leq \sqrt{\frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{3N(1-\gamma)}. \quad (143)$$

948 Using the previous equation and also (141), it results in using notation $2 \log(\frac{18SAN}{\delta}) = L$,

$$\begin{aligned} g_{s,a}(\alpha_P^\lambda, V) &\stackrel{(a)}{\leq} \sup_{\alpha_P^\lambda \in N_{\varepsilon_1}} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_{\alpha_P^\lambda} \right| + \varepsilon_1 \\ &\stackrel{(b)}{\leq} \sqrt{\frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{3N(1-\gamma)} + \varepsilon_1 \end{aligned} \quad (144)$$

$$\begin{aligned} &\stackrel{(c)}{\leq} \sqrt{\frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{\log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N(1-\gamma)} \\ &\stackrel{(d)}{\leq} 2\sqrt{\frac{L}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{L}{N(1-\gamma)} \end{aligned} \quad (145)$$

$$\begin{aligned} &\leq 2\sqrt{\frac{L}{N}} \|V\|_\infty + \frac{L}{N(1-\gamma)} \\ &\leq 3\sqrt{\frac{L}{(1-\gamma)^2 N}} \end{aligned} \quad (146)$$

949 where (a) is because the optimal α^* falls into the ε_1 -ball centered around some point inside N_{ε_1} and
 950 $g_{s,a}(\alpha_P^\lambda, V)$ is 1-Lipschitz with regard to λ and ω , (b) is due to Eq. (143), (c) arises from taking
 951 $\varepsilon_1 = \frac{\log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{3N(1-\gamma)}$, (d) is verified by $|N_{\varepsilon_1}| \leq \left(\frac{3\|1\|_*}{\varepsilon_1(1-\gamma)} \right)^2 \leq 9N \|1\|$ and that variance of a ceiling
 952 function of a vector is smaller than the variance of non-ceiling vector, and the last inequality comes
 953 from the fact $\|V^{*,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$ and taking $N \geq 2 \log(\frac{18SAN\|1\|_*}{\delta}) = L$.

954 Contrary to the previous term, the second term $g_{s,a}(\alpha_P^\lambda, V)$ is more difficult as we need concentration,
 955 but there is an extra dependency in the data through the parameter α_P^λ . We need to decouple this
 956 problem using absorbing MDPs. Then it leads to

$$g_{s,a}(\alpha^{\lambda,\omega}, V) \quad (147)$$

$$= \left| \max_{\substack{\mu_{\hat{P}_{s,a}}^{\lambda,\omega} \in \mathcal{M}_{\hat{P}_{s,a}}^{\lambda,\omega}}} \left(P_{s,a}^0 - \hat{P}_{s,a}^0 \right) (V - \mu_{\hat{P}_{s,a}}^{\lambda,\omega}) \right| \quad (148)$$

$$= \left| \max_{\mu \in \mathcal{M}_{\hat{P}_{s,a}}^{\lambda,\omega}} \left(P_{s,a}^0 - \hat{P}_{s,a}^0 \right) (V - \mu_{P_{s,a}^0}^{\lambda,\omega}) + \left(P_{s,a}^0 - \hat{P}_{s,a}^0 \right) (\mu_{P_{s,a}^0}^{\lambda,\omega} - \mu_{\hat{P}_{s,a}^0}^{\lambda,\omega}) \right| \quad (149)$$

$$\leq \left| \max_{\substack{\mu_{P_{s,a}^0}^{\lambda,\omega} \in \mathcal{M}_{P_{s,a}^0}^{\lambda,\omega}}} \left(P_{s,a}^0 - \hat{P}_{s,a}^0 \right) (V - \mu_{P_{s,a}^0}^{\lambda,\omega}) + \max_{\substack{\mu_{\hat{P}_{s,a}^0}^{\lambda,\omega} \in \mathcal{M}_{\hat{P}_{s,a}^0}^{\lambda,\omega}}} \left(P_{s,a}^0 - \hat{P}_{s,a}^0 \right) (\mu_{P_{s,a}^0}^{\lambda,\omega} - \mu_{\hat{P}_{s,a}^0}^{\lambda,\omega}) \right| \quad (150)$$

957 In the first equality, we add the term $\mu_{P_{s,a}^0}^{\lambda,\omega}$ to retrieve the previous concentration problem, fixing $P_{s,a}^0$
 958 and optimizing λ, ω . In the second, we extend the max using triangular inequality. The first term in
 959 the last equality is exactly the term we have controlled previously, while the second one needs more
 960 attention. We decouple the dependency of the data, and then controlling the difference between the μ .
 961 Then using the characterization of the optimal μ from equation (47):

$$\left(P_{s,a}^0 - \hat{P}_{s,a}^0 \right) (\mu_{P_{s,a}^0}^{\lambda,\omega} - \mu_{\hat{P}_{s,a}^0}^{\lambda,\omega}) = \sum_{s'} \lambda \left(P_{s,a}^0(s') - \hat{P}_{s,a}^0(s') \right) (\nabla \|P_{s,a}^0\| - \nabla \|\hat{P}_{s,a}^0\|)$$

962 Here we assume that the subgradient are gradient as we assume that the norm is C^2 . The question
 963 that arises is whether the gradient if the norm is Lipschitz. Assuming that the norm is C^2 , using
 964 Mean value theorem, we know that

$$\left\| (\nabla \|P_{s,a}^0\| - \nabla \|\hat{P}_{s,a}^0\|) \right\|_2 \leq \sup_{x \in \Delta(S)} \|\nabla^2 \|x\|\|_2 \left\| (P_{s,a}^0 - \hat{P}_{s,a}^0) \right\|_2.$$

965 As the norm is C^2 , is continuous and as the simplex is bounded, this quantity exists according to
 966 Extreme value theorem. It is possible to compute this contact depending on S for explicit norm such
 967 as L_p . Indeed, for L_2 :

$$\nabla^2 \|x\|_2 = \frac{(I - \frac{x \otimes x}{\|x\|_2^2})}{\|x\|_2} \leq \frac{1}{\|x\|_2} I \leq \frac{1}{\min_{x \in \Delta(S)} \|x\|_2} I = \sqrt{S}$$

968 where \otimes is the Kronecker product. So we have an upper bound independently of x . For $L_p = \|x\|_p$
 969 norms, $p \geq 2$, we have simple taking derivative twice:

$$\nabla^2 \|x\|_p = \frac{p-1}{L_p} (\mathcal{A}^{p-2} - g_p g_p^T)$$

970 with

$$\mathcal{A} = \text{Diag} \left(\frac{\text{abs}(x)}{L_p} \right)$$

$$g_p = \mathcal{A}^{p-2} \left(\frac{x}{L_p} \right).$$

971 where Diag is the diagonal matrix. However, as $x \leq L_p$, $\mathcal{A} \leq I$, we get

$$H \leq \frac{p-1}{\|x\|_p} \leq (p-1)S^{1/q} = C_S \quad (151)$$

972 where the $1/L_p$ is minimized for the uniform distribution. Then using Cauchy Swartz inequality, it
 973 holds

$$\left(P_{s,a}^0 - \hat{P}_{s,a}^0 \right) (\mu_{P_{s,a}^0}^{\lambda,\omega} - \mu_{\hat{P}_{s,a}^0}^{\lambda,\omega}) \leq \lambda \left\| \left(P_{s,a}^0 - \hat{P}_{s,a}^0 \right) \right\|_2^2. \quad (152)$$

974 Then the question is how to bound the quantity $\left\| \left(P_{s,a}^0 - \hat{P}_{s,a}^0 \right) \right\|_2^2$. To do so, we will use Mac
 975 Diarmid inequality.

976 **Definition 3.** *Bounded difference property*

977 A function $f : \mathcal{X}_1 \times \dots \times \mathcal{X}_n \rightarrow \mathbb{R}$ satisfies the bounded difference property if for each $i = 1, \dots, n$
 978 the change of coordinate from s_i to s'_i may change the value of the function at most on c_i

$$\forall i \in [n] : \sup_{x'_i \in \mathcal{X}_i} |f(x_1, \dots, x_i, \dots, x_n) - f(x_1, \dots, x'_i, \dots, x_n)| \leq c_i$$

979 In our case, we consider $f(X_1, \dots, X_n) = \|\sum_{k=1}^n X_k\|_2$. Then we can notice that by triangle
 980 inequality for any x_1, \dots, x_n and x'_k with $X_{i,s'} = P_{i,s,a}^0(s') - P_{s,a}^0(s')$ (index i holds for index of
 981 sample generated from the generative model) that

$$\begin{aligned} f(x_1, \dots, x_k, \dots, x_n) &= \|x_1 + \dots + x_n\|_2 \leq \|x_1 + \dots + x_n - x_k + x'_k\|_2 + \|x_k - x'_k\|_2 \\ &\leq f(x_1, \dots, x'_k, \dots, x_n) + 2 \end{aligned}$$

982 **Theorem 5.** (*McDiarmid's inequality*). *McDiarmid et al. [1989]* Let $f : \mathcal{X}_1 \times \dots \times \mathcal{X}_n \rightarrow \mathbb{R}$ be a
 983 function satisfying the bounded difference property with bounds c_1, \dots, c_n . Consider independent
 984 random variables $X_1, \dots, X_n, X_i \in \mathcal{X}_i$ for all i . Then for any $t > 0$

$$\mathbb{P}[f(X_1, \dots, X_n) - \mathbb{E}[f(X_1, \dots, X_n)] \geq t] \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n c_i^2}\right)$$

985 Using McDiarmid's inequality and union bound, we can bound the term as here

$$\left\| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) \right\|_2^2 - \mathbb{E} \left[\left\| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) \right\|_2^2 \right] \leq \frac{2N \log(|S||A|/\delta)}{N^2}$$

986 with probability $1 - \delta/(|S||A|)$. Moreover, the additional term can be bounded as follows:

$$\mathbb{E} \left[\left\| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) \right\|_2^2 \right] = \mathbb{E} \left[\sum_{s'} (P_{s,a}^0(s') - P_{s,a}^0(s'))^2 \right] = \mathbb{E} \left[\sum_{s'} \left(\frac{1}{N} \sum_i X_{i,s'} \right)^2 \right]$$

987 with $X_{i,s'} = P_{i,s,a}^0(s') - P_{s,a}^0(s')$ is one sample sampled from the generative model. Then

$$\begin{aligned} \mathbb{E} \left[\left\| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) \right\|_2^2 \right] &= \frac{1}{N^2} \sum_{s'} \text{Var} \left(\sum_i X_{i,s} \right) \stackrel{a}{=} \frac{1}{N^2} \sum_i \sum_{s'} \text{Var}(X_{i,s}) \\ &= \frac{1}{N^2} \sum_i \mathbb{E} \left(\sum_{s'} X_{i,s}^2 \right) \leq \frac{4}{N} \end{aligned}$$

988 where (a) the last equality comes from the independence of the random variables and where the last
 989 inequality comes from the fact the maximum of two elements in the simplex is bounded by 2. Finally,
 990 regrouping the two terms, we obtain with probability $1 - \delta/(|S||A|)$:

$$\begin{aligned} \left\| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) \right\|_2^2 &\leq \frac{2N \log(|S||A|/(\delta))}{N^2} + \frac{4}{N} = \frac{8 \log(|S||A|/(\delta))}{N} + \frac{4}{N} \\ &\leq \frac{6 \log(|S||A|/(\delta))}{N} = \frac{L'}{N} \end{aligned}$$

991 with $L' = 6 \log(|S||A|/(\delta))$. Finally, plugging the previous equation in (152):

$$\max_{\mu \in \mu_{\widehat{P}_{s,a}^0}^\lambda} \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) (\mu_{P_{s,a}^0}^\lambda - \mu) \leq \max_\lambda \left\| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) \right\|_2^2 C_S \lambda.$$

992 This term can be easily controlled by taking the supremum over λ which is a 1 dimensional parameter.
 993 Then we can bound $\lambda \in [0, H \|1\|_*]$. Indeed,

$$\lambda^* = \|V - \mu^* - \eta\|_* \leq \|V\|_* \leq H \|1\|_*.$$

994 Finally, we obtain:

$$\max_{\lambda} \left\| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) \right\|_2^2 C_S \lambda \leq \frac{L' C_S \|1\|_*}{N(1-\gamma)}.$$

995 Regrouping all terms:

$$\begin{aligned} g_{s,a}(\alpha_{\widehat{P}}^\lambda, V) &\leq \left| \max_{\mu_{P_{s,a}^0}^\lambda \in \mathcal{M}_{P_{s,a}^0}^\lambda} \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) (V - \mu_{P_{s,a}^0}^\lambda) + \max_{\mu_{\widehat{P}_{s,a}^0}^\lambda \in \mathcal{M}_{\widehat{P}_{s,a}^0}^\lambda} \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) (\mu_{P_{s,a}^0}^\lambda - \mu_{\widehat{P}_{s,a}^0}^\lambda) \right| \\ &\leq 2\sqrt{\frac{L}{N}} \sqrt{\text{Var}(V)} + \frac{L' C_S \|1\|_*}{N(1-\gamma)} + \frac{L}{N(1-\gamma)} \leq 2\sqrt{\frac{L}{N}} \sqrt{\text{Var}(V)} + \frac{3L C_S \|1\|_*}{N(1-\gamma)} \end{aligned} \quad (153)$$

$$(154)$$

996 We can recognize that the second term is a second order term as long as $N \geq (C_S \|1\|_*)^2$, we can
 997 regroup the two terms. Finally, as $g_{s,a}(\alpha_{\widehat{P}}^\lambda, V) \geq g_{s,a}(\alpha_P^\lambda, V)$, we obtain

$$\left| P_{s,a}^{\pi,V} V - \widehat{P}_{s,a}^{\pi,V} V \right| \leq 2\sqrt{\frac{L}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{3L C_S \|1\|_*}{N(1-\gamma)} \quad (155)$$

998 It is important to note that the geometry of the norm is present in the second order term $\frac{3L C_S \|1\|_*}{N(1-\gamma)}$
 999 but this term is negligible as it is proportional to $1/N$ with regard to the variance term in $1/\sqrt{N}$.
 1000 Moreover, note that the quantity $C_S \|1\|_* = S$ for L_2 norms.

1001 For the specific case of TV which is not C^2 smooth, this lemma still holds as in (141), we only need
 1002 to control one term without the dependency on data in the supremum as $\alpha_{\widehat{P}}^\lambda$ reduces to a scalar α
 1003 which does not depend on P . Then extra decomposition using smoothness of the norm is not needed,
 1004 as the only remaining term in the max in (141) is the left hand side term.

1005 For the s -rectangular case, the first equation can be rewritten simply factorizing by $\pi(a|s)$ using
 1006 lemma 4.

$$\begin{aligned} \left| P_{s,a}^{\pi,V} V - \widehat{P}_{s,a}^{\pi,V} V \right| &= \left| \sum_a \pi(a|s) \max_{\mu_{P_{s,a}^0}^\lambda \in \mathcal{M}_{P_{s,a}^0}^\lambda} \left\{ P_{s,a}^0 (V - \mu) - \sigma(\text{sp}((V - \mu)_*)) \right\} \right. \\ &\quad \left. - \max_{\mu_{\widehat{P}_{s,a}^0}^\lambda \in \mathcal{M}_{\widehat{P}_{s,a}^0}^\lambda} \left\{ \widehat{P}_{s,a}^0 (V - \mu_{\widehat{P}_{s,a}^0}^\lambda) - \sigma(\text{sp}((V - \mu_{\widehat{P}_{s,a}^0}^\lambda)_*)) \right\} \right| \end{aligned} \quad (156)$$

$$\leq \sum_a \pi(a|s) \left(2\sqrt{\frac{L}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{L C_S \|1\|_*}{N(1-\gamma)} \right) \quad (157)$$

$$= 2\sqrt{\frac{L}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{3L C_S \|1\|_*}{N(1-\gamma)} \quad (158)$$

1007 using sa -rectangular results, which gives the result.

1008 Combining this lemma with a matrix notation, one has with probability $1 - \delta$:

$$\left| \underline{\widehat{P}}^{\pi^*,V} V^{\pi^*,\sigma} - \underline{P}^{\pi^*,V} V^{\pi^*,\sigma} \right| \leq 2\sqrt{\frac{L}{N}} \sqrt{\text{Var}_{P^*}(V^{*,\sigma})} + \frac{3L C_S \|1\|_*}{N(1-\gamma)} \quad (159)$$

$$(160)$$

1009

□

1010 **9.3.4 Proof of Lemma 9**

1011 Using the same argument as in (209), it holds that for any α^* solution of (??) or (53)

$$\left(I - \gamma \widehat{P}^{\pi^*, V}\right)^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, V}}(V^{*, \sigma})} = \sqrt{\frac{1}{1 - \gamma} \sum_{t=0}^{\infty} \gamma^t \left(\widehat{P}^{\pi^*, V}\right)^t \text{Var}_{\widehat{P}^{\pi^*, V}}([V^{*, \sigma}]_{\alpha^{**}})}. \quad (161)$$

1012 Then we can control $\text{Var}_{\widehat{P}^{\pi^*, V}}(V^{*, \sigma})$. Defining $V' := V^{*, \sigma} - \eta 1$, $\eta \in \mathbb{R}$, we use Bellman's equation
1013 in (32) which lead to

$$V' = V^{*, \sigma} - \eta 1 \leq V^{*, \sigma} - \eta 1 = r_{\pi^*} + \gamma \underline{P}^{\pi^*, V} V^{*, \sigma} - \eta 1 \quad (162)$$

$$= r_{\pi^*} + \gamma \underline{P}^{\pi^*, V} [V^{*, \sigma} - \gamma \sigma \text{sp}(V^{*, \sigma})_* - \eta 1] \quad (163)$$

$$= r'_{\pi^*} + \gamma \widehat{P}^{\pi^*, V} V' + \gamma \left(\underline{P}^{\pi^*, V} - \widehat{P}^{\pi^*, V}\right) V^{*, \sigma} - \gamma \sigma \text{sp}([V^{*, \sigma}]_*) \quad (164)$$

$$= r'_{\pi^*} + \gamma \widehat{P}^{\pi^*, V} V' + \gamma \left(\underline{P}^{\pi^*, V} - \widehat{P}^{\pi^*, V}\right) V^{*, \sigma} \quad (165)$$

$$\leq r'_{\pi^*} + \gamma \widehat{P}^{\pi^*, V} V' + \gamma \left(\underline{P}^{\pi^*, V} - \widehat{P}^{\pi^*, V}\right) V^{*, \sigma} \quad (166)$$

1014 where in the second line we use Lemma 3. and we define $r'_{\pi^*} = r_{\pi^*} - (1 - \gamma)\eta < r_{\pi^*} < 1$. We
1015 obtain the same result in s -rectangular case using lemma 4 instead. Then

$$\begin{aligned} \text{Var}_{\widehat{P}^{\pi^*, V}}([V^{*, \sigma}]) &\stackrel{(a)}{=} \text{Var}_{\widehat{P}^{\pi^*, V}}(V') = \widehat{P}^{\pi^*, V}(V' \circ V') - \left(\widehat{P}^{\pi^*, V} V'\right) \circ \left(\widehat{P}^{\pi^*, V} V'\right) \\ &= \widehat{P}^{\pi^*, V}(V' \circ V') - \left(\widehat{P}^{\pi^*, V} V'\right) \circ \left(\widehat{P}^{\pi^*, V} V'\right) \\ &\stackrel{(b)}{\leq} \widehat{P}^{\pi^*, V}(V' \circ V') - \frac{1}{\gamma^2} \left(V' - r'_{\pi^*} - \gamma \left(\underline{P}^{\pi^*, V} - \widehat{P}^{\pi^*, V}\right) V^{*, \sigma}\right)^{\circ 2} \\ &= \widehat{P}^{\pi^*, V}(V' \circ V') - \frac{1}{\gamma^2} V' \circ V' + \frac{2}{\gamma^2} V' \circ \left(r'_{\pi^*} + \gamma \left(\underline{P}^{\pi^*, V} - \widehat{P}^{\pi^*, V}\right) V^{*, \sigma}\right) \\ &\quad - \frac{1}{\gamma^2} \left(r'_{\pi^*} + \gamma \left(\underline{P}^{\pi^*, V} - \widehat{P}^{\pi^*, V}\right) V^{*, \sigma}\right)^{\circ 2} \\ &\stackrel{(c)}{\leq} \widehat{P}^{\pi^*, V}(V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} 1 \end{aligned} \quad (167)$$

$$+ \frac{2}{\gamma} \|V'\|_{\infty} \left| \left(\underline{P}^{\pi^*, V} - \widehat{P}^{\pi^*, V}\right) V^{*, \sigma} \right| \quad (168)$$

$$\leq \widehat{P}^{\pi^*, V}(V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} 1 \quad (169)$$

$$+ \frac{2}{\gamma} \|V'\|_{\infty} \left(2 \sqrt{\frac{L}{(1 - \gamma)^2 N}} + \frac{3C_S \|1\|_* L}{N(1 - \gamma)}\right) 1, \quad (170)$$

1016 where (a) holds by the fact that $\text{Var}_{P_{\pi}}(V - c1) = \text{Var}_{P_{\pi}}(V)$ for any scalar c and $V \in \mathbb{R}^S$, (b) follows
1017 from (166), (c) arises from $\frac{1}{\gamma^2} V' \circ V' \geq \frac{1}{\gamma} V' \circ V'$ and $-1 \leq r_{\pi^*} - (1 - \gamma)V_{\min} 1 = r'_{\pi^*} \leq r_{\pi^*} \leq 1$,

1018 and the last inequality holds by Lemma 8. Plugging (170) into (161) leads to

$$(I - \gamma \widehat{P}^{\pi^*, V})^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, V}}(V^*, \sigma)} \quad (171)$$

$$\leq \sqrt{\frac{1}{1-\gamma}} \left(\sum_{t=0}^{\infty} \gamma^t (\widehat{P}^{\pi^*, V})^t \left(\widehat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} 1 \right) \right) \quad (172)$$

$$+ \frac{2}{\gamma} \|V'\|_{\infty} \left(2 \sqrt{\frac{L}{(1-\gamma)^2 N}} + \frac{3C_S \|1\|_* L}{N(1-\gamma)} \right) 1 \Big)^{1/2}$$

$$\stackrel{(i)}{\leq} \sqrt{\frac{1}{1-\gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t (\widehat{P}^{\pi^*, V})^t \left(\widehat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' \right) \right|} \\ + \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (\widehat{P}^{\pi^*, V})^t \left(\frac{2}{\gamma^2} \|V'\|_{\infty} 1 + \frac{2}{\gamma} \|V'\|_{\infty} \left(2 \sqrt{\frac{L}{(1-\gamma)^2 N}} + \frac{3C_S \|1\|_* L}{N(1-\gamma)} \right) 1 \right)}$$

$$\leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t (\widehat{P}^{\pi^*, V})^t \left[\widehat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' \right] \right|} \quad (173)$$

$$+ \sqrt{\frac{\left(2 + 2 \left(2 \sqrt{\frac{L}{(1-\gamma)^2 N}} + \frac{3C_S \|1\|_* L}{N(1-\gamma)} \right) \right) \|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1, \quad (174)$$

1019 where (i) holds by the triangle inequality. Therefore, the remainder of the proof shall focus on the
1020 first term, which follows

$$\left| \sum_{t=0}^{\infty} \gamma^t (\widehat{P}^{\pi^*, V})^t \left(\widehat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' \right) \right| \\ = \left| \left(\sum_{t=0}^{\infty} \gamma^t (\widehat{P}^{\pi^*, V})^{t+1} - \sum_{t=0}^{\infty} \gamma^{t-1} (\widehat{P}^{\pi^*, V})^t \right) (V' \circ V') \right| \leq \frac{1}{\gamma} \|V'\|_{\infty}^2 1 \quad (175)$$

1021 by recursion. Inserting (175) back to (174) leads to

$$\left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, V}}(V^*, \sigma)_{\alpha^*}} \\ \leq \sqrt{\frac{\|V'\|_{\infty}^2}{\gamma(1-\gamma)} 1} + 3 \sqrt{\frac{\left(1 + \left(\sqrt{\frac{L}{(1-\gamma)^2 N}} + \frac{C_S \|1\|_* L}{N(1-\gamma)} \right) \right) \|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1 \\ \leq 4 \sqrt{\frac{\left(1 + \left(\sqrt{\frac{L}{(1-\gamma)^2 N}} + \frac{C_S \|1\|_* L}{N(1-\gamma)} \right) \right) \|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1 \quad (176)$$

$$\leq 4 \sqrt{\frac{\left(1 + \left(1 \sqrt{\frac{L}{(1-\gamma)^2 N}} + \frac{C_S \|1\|_* L}{N(1-\gamma)} \right) \right) \|V'\|_*$$

1022 Taking the infimum over η in the right-hand side, recall $V' := V^{*,\sigma} - \eta 1$, we obtain the definition of
 1023 the span semi norm.

$$\begin{aligned} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, V}}(V^{*,\sigma})_{\alpha^*}} &\leq 4 \sqrt{\frac{\left(1 + \left(\sqrt{\frac{L}{(1-\gamma)^2 N} + \frac{C_S \|1\|_* L}{N(1-\gamma)}} \right) \text{sp}(V^{*,\sigma})_* \right)}{(1-\gamma)^2 \gamma^2}} 1 \\ &\leq 4 \sqrt{\frac{\left(1 + \left(\sqrt{\frac{L}{(1-\gamma)^2 N} + \frac{C_S \|1\|_* L}{N(1-\gamma)}} \right) \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, C_g \sigma\}}} 1 \end{aligned} \quad (178)$$

$$\leq 4 \sqrt{\frac{\left(1 + \left(\sqrt{\frac{L}{(1-\gamma)^2 N} + \frac{C_S \|1\|_* L}{N(1-\gamma)}} \right) \right)}{\gamma^3 (1-\gamma)^3}} 1, \quad (179)$$

1024 where the penultimate inequality follows from applying Lemma 5 with $P = P^0$ and $\pi = \pi^*$:

$$\text{sp}(V^{*,\sigma})_* \leq \frac{1}{\gamma \max\{1-\gamma, C_g \sigma\}}.$$

1025 or with an extra factor for s rectangular assumptions.

$$\text{sp}(V^{*,\sigma})_* \leq \frac{1}{\gamma \max\{1-\gamma, \min_s \|\pi_s\|_* \tilde{\sigma} C_g\}}.$$

1026 9.3.5 Proof of Lemma 10

1027 In this proof, we will sa -rectangular notations, especially $\alpha_{s,a}^{**}$ but it holds also for α_s^{**} and s -
 1028 rectangular case. For any $(s, a) \in \mathcal{S} \times \mathcal{A}$, using the results in (141), for both sa -rectangular case:

$$\left| \widehat{P}_{s,a}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - P_{s,a}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right| \leq \max \left\{ \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha_{P_{s,a}}^{\lambda, \omega^*}} \right|, \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha_{\widehat{P}_{s,a}}^{\lambda, \omega^*}} \right| \right\} \quad (180)$$

1029 The first term in this max can be bounded using:

$$\begin{aligned} &\left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha_{P_{s,a}}^{\lambda, \omega^*}} \right| \quad (181) \\ &\stackrel{(a)}{\leq} \left(\left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha_{P_{s,a}}^{\lambda, \omega^*}} \right| + \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) \left([\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha_{P_{s,a}}^{\lambda, \omega^*}} - [\widehat{V}^{*, \sigma}]_{\alpha_{P_{s,a}}^{\lambda, \omega^*}} \right) \right| \right) \\ &\leq \left(\left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{*, \sigma}]_{\alpha_{P_{s,a}}^{\lambda, \omega^*}} \right| + \left\| P_{s,a}^0 - \widehat{P}_{s,a}^0 \right\|_1 \left\| [\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha_{P_{s,a}}^{\lambda, \omega^*}} - [\widehat{V}^{*, \sigma}]_{\alpha_{P_{s,a}}^{\lambda, \omega^*}} \right\|_\infty \right) \\ &\stackrel{(b)}{\leq} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{*, \sigma}]_{\alpha_{P_{s,a}}^{\lambda, \omega^*}} \right| + 2 \left\| \widehat{V}^{\widehat{\pi}, \sigma} - \widehat{V}^{*, \sigma} \right\|_\infty \\ &\stackrel{(c)}{\leq} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{*, \sigma}]_{\alpha_{P_{s,a}}^{\lambda, \omega^*}} \right| + 2\varepsilon_{\text{opt}} \end{aligned} \quad (182)$$

1030 where (a) comes from the triangle inequality, and (b) comes from $\|P_{s,a}^0 - \widehat{P}_{s,a}^0\|_1 \leq 2$ and
 1031 $\|[\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha_{P_{s,a}}^{\lambda, \omega^*}} - [\widehat{V}^{*, \sigma}]_{\alpha_{P_{s,a}}^{\lambda, \omega^*}}\|_\infty \leq \|\widehat{V}^{\widehat{\pi}, \sigma} - \widehat{V}^{*, \sigma}\|_\infty$, and (c) follows from the definition of the
 1032 optimization error in (55). The second term of the max can be controlled in the same manner, i.e.:

$$\left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha_{\widehat{P}_{s,a}}^{\lambda, \omega^*}} \right| \leq \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{*, \sigma}]_{\alpha_{\widehat{P}_{s,a}}^{\lambda, \omega^*}} \right| + 2\varepsilon_{\text{opt}} \quad (183)$$

$$\leq \max_{\mu_{P_{s,a}^0}^\lambda \in \mathcal{M}_{P_{s,a}^0}^\lambda} \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) (\widehat{V}^{*, \sigma} - \mu_{P_{s,a}^0}^\lambda) + \max_{\mu_{\widehat{P}_{s,a}^0}^\lambda \in \mathcal{M}_{\widehat{P}_{s,a}^0}^\lambda} \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) (\mu_{P_{s,a}^0}^\lambda - \mu_{\widehat{P}_{s,a}^0}^\lambda) \quad (184)$$

$$+ 2\varepsilon_{\text{opt}} \quad (185)$$

1033 where the last inequality follow the decomposition of (147). Finally, to control the remaining term

$$\max_{\mu_{P_{s,a}^0}^\lambda \in \mathcal{M}_{P_{s,a}^0}^\lambda} \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) \left(\widehat{V}^{*,\sigma} - \mu_{P_{s,a}^0}^\lambda \right) = \max_{\alpha_P^\lambda \in \mathcal{A}_P^\lambda} \left\{ \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_{\alpha_P^\lambda} \right\} \quad (186)$$

1034 (185) for any given $\alpha \in [0, \alpha_{P_{s,a}^0}^{\lambda, \omega^*} [C [0, \frac{1}{1-\gamma}]^S$ in the variational family with one parameter λ , with
 1035 the dependency between $\widehat{V}^{*,\sigma}$ and \widehat{P}^0 , we resort to the following leave-one-out argument or absorbing
 1036 MDPs used in [Agarwal et al., 2020, Li et al., 2022b, Shi and Chi, 2022, Clavier et al., 2023]. To
 1037 begin, we create a collection of auxiliary RMDPs that exhibit the intended statistical independence
 1038 between robust value functions and the estimated nominal transition kernel. These auxiliary RMDPs
 1039 are designed to be minimally distinct from the initial RMDPs, subsequently, we manage to control
 1040 the relevant term within these auxiliary RMDPs and demonstrate that its value closely approximates
 1041 the target quantity for the desired RMDP. Recall that the empirical infinite-horizon robust MDP $\widehat{\mathcal{M}}_{\text{rob}}$
 1042 is defined using the nominal transition kernel \widehat{P}^0 . Inspired by Agarwal et al. [2020], we can construct
 1043 an auxiliary absorbing robust MDP $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ for each state s and any non-negative scalar $u \geq 0$, so
 1044 that it is the same as $\widehat{\mathcal{M}}_{\text{rob}}$ except for the transition properties in state s . These auxiliary MDPS are
 1045 called absorbing MDPs are have been used for the first time in the context of RMDPS in Clavier et al.
 1046 [2023]. Defining the reward function and nominal transition kernel of $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ as $P^{s,u}$ and $r^{s,u}$, which
 1047 are expressed as follows using the same notation as Shi et al. [2023]:

$$\begin{cases} r^{s,u}(s, a) = u & \forall a \in \mathcal{A}, \\ r^{s,u}(\tilde{s}, a) = r(\tilde{s}, a) & \forall (\tilde{s}, a) \in \mathcal{S} \times \mathcal{A} \text{ and } \tilde{s} \neq s. \end{cases} \quad (187)$$

1048

$$\begin{cases} P^{s,u}(s' | s, a) = \mathbb{1}(s' = s) & \forall (s', a) \in \mathcal{S} \times \mathcal{A}, \\ P^{s,u}(\cdot | \tilde{s}, a) = \widehat{P}^0(\cdot | \tilde{s}, a) & \forall (\tilde{s}, a) \in \mathcal{S} \times \mathcal{A} \text{ and } \tilde{s} \neq s, \end{cases} \quad (188)$$

1049 Nominal transition probability at state s of the auxiliary $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ never leaves state s once entered,
 1050 which gives the name absorbing to these auxiliary RMPDs. Finally, we define the robust Bellman
 1051 operator $\widehat{\mathcal{T}}_{s,u}^\sigma(\cdot)$ associated $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ as

$$\widehat{\mathcal{T}}_{s,u}^\sigma(Q)(\tilde{s}, a) = r^{s,u}(\tilde{s}, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(P_{\tilde{s}, a}^{s,u})} \mathcal{P}V, \quad \text{with } V(\tilde{s}) = \max_a Q(\tilde{s}, a). \quad (189)$$

1052 in sa -rectangular case and with stochastic policy in s -rectangular case. Using these auxiliary RMDPs
 1053 we can remark equivalence between $\widehat{\mathcal{M}}_{\text{rob}}$ and the auxiliary RMDP $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ fixed-point. First, $\widehat{Q}^{*,\sigma}$
 1054 is the unique-fixed point of $\widehat{\mathcal{T}}^\sigma(\cdot)$ with associated value $\widehat{V}^{*,\sigma}$. We will show that the robust value
 1055 function $\widehat{V}_{s,u}^{*,\sigma}$ obtained from the fixed point of $\widehat{\mathcal{T}}_{s,u}^\sigma(\cdot)$ is the same as the the robust value function
 1056 $\widehat{V}^{*,\sigma}$ derived from $\widehat{\mathcal{T}}^\sigma(\cdot)$, as long as we choose u as

$$u^* := u^*(s) = \widehat{V}^{*,\sigma}(s) - \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(e_s)} \mathcal{P}\widehat{V}^{*,\sigma}. \quad (190)$$

1057 with e_s is the s -th standard basis vector in \mathbb{R}^S . This assertion is verified as:

1058 • **First for state $s' \neq s$, for all $a \in \mathcal{A}$:** it holds

$$\begin{aligned} r^{s,u^*}(s', a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(P_{s', a}^{s, u^*})} \mathcal{P}\widehat{V}^{*,\sigma} &= r(s', a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(\widehat{P}_{s', a}^0)} \mathcal{P}\widehat{V}^{*,\sigma} \\ &= \widehat{\mathcal{T}}^\sigma(\widehat{Q}^{*,\sigma})(s', a) = \widehat{Q}^{*,\sigma}(s', a), \end{aligned} \quad (191)$$

1059 where the first equality holds because of (187) and (188), and the last inequality comes
 1060 from that $\widehat{Q}^{*,\sigma}$ is the fixed point of $\widehat{\mathcal{T}}^\sigma(\cdot)$ (see Lemma 8.3) and the definition of the robust
 1061 Bellman operator in (13).

1062 • **Then for state s , for any $a \in \mathcal{A}$:**

$$\begin{aligned} r^{s,u^*}(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s, a}^{s, u^*})} \mathcal{P}\widehat{V}^{*,\sigma} &= u^* + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(e_s)} \mathcal{P}\widehat{V}^{*,\sigma} \\ &= \widehat{V}^{*,\sigma}(s) - \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(e_s)} \mathcal{P}\widehat{V}^{*,\sigma} + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\text{sa}, \sigma}(e_s)} \mathcal{P}\widehat{V}^{*,\sigma} = \widehat{V}^{*,\sigma}(s), \end{aligned} \quad (192)$$

1063 using in the first equality is the definition of $P_{s,a}^{s,u^*}$ in (188) and where we use the definition
 1064 of u^* in (190) in the second one.

1065 Finally, we have proved that there exists a fixed point $\widehat{Q}_{s,u^*}^{*,\sigma}$ of the operator $\widehat{\mathcal{T}}_{s,u^*}^\sigma(\cdot)$ by taking

$$\begin{cases} \widehat{Q}_{s,u^*}^{*,\sigma}(s, a) = \widehat{V}^{*,\sigma}(s) & \forall a \in \mathcal{A}, \\ \widehat{Q}_{s,u^*}^{*,\sigma}(s', a) = \widehat{Q}^{*,\sigma}(s', a) & \forall s' \neq s \text{ and } a \in \mathcal{A}. \end{cases} \quad (193)$$

1066 we have confirmed the existence of a fixed point of the operator $\widehat{\mathcal{T}}_{s,u^*}^\sigma(\cdot)$ with corresponding value
 1067 function $\widehat{V}_{s,u^*}^{*,\sigma}$ that coincide with $\widehat{V}^{*,\sigma}$. Note that the corresponding properties between $\widehat{\mathcal{M}}_{\text{rob}}$ and
 1068 $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ in Step 1 and Step 2 hold in fact for any uncertainty set and s - or sa -rectangular assumptions.
 1069 Equipped with these fixed point equalities, we can use concentration inequalities to show this lemma.

1070 **Concentration inequality using an ε -net for all reward values u .** First we can verify that

$$0 \leq u^* \leq [\widehat{V}^{*,\sigma}(s)]_{\alpha_{P_{s,a}}^{\lambda,\omega^*}} \leq \widehat{V}^{*,\sigma}(s) \leq \frac{1}{1-\gamma}. \quad (194)$$

1071 We first construct a N_{ε_2} -net over the interval $[0, 1/(1-\gamma)]$, where $|N_{\varepsilon_2}|$ the size of the net can be
 1072 controlled by $|N_{\varepsilon_2}| \leq \frac{3}{\varepsilon_2(1-\gamma)}$ [Vershynin, 2018]. The only parameter that vary is λ in the variation
 1073 family $\alpha_{P_{sa}}^\lambda$ so we have 1-dimensional control and not a vector in \mathbb{R}^S . Then similarly to Lemma 8.3,
 1074 it holds that for each $u \in N_{\varepsilon_2}$, there exists a unique fixed point $\widehat{Q}_{s,u}^{*,\sigma}$ of the operator $\widehat{\mathcal{T}}_{s,u}^\sigma(\cdot)$, which
 1075 satisfies $0 \leq \widehat{Q}_{s,u}^{*,\sigma} \leq \frac{1}{1-\gamma} \cdot 1$. Consequently, the corresponding robust value function can be upper
 1076 bounded by $\|\widehat{V}_{s,u}^{*,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$. Using (188) and (187) by construction for all $u \in N_{\varepsilon_2}$, $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ is
 1077 statistically independent of $\widehat{P}_{s,a}^0$. This independence indicates that $[\widehat{V}_{s,u}^{*,\sigma}]_\alpha$ and $\widehat{P}_{s,a}^0$ are independent
 1078 for a fixed α . Using (145) and (146) and taking the union bound over all $(s, a, \alpha) \in \mathcal{S} \times \mathcal{A} \times N_{\varepsilon_1}$,
 1079 $u \in N_{\varepsilon_2}$ gives that, with probability at least $1 - \delta$, it holds for all $(s, a, u) \in \mathcal{S} \times \mathcal{A} \times N_{\varepsilon_2}$ that

$$\max_{\substack{\alpha^{\lambda,\omega} \in \mathcal{A}_{P_{sa}}^{\lambda,\omega}}} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}_{s,u}^{*,\sigma}]_{\alpha_{P_{sa}}^{\lambda,\omega^*}} \right| \leq 2 \sqrt{\frac{2 \log\left(\frac{18\|1\|_* S A N |N_{\varepsilon_2}|}{\delta}\right)}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}_{s,u}^{*,\sigma})} \quad (195)$$

$$\begin{aligned} &+ \varepsilon_2 \\ &\leq 2 \sqrt{\frac{2 \log\left(\frac{18\|1\|_* S A N |N_{\varepsilon_2}|}{\delta}\right)}{(1-\gamma)^2 N}} + \varepsilon_2, \end{aligned} \quad (196)$$

1080 Finally, we use **uniform concentration** to obtain the lemma. Recalling that $u^* \in [0, \frac{1}{1-\gamma}]$ (see
 1081 (194)), we can always find some $\bar{u} \in N_{\varepsilon_2}$ such that $|\bar{u} - u^*| \leq \varepsilon_2$. Consequently, plugging in the
 1082 operator $\widehat{\mathcal{T}}_{s,\bar{u}}^\sigma(\cdot)$ in (189) yields

$$\forall Q \in \mathbb{R}^{SA} : \quad \left\| \widehat{\mathcal{T}}_{s,\bar{u}}^\sigma(Q) - \widehat{\mathcal{T}}_{s,u^*}^\sigma(Q) \right\|_\infty = |\bar{u} - u^*| \leq \varepsilon_2$$

1083 We can then remark that the fixed points of $\widehat{\mathcal{T}}_{s,\bar{u}}^\sigma(\cdot)$ and $\widehat{\mathcal{T}}_{s,u^*}^\sigma(\cdot)$ obey

$$\begin{aligned} \left\| \widehat{Q}_{s,\bar{u}}^{*,\sigma} - \widehat{Q}_{s,u^*}^{*,\sigma} \right\|_\infty &= \left\| \widehat{\mathcal{T}}_{s,\bar{u}}^\sigma(\widehat{Q}_{s,\bar{u}}^{*,\sigma}) - \widehat{\mathcal{T}}_{s,u^*}^\sigma(\widehat{Q}_{s,u^*}^{*,\sigma}) \right\|_\infty \\ &\leq \left\| \widehat{\mathcal{T}}_{s,\bar{u}}^\sigma(\widehat{Q}_{s,\bar{u}}^{*,\sigma}) - \widehat{\mathcal{T}}_{s,\bar{u}}^\sigma(\widehat{Q}_{s,u^*}^{*,\sigma}) \right\|_\infty + \left\| \widehat{\mathcal{T}}_{s,\bar{u}}^\sigma(\widehat{Q}_{s,u^*}^{*,\sigma}) - \widehat{\mathcal{T}}_{s,u^*}^\sigma(\widehat{Q}_{s,u^*}^{*,\sigma}) \right\|_\infty \\ &\leq \gamma \left\| \widehat{Q}_{s,\bar{u}}^{*,\sigma} - \widehat{Q}_{s,u^*}^{*,\sigma} \right\|_\infty + \varepsilon_2, \end{aligned}$$

1084 where we use that the operator $\widehat{\mathcal{T}}_{s,u}^\sigma(\cdot)$ is a γ -contraction. It gives that:

$$\left\| \widehat{Q}_{s,\bar{u}}^{*,\sigma} - \widehat{Q}_{s,u^*}^{*,\sigma} \right\|_\infty \leq \frac{\varepsilon_2}{(1-\gamma)} \quad \text{and} \quad \left\| \widehat{V}_{s,\bar{u}}^{*,\sigma} - \widehat{V}_{s,u^*}^{*,\sigma} \right\|_\infty \leq \left\| \widehat{Q}_{s,\bar{u}}^{*,\sigma} - \widehat{Q}_{s,u^*}^{*,\sigma} \right\|_\infty \leq \frac{\varepsilon_2}{(1-\gamma)}. \quad (197)$$

1085 Finally to control the first term in (185), using the identity $\widehat{V}^{*,\sigma} = \widehat{V}_{s,u^*}^{*,\sigma}$ or fixed point relation
 1086 between the two RMPDS, established in previous step of the proof gives that: for all $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned}
 & \max_{\alpha_{P_{s,a}}^{\lambda,\omega} \in A_{P_{s,a}}^{\lambda,\omega}} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{*,\sigma}]_{\alpha_{P_{s,a}}^{\lambda,\omega}} \right| \\
 & \leq \max_{\alpha_{P_{s,a}}^{\lambda,\omega} \in A_{P_{s,a}}^{\lambda,\omega}} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{*,\sigma}]_{\alpha_{P_{s,a}}^{\lambda,\omega}} \right| \\
 & \stackrel{(a)}{\leq} \max_{\alpha_{P_{s,a}}^{\lambda,\omega} \in A_{P_{s,a}}^{\lambda,\omega}} \left\{ \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}_{s,\bar{u}}^{*,\sigma}]_{\alpha_{P_{s,a}}^{\lambda,\omega}} \right| + \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) \left([\widehat{V}_{s,\bar{u}}^{*,\sigma}]_{\alpha_{P_{s,a}}^{\lambda,\omega}} - [\widehat{V}_{s,u^*}^{*,\sigma}]_{\alpha_{P_{s,a}}^{\lambda,\omega}} \right) \right| \right\} \\
 & \stackrel{(b)}{\leq} \max_{\alpha_{P_{s,a}}^{\lambda,\omega} \in A_{P_{s,a}}^{\lambda,\omega}} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}_{s,\bar{u}}^{*,\sigma}]_{\alpha_{P_{s,a}}^{\lambda,\omega}} \right| + \frac{2\varepsilon_2}{(1-\gamma)} \\
 & \stackrel{(c)}{\leq} \frac{2\varepsilon_2}{(1-\gamma)} + \varepsilon_2 + 2\sqrt{\frac{2\log(\frac{18\|1\|_*S\mathcal{A}N|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}_{s,\bar{u}}^{*,\sigma})} + \frac{4\log(\frac{18\|1\|_*S\mathcal{A}N|N_{\varepsilon_2}|}{\delta})}{3N(1-\gamma)} \\
 & \leq \frac{3\varepsilon_2}{(1-\gamma)} + 2\sqrt{\frac{2\log(\frac{18\|1\|_*S\mathcal{A}N|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma})} + \frac{4\log(\frac{18\|1\|_*S\mathcal{A}N|N_{\varepsilon_2}|}{\delta})}{3N(1-\gamma)} \\
 & \quad + 2\sqrt{\frac{2\log(\frac{18\|1\|_*S\mathcal{A}N|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\left| \text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma}) - \text{Var}_{P_{s,a}^0}(\widehat{V}_{s,\bar{u}}^{*,\sigma}) \right|} \\
 & \stackrel{(d)}{\leq} \frac{3\varepsilon_2}{(1-\gamma)} + 2\sqrt{2\frac{\log(\frac{18\|1\|_*S\mathcal{A}N|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma})} + 2\sqrt{\frac{4\varepsilon_2\log(\frac{18\|1\|_*S\mathcal{A}N|N_{\varepsilon_2}|}{\delta})}{N(1-\gamma)^2}}
 \end{aligned} \tag{198}$$

$$\leq 2\sqrt{\frac{L''}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma})} + \frac{14\log(\frac{54\|1\|_*S\mathcal{A}N|N_{\varepsilon_2}|}{\delta})}{N(1-\gamma)} \tag{199}$$

$$\leq 16\sqrt{\frac{L''}{(1-\gamma)^2N}}, \tag{200}$$

1087 with $L'' = \log\left(\frac{54\|1\|_*S\mathcal{A}N^2}{(1-\gamma)\delta}\right)$ where (a) comes from triangular inequality, (b) is due (197), for any
 1088 $\alpha \in \mathbb{R}^S$

$$\begin{aligned}
 \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) \left([\widehat{V}_{s,\bar{u}}^{*,\sigma}]_{\alpha} - [\widehat{V}_{s,u^*}^{*,\sigma}]_{\alpha} \right) \right| & \leq \left\| P_{s,a}^0 - \widehat{P}_{s,a}^0 \right\|_1 \left\| [\widehat{V}_{s,\bar{u}}^{*,\sigma}]_{\alpha} - [\widehat{V}_{s,u^*}^{*,\sigma}]_{\alpha} \right\|_{\infty} \\
 & \leq 2 \left\| \widehat{V}_{s,\bar{u}}^{*,\sigma} - \widehat{V}_{s,u^*}^{*,\sigma} \right\|_{\infty} \leq \frac{2\varepsilon_2}{(1-\gamma)},
 \end{aligned} \tag{201}$$

1089 (c) follows from (195), (d) holds using Lemma 1 with (197). Here, the two last inequalities hold by
 1090 letting $\varepsilon_2 = \frac{2\log(\frac{18\|1\|_*S\mathcal{A}N|N_{\varepsilon_2}|}{\delta})}{N}$, which gives $|N_{\varepsilon_2}| \leq \frac{3}{\varepsilon_2(1-\gamma)} \leq \frac{3N}{1-\gamma}$, and the last inequality holds
 1091 by the fact $\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma}) \leq \|\widehat{V}^{*,\sigma}\|_{\infty} \leq \frac{1}{1-\gamma}$ and letting $N \geq 2\log\left(\frac{54\|1\|_*S\mathcal{A}N^2}{(1-\gamma)\delta}\right) = L''$.

1092 Rewriting (180), the first term of the max is controlled.

$$\max \left\{ \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\widehat{\pi},\sigma}]_{\alpha_{P_{s,a}}^{\lambda^*}} \right|, \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\widehat{\pi},\sigma}]_{\alpha_{\widehat{P}_{s,a}^0}^{\lambda^*}} \right| \right\}$$

1093 The second term can be controlled by the same term as the first one plus an additional term with

$$\begin{aligned}
 & \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\widehat{\pi},\sigma}]_{\alpha_{\widehat{P}_{s,a}^0}^{\lambda^*}} \right| \leq \\
 & \left| \max_{\mu_{P_{s,a}^0}^{\lambda} \in \mathcal{M}_{P_{s,a}^0}^{\lambda}} \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) (\widehat{V}^{*,\sigma} - \mu_{P_{s,a}^0}^{\lambda}) + \max_{\mu_{\widehat{P}_{s,a}^0}^{\lambda} \in \mathcal{M}_{\widehat{P}_{s,a}^0}^{\lambda}} \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) (\mu_{P_{s,a}^0}^{\lambda} - \mu_{\widehat{P}_{s,a}^0}^{\lambda}) \right|
 \end{aligned}$$

1094 and similarly to previous lemma in (153), the residual or term in the right in the previous equation
 1095 can be controlled with $\frac{L' C_S \|1\|_*}{N(1-\gamma)}$. Finally, putting (199) and (200) back into Equation (185) and using
 1096 Eq. (200) with probability at least $1 - \delta$ we obtain

$$\begin{aligned}
 \left| \widehat{P}_{s,a}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - P_{s,a}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right| &\leq \max_{\alpha_{P_{s,a}}^{\lambda, \omega} \in \mathcal{A}_{P_{s,a}}^{\lambda, \omega}} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\star, \sigma}]_{\alpha_{P_{s,a}}^{\lambda, \omega}} \right| + 2\varepsilon_{\text{opt}} \\
 &\leq 2\sqrt{\frac{L'}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{\star, \sigma})} + 2\varepsilon_{\text{opt}} + \frac{14L'' C_S \|1\|_*}{N(1-\gamma)} \\
 &\leq 2\sqrt{\frac{L''}{(1-\gamma)^2 N}} + 2\varepsilon_{\text{opt}} + \frac{14L'' C_S \|1\|_*}{N(1-\gamma)}, \tag{202}
 \end{aligned}$$

1097 $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$. Using matrix form we obtain finally:

$$\begin{aligned}
 \left| \widehat{\underline{P}}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - \underline{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right| &\leq 2\sqrt{\frac{L''}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{\star, \sigma})} 1 + 2\varepsilon_{\text{opt}} 1 \\
 &\leq 2\sqrt{\frac{L''}{(1-\gamma)^2 N}} 1 + 2\varepsilon_{\text{opt}} 1 + \frac{14L'' C_S \|1\|_*}{N(1-\gamma)} \tag{203}
 \end{aligned}$$

1098 The proof is similar in the s -rectangular case, factorising by $\pi(a|s)$, like in 8. Moreover, the proof
 1099 is similar for TV without the geometric term depending on C_S .

1100 9.3.6 Proof of Lemma 11

1101 We always use the same manner as in Appendix 9.3.4. Similarly to (161), it holds:

$$\left(I - \gamma \underline{P}^{\widehat{\pi}, \widehat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi}, \widehat{V}}}(\widehat{V}^{\widehat{\pi}, \sigma})} \leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left(\underline{P}^{\widehat{\pi}, \widehat{V}} \right)^t \text{Var}_{\underline{P}^{\widehat{\pi}, \widehat{V}}}(\widehat{V}^{\widehat{\pi}, \sigma})}. \tag{204}$$

1102 In order to upper bound $\text{Var}_{\underline{P}^{\widehat{\pi}, \widehat{V}}}(\widehat{V}^{\widehat{\pi}, \sigma})$, we define $V' := \widehat{V}^{\widehat{\pi}, \sigma} - \eta 1$ for any α^* solving a dual
 1103 optimization problem with $\eta \in \mathbb{R}$. Using as (168), it holds

$$\begin{aligned}
 \text{Var}_{\underline{P}^{\widehat{\pi}, \widehat{V}}}(\widehat{V}^{\widehat{\pi}, \sigma}) &\leq \underline{P}^{\widehat{\pi}, \widehat{V}} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} 1 + \frac{2}{\gamma} \|V'\|_{\infty} \left| \left(\widehat{\underline{P}}^{\widehat{\pi}, \widehat{V}} - \underline{P}^{\widehat{\pi}, \widehat{V}} \right) \widehat{V}^{\widehat{\pi}, \sigma} \right| \\
 &\leq \underline{P}^{\widehat{\pi}, \widehat{V}} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} 1 + \frac{2}{\gamma} \|V'\|_{\infty} \left(2\sqrt{\frac{L''}{(1-\gamma)^2 N}} + 2\varepsilon_{\text{opt}} + \frac{14L'' C_S \|1\|_*}{N(1-\gamma)} \right) 1, \tag{205}
 \end{aligned}$$

1104 where the last inequality makes use of Lemma 10. Plugging (205) back into (204) leads to

$$\begin{aligned}
(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}})^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} &\stackrel{(a)}{\leq} \sqrt{\frac{1}{1-\gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t (\underline{P}^{\hat{\pi}, \hat{V}})^t \left(\underline{P}^{\hat{\pi}, \hat{V}} (V' \circ V') - \frac{1}{\gamma} V' \circ V' \right) \right|} \\
&+ \sqrt{\frac{1}{(1-\gamma)^2 \gamma^2} \left(2\sqrt{\frac{L''}{(1-\gamma)^2 N}} + 2\varepsilon_{\text{opt}} + \frac{14L''C_S \|1\|_*}{N(1-\gamma)} \right) \|V'\|_{\infty}} \\
&\stackrel{(b)}{\leq} \sqrt{\frac{\|V'\|_{\infty}^2}{\gamma(1-\gamma)}} 1 + \sqrt{\frac{\left(2\sqrt{\frac{L''}{(1-\gamma)^2 N}} + 2\varepsilon_{\text{opt}} + \frac{14L''C_S \|1\|_*}{N(1-\gamma)} \right) \|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1 \\
&\stackrel{(c)}{\leq} \sqrt{\frac{\|V'\|_{\infty}^2}{\gamma(1-\gamma)}} 1 + 5\sqrt{\frac{\left(1 + \varepsilon_{\text{opt}} + \frac{L''C_S \|1\|_*}{N(1-\gamma)} \right) \|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1 \\
&\leq 6\sqrt{\frac{\left(1 + \varepsilon_{\text{opt}} + \frac{L''C_S \|1\|_*}{N(1-\gamma)} \right) \|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1, \tag{206} \\
&\leq 6\sqrt{\left(1 + \varepsilon_{\text{opt}} + \frac{L''C_S \|1\|_*}{N(1-\gamma)} \right) \frac{\|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1, \tag{207}
\end{aligned}$$

1105 where (a) is the same as (174), (b) holds by repeating the argument of (175), (c) follows by taking
1106 $N \geq \frac{L''}{(1-\gamma)^2}$ and then the last inequality holds by $\|V'\|_{\infty} \leq \|V^{*,\sigma}\|_{\infty} \leq \frac{1}{1-\gamma}$. Then taking the
1107 infimum over η in the right-hand side of the equation in the definition of V' and using $\text{sp}(\cdot)_{\infty} \leq \|\cdot\|_*$
1108 gives

$$(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}})^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} \leq 6\sqrt{\left(1 + \varepsilon_{\text{opt}} + \frac{L''C_S \|1\|_*}{N(1-\gamma)} \right) \frac{\text{sp}(V)_{\infty}}{(1-\gamma)^2 \gamma^2}} 1$$

1109 Finally, applying Lemma 5 with $P = \hat{P}^0$ and $\pi = \hat{\pi}$ yields

$$\text{sp}(\hat{V}^{\hat{\pi}, \sigma})_* \leq \frac{1}{\gamma \max\{1-\gamma, \gamma C_g \sigma\}}, \tag{208}$$

1110 for sa -rectangular or

$$\text{sp}(\hat{V}^{\hat{\pi}, \sigma})_* \leq \frac{1}{\gamma \max\{1-\gamma, \min_s \|\hat{\pi}\|_* \tilde{\sigma}\}}$$

1111 which can be inserted into (207) and gives in sa -rectangular case:

$$\begin{aligned}
(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}})^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} &\leq 6\sqrt{\frac{\left(1 + \varepsilon_{\text{opt}} + \frac{L''C_S \|1\|_*}{N(1-\gamma)} \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, \sigma\}}} 1 \\
&\leq 6\sqrt{\frac{\left(1 + \varepsilon_{\text{opt}} + \frac{L''C_S \|1\|_*}{N(1-\gamma)} \right)}{(1-\gamma)^3 \gamma^3}} 1.
\end{aligned}$$

1112 where first inequalities comes from that we can bound it Eq. left-hand side of equation (207) by
1113 $\|V'\|_{\infty} \leq \|V^{*,\sigma}\|_{\infty} \leq \frac{1}{1-\gamma}$. Proof for s -rectangular is similar, but requires adding an extra factor
1114 depending on the norm of the current policy and we have:

$$\begin{aligned}
(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}})^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} &\leq 6\sqrt{\frac{\left(1 + \varepsilon_{\text{opt}} + \frac{L''C_S \|1\|_*}{N(1-\gamma)} \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, C_g \tilde{\sigma} \min_s \|\hat{\pi}\|_{\infty}\}}} 1 \\
&\leq 6\sqrt{\frac{\left(1 + \varepsilon_{\text{opt}} + \frac{L''C_S \|1\|_*}{N(1-\gamma)} \right)}{(1-\gamma)^3 \gamma^2}} 1.
\end{aligned}$$

1115 **9.3.7 Proof of Lemma 7**

1116 Observing that each row of P_π belongs to $\Delta(S)$, it can be directly verified that each row of $(1 -$
 1117 $\gamma)(I - \gamma P_\pi)^{-1}$ falls into $\Delta(S)$. As a result,

$$\begin{aligned} (I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} &= \frac{1}{1-\gamma} (1-\gamma) (I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} \\ &\stackrel{(a)}{\leq} \frac{1}{1-\gamma} \sqrt{(1-\gamma) (I - \gamma P_\pi)^{-1} \text{Var}_{P_\pi}(V^{\pi,P})} \\ &= \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \text{Var}_{P_\pi}(V^{\pi,P})}, \end{aligned} \quad (209)$$

1118 where (a) is due to Jensen's inequality. Then for any $\eta \in \mathbb{R}^+$, $V' := V^{\pi,P} - \eta \mathbf{1}$ for any α solving a
 1119 dual optimization problem, we can upper bound $\text{Var}_{P_\pi}(V^{\pi,P})$:

$$\begin{aligned} \text{Var}_{P_\pi}(V^{\pi,P}) &\stackrel{(i)}{=} \text{Var}_{P_\pi}(V') = P_\pi(V' \circ V') - (P_\pi V') \circ (P_\pi V') \\ &\stackrel{(ii)}{\leq} P_\pi(V' \circ V') - \frac{1}{\gamma^2} (V' - r_\pi + (1-\gamma)\eta \mathbf{1}) \circ (V' - r_\pi + (1-\gamma)\eta \mathbf{1}) \\ &= P_\pi(V' \circ V') - \frac{1}{\gamma^2} V' \circ V' + \frac{2}{\gamma^2} V' \circ (r_\pi - (1-\gamma)\eta \mathbf{1}) - \frac{1}{\gamma^2} (r_\pi - (1-\gamma)\eta \mathbf{1}) \circ (r_\pi - (1-\gamma)\eta \mathbf{1}) \\ &\leq P_\pi(V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_\infty \mathbf{1} \leq P_\pi(V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_\infty \mathbf{1}, \end{aligned} \quad (210)$$

1120 where (i) holds by the fact that $\text{Var}_{P_\pi}(V^{\pi,P} - b \mathbf{1}) = \text{Var}_{P_\pi}([V^{\pi,P}])$ for any scalar b and $V^{\pi,P} \in \mathbb{R}^S$,
 1121 (ii) follows from $V' \leq r_\pi + \gamma P_\pi V^{\pi,P} - \eta \mathbf{1} = r_\pi - (1-\gamma)\eta \mathbf{1} + \gamma P_\pi V'$, and the last line arises
 1122 from $\frac{1}{\gamma^2} V' \circ V' \geq \frac{1}{\gamma} V' \circ V'$ and $\|r_\pi - (1-\gamma)\eta \mathbf{1}\|_\infty \leq 1$. for $\eta \in [0, 1/(1-\gamma)]$ [Plugging (210)
 1123 back to (209) leads to

$$\begin{aligned} (I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} &\leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \left(P_\pi(V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_\infty \mathbf{1} \right)} \\ &\stackrel{(i)}{\leq} \sqrt{\frac{1}{1-\gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \left(P_\pi(V' \circ V') - \frac{1}{\gamma} V' \circ V' \right) \right|} + \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \frac{2}{\gamma^2} \|V'\|_\infty \mathbf{1}} \\ &\leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\left| \left(\sum_{t=0}^{\infty} \gamma^t (P_\pi)^{t+1} - \sum_{t=0}^{\infty} \gamma^{t-1} (P_\pi)^t \right) (V' \circ V') \right|} + \sqrt{\frac{2\|V'\|_\infty \mathbf{1}}{\gamma^2(1-\gamma)^2}} \\ &\stackrel{(ii)}{\leq} \sqrt{\frac{\|V'\|_\infty^2 \mathbf{1}}{\gamma(1-\gamma)}} + \sqrt{\frac{2\|V'\|_\infty \mathbf{1}}{\gamma^2(1-\gamma)^2}} \\ &\leq \sqrt{\frac{8\|V'\|_\infty \mathbf{1}}{\gamma^2(1-\gamma)^2}}, \end{aligned} \quad (211)$$

$$(212)$$

where (i) holds by the triangle inequality, (ii) holds by following recursion, and the last inequality holds by $\|V'\|_\infty \leq \frac{1}{1-\gamma}$. Then taking the minimum over η in the right-hand side of the equation gives the result.

$$(I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} \leq \sqrt{\frac{8 \text{sp}(V^{\pi,P})_\infty}{\gamma^2(1-\gamma)^2}}$$

1124 However, we also $\|V'\|_\infty \leq \|V^{\pi,P}\|_\infty \leq \frac{1}{1-\gamma}$ in (211). So finally, the result is

$$(I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} \leq \sqrt{\frac{8}{\gamma^2(1-\gamma)^2} \min\{\text{sp}([V^{\pi,P})_\infty, \frac{1}{1-\gamma}\}\mathbf{1}}.$$

1125 10 Proof of Theorem 2

1126 In this section, we focus on the scenarios in the uncertainty sets are constructed with (s, a) -
 1127 rectangularity condition with some general norms. Towards this, we firstly observe that for the
 1128 two limiting cases ℓ_1 norm and ℓ_∞ norm, one has $\|p_1 - p_2\|_1 \leq 2$ and $\|p_1 - p_2\|_\infty \leq 1$ for any two
 1129 probability distribution $p_1, p_2 \in \mathbb{R}^S$. Namely, the accessible ranges of the uncertainty level σ for ℓ_1
 1130 norm and ℓ_∞ norm are $(0, 2]$ and $(0, 1]$, respectively. In addition, we have

$$\forall p_1, p_2 \in \mathbb{R}^S : \quad \|p_1 - p_2\|_\infty \leq \|p_1 - p_2\| \leq \|p_1 - p_2\|_1 \quad (213)$$

1131 for any norm $\|\cdot\|$. It indicates that the accessible range of the uncertainty level $\sigma_{\|\cdot\|}$ for any given
 1132 norm $\|\cdot\|$ is between $(0, \sigma_{\|\cdot\|}^{\max}]$, where $1 \leq \sigma_{\|\cdot\|}^{\max} \leq 2$.

1133 To continue, we specify the definition of the uncertainty set with sa -rectangularity condition with
 1134 some given general norm $\|\cdot\|$ as below: for any nominal transition kernel $P \in \mathbb{R}^{SA \times S}$,

$$\mathcal{U}_{\|\cdot\|}^\sigma(P) := \mathcal{U}_{\|\cdot\|}^\sigma(P) = \otimes \mathcal{U}_p^\sigma(P_{s,a}), \quad \mathcal{U}_{\|\cdot\|}^\sigma(P_{s,a}) := \left\{ P'_{s,a} \in \Delta(S) : \|P'_{s,a} - P_{s,a}\| \leq \sigma_{\|\cdot\|} \right\}. \quad (214)$$

1135 Then, we recall the assumption of the uncertainty radius $\sigma_{\|\cdot\|} \in (0, \sigma_{\|\cdot\|}^{\max}(1 - c_0)]$ with $0 < c_0 < 1$.

1136 Then, resorting to the same class of hard MDPs in [Shi et al., 2023, Section C.1], we can complete
 1137 the proof by directly following the same proof pipeline of Shi et al. [2023, Section C] by replacing σ
 1138 with $\sigma_{\|\cdot\|}^{\max} \sigma_{\|\cdot\|}$.

1139 11 Proof of Theorem 4

1140 Developing the lower bound for the cases with s -rectangular uncertainty set involves several new
 1141 challenges compared to that of (s, a) -rectangular cases. Specifically, the first challenge is that the
 1142 optimal policy can be stochastic and hard to be characterized with a closed form for the RMDPs with
 1143 a s -rectangular uncertainty set, rather than deterministic policies in (s, a) -rectangular cases. Such
 1144 richer and smoother class of optimal policies makes slightly changing the transition kernel generally
 1145 could only leads to a smoothly changed stochastic optimal policy instead of a completely different
 1146 one. Such reduced changing of optimal policy further gives smaller performance gap, thus challenges
 1147 of a tighter lower bound. Second, most of the hard instances in the literature are constructed as SA
 1148 states with a constant number of action spaces without loss of generality. While when it comes to
 1149 s -rectangular uncertainty set, the action space size becomes important and can't be assumed as a
 1150 constant anymore. So a new class of instances are required.

1151 To address these challenges, in this section, we construct a new set of hard RMDP instances for two
 1152 limiting cases: ℓ_1 norm and ℓ_∞ norm.

1153 11.1 Construction of the hard problem instances

1154 Before proceeding, we introduce two useful sets related to the state space and action space as below:

$$\mathcal{S} = \{0, 1, \dots, S\}, \quad \text{and} \quad \mathcal{A} = \{0, 1, \dots, A - 1\}.$$

1155 In this section, we construct a set of RMDPs termed as $\mathcal{M}_{\ell_\infty}$, which consists of $S(A - 1)$ components
 1156 including $S(A - 1)$ components, each associates with some different state-action pair. Specifically, it
 1157 is defined as

$$\mathcal{M}_{\ell_\infty} := \left\{ \mathcal{M}_\theta = (\mathcal{S}, \mathcal{A}, \mathcal{U}^\sigma(P^\theta), r, \gamma) \mid \theta \in \Theta = \{(i, j) : (i, j) \in \mathcal{S} \times \mathcal{A} \setminus \{0\}\} \right\}. \quad (215)$$

1158 We introduce the detailed definition of $\mathcal{M}_{\ell_\infty}$ by introducing several key components of it sequentially.
 1159 In particular, for any RMDP $\mathcal{M}_\theta \in \mathcal{M}_{\ell_\infty}$, the state space is of size $2S$, which includes two classes
 1160 of states $\mathcal{X} = \{x_0, x_1, \dots, x_{S-1}\}$ and $\mathcal{Y} = \{y_0, y_1, \dots, y_{S-1}\}$. The action space for each state is
 1161 \mathcal{A} of A possible actions. So we have totally $2S$ states and there is in total $2SA$ state-action pairs.

1162 Armed with the above definitions, we can first introduce the following nominal transition kernel: for
 1163 all $(s, a) \in \mathcal{X} \cup \mathcal{Y} \times \mathcal{A}$

$$P^{(0,0)}(s' | s, a) = \begin{cases} p\mathbb{1}(s' = y_i) + (1-p)\mathbb{1}(s' = x_i) & \text{if } s = x_i, a = 0, \quad \forall i \in \mathcal{S} \\ q\mathbb{1}(s' = y_i) + (1-q)\mathbb{1}(s' = x_i) & \text{if } s = x_i, a \neq 0, \quad \forall i \in \mathcal{S} \\ \mathbb{1}(s' = s) & \text{if } s \in \mathcal{Y} \end{cases} \quad (216)$$

1164 Here, p and q are set according to

$$0 \leq p \leq 1 \quad \text{and} \quad 0 \leq q = p - \Delta \quad (217)$$

1165 for some p and $\Delta > 0$ that will be introduced momentarily.

1166 Then we introduce the $S(A-1)$ components inside \mathcal{M}_∞ . Namely, for any $(i, j) \in \mathcal{S} \times \mathcal{A} \setminus \{0\}$, the
 1167 nominal transition kernel of $\mathcal{M}_{(i,j)}$ is specified as

$$P^{(i,j)}(s' | s, a) = \begin{cases} p\mathbb{1}(s' = y_i) + (1-p)\mathbb{1}(s' = x_i) & \text{if } s = x_i, a = j \\ q\mathbb{1}(s' = y_i) + (1-q)\mathbb{1}(s' = x_i) & \text{if } s = x_i \in \mathcal{X}, a = 0 \\ P^{(0,0)}(s' | s, a) & \text{otherwise} \end{cases} \quad (218)$$

1168 In words, the nominal transition kernel of each variant $\mathcal{M}_{(i,j)}$ only differs slightly from that of the
 1169 basic nominal transition kernel $P^{(0,0)}$ when $s = x_i$ and $a = \{0, j\}$, which makes all the components
 1170 inside $\mathcal{M}_{\ell_\infty}$ closed to each other.

1171 In addition, the reward function is defined as

$$\forall a \in \mathcal{A} : \quad r(s, a) = \begin{cases} 1 & \text{if } s \in \mathcal{Y} \\ 0 & \text{otherwise.} \end{cases} \quad (219)$$

1172 **Uncertainty set of the transition kernels.** Recall the following useful notation for any transition
 1173 probability P , i.e., the transition vector associated with some state s is denoted as:

$$P_s := P(\cdot, \cdot | s) \in \mathbb{R}^{1 \times SA}, \quad P_s^\theta := P^\theta(\cdot, \cdot | s) \in \mathbb{R}^{1 \times SA}. \quad (220)$$

1174 With this in hand, the uncertainty set (definition in (5)) with ℓ_∞ norm for any P^θ with $\theta \in \Theta$ can be
 1175 represented as:

$$\mathcal{U}_\infty^{s, \tilde{\sigma}}(P_s^\theta) := \mathcal{U}_{\|\cdot\|}^{s, \tilde{\sigma}}(P_s^\theta) = \left\{ P'_s \in \Delta(\mathcal{S})^{\mathcal{A}} : \|P'_s - P_s^\theta\| \leq \tilde{\sigma} = \sigma \|1\|_\infty = \sigma \right\}. \quad (221)$$

1176 So without loss of generality, we set the radius $\sigma \in (0, (1-c_0)]$ with $0 < c_0 < 1$. Before proceeding,
 1177 we observe that as the uncertainty set above is defined with respect to ℓ_∞ , it directly implies that for
 1178 each $(s, a) \in \mathcal{S} \times \mathcal{A}$, the uncertainty set is independent and can be decomposed as

$$\mathcal{U}_\infty^{s, \tilde{\sigma}}(P_s^\theta) = \otimes \mathcal{U}_{\|\cdot\|}^{s, \tilde{\sigma}}(P_{s,a}^\theta) = \left\{ P'_{s,a} \in \Delta(\mathcal{S}) : \|P'_{s,a} - P_{s,a}^\theta\| \leq \sigma \right\}. \quad (222)$$

1179 Notably, this indicates that using s -rectangular uncertainty set with ℓ_∞ norm as the divergence
 1180 function is analogous to the case of using (s, a) -rectangular uncertainty set with ℓ_∞ norm. As a
 1181 result, we follow the pipeline of the prior art Shi et al. [2023, Section C] which established the
 1182 minimax-optimal lower bound for (s, a) -rectangular RMDPs with TV distance, which is analogous
 1183 to the ℓ_∞ case. Towards this, we set p, q, Δ as the same as the ones in Shi et al. [2023, Section C.1],
 1184 where we recall the expressions of p, q, Δ for self-contained as below: taking $c_1 := \frac{c_0}{2}$,

$$p = (1 + c_1) \max\{1 - \gamma, \sigma\} \quad \text{and} \quad \Delta \leq c_1 \max\{1 - \gamma, \sigma\}, \quad (223)$$

1185 which ensure several facts:

$$0 \leq p \leq 1 \quad \text{and} \quad p \geq q \geq \max\{1 - \gamma, \sigma\}. \quad (224)$$

1186 **Value functions and optimal policies.** For each RMDP instance $\mathcal{M}_\theta \in \mathcal{M}_{\ell_\infty}$, with some abuse
 1187 of notation, we denote π_θ^* as the optimal policy. In addition, let $V_\theta^{\pi, \sigma}$ (resp. $V_\theta^{*, \sigma}$) represent the
 1188 corresponding robust value function of any policy π (resp. π_θ^*) with uncertainty level σ . Armed with
 1189 these notations, the following lemma shows some essential properties concerning the value functions
 1190 and optimal policies; the proof is postponed to Appendix 11.3.

1191 **Lemma 12.** Consider any $\mathcal{M}_\theta \in \mathcal{M}_{\ell_\infty}$ and any policy π , one has

$$\forall (i, j) \in \Theta : \quad V_{(i,j)}^{\pi, \sigma}(x_i) \leq \frac{\gamma(z_{(i,j)}^\pi - \sigma)}{(1-\gamma) \left(1 + \frac{\gamma(z_{(i,j)}^\pi - \sigma)}{1-\gamma(1-\sigma)}\right) (1-\gamma(1-\sigma))}, \quad (225)$$

1192 where $z_{(i,j)}^\pi$ is defined as

$$\forall (i, j) \in \Theta : \quad z_{(i,j)}^\pi := p\pi(j|x_i) + q[1 - \pi(j|x_i)]. \quad (226)$$

1193 In addition, the robust optimal value functions and the robust optimal policies satisfy

$$\forall (i, j) \in \Theta, s \in \mathcal{X} : \quad V_{(i,j)}^{*, \sigma}(s) = \frac{\gamma(p - \sigma)}{(1-\gamma) \left(1 + \frac{\gamma(p - \sigma)}{1-\gamma(1-\sigma)}\right) (1-\gamma(1-\sigma))} \quad (227)$$

1194 and

$$\pi_{(i,j)}^*(j|x_i) = 1 \quad \text{and} \quad \pi_{(i,j)}^*(0|s) = 1 \quad \forall s \in \mathcal{X} \setminus \{x_i\}. \quad (228)$$

1195 In words, this lemma shows that for any RMDP $\mathcal{M}_{(i,j)}$, the optimal policy on state x_i satisfies
 1196 $\pi_{(i,j)}^*(j|x_i) = 1$ and will focus on $a = 0$ for all other states $s \in \mathcal{X} \setminus \{x_i\}$.

1197 11.2 Establishing the minimax lower bound

1198 **Step 1: converting the goal to estimate (i, j) .** Now we are in position to derive the lower bound.
 1199 Recall the goal is to control the following quantity associated with any policy estimator $\hat{\pi}$ based on
 1200 the dataset with in total N_{all} samples:

$$\max_{(i,j) \in \Theta} \mathbb{P}_{(i,j)} \left\{ \max_{s \in \mathcal{X} \cup \mathcal{Y}} \left(V_{(i,j)}^{*, \sigma}(s) - V_{(i,j)}^{\hat{\pi}, \sigma}(s) \right) \right\} \geq \max_{(i,j) \in \Theta} \mathbb{P}_{(i,j)} \left\{ \max_{s \in \mathcal{X}} \left(V_{(i,j)}^{*, \sigma}(s) - V_{(i,j)}^{\hat{\pi}, \sigma}(s) \right) \right\}. \quad (229)$$

1201 To do so, we can invoke a key claim in Shi et al. [2023] here since our problem setting can be reduced
 1202 to the same one in Shi et al. [2023]: With $\varepsilon \leq \frac{c_1}{32(1-\gamma)}$, letting

$$\Delta = 32(1-\gamma) \max\{1-\gamma, \sigma\} \varepsilon \leq c_1 \max\{1-\gamma, \sigma\} \quad (230)$$

1203 which satisfies (223), it leads to that for any policy $\hat{\pi}$ and all $(i, j) \in \Theta$,

$$\begin{aligned} V_{(i,j)}^{*, \sigma}(x_i) - V_{(i,j)}^{\hat{\pi}, \sigma}(x_i) &\geq 2\varepsilon(1 - \hat{\pi}(j|x_i)), \\ \forall s \in \mathcal{X} \setminus \{x_i\} : \quad V_{(i,j)}^{*, \sigma}(s) - V_{(i,j)}^{\hat{\pi}, \sigma}(s) &\geq 2\varepsilon(1 - \hat{\pi}(0|s)). \end{aligned} \quad (231)$$

1204 Before continuing, we introduce a useful notation for the subset of Θ excluding the cases with state i
 1205 is selected:

$$\forall i \in \mathcal{S} : \quad \Theta_{-i} = \Theta \setminus \{(i', j) : i' = i, j \in \mathcal{A} \setminus \{0\}\}. \quad (232)$$

1206 Armed with the above facts and notations, we first suppose there exists a policy $\hat{\pi}$ such that for some
 1207 $(i, j) \in \Theta$,

$$\mathbb{P}_{(i,j)} \left\{ V_{(i,j)}^{*, \sigma}(x_i) - V_{(i,j)}^{\hat{\pi}, \sigma}(x_i) \leq \varepsilon \right\} \geq \frac{3}{4}. \quad (233)$$

1208 which in view of (231) indicates that we necessarily have $\hat{\pi}(j|x_i) \geq \frac{1}{4}$ with probability at least $\frac{3}{4}$.

1209 As a result, taking

$$j' = \arg \max_{a \in \mathcal{A}} \hat{\pi}(a | x_i), \quad (234)$$

1210 we are motivated to construct the following estimate of θ :

$$\hat{\theta} \begin{cases} = (i, j') & \text{if } j' > 0 \\ \in \mathcal{G}_{-w} & \text{if } j' = 0, \end{cases} \quad (235)$$

1211 which satisfies

$$\mathbb{P}_{(i,j)} \{\hat{\theta} = (i, j)\} \geq \mathbb{P}_{(i,j)} \{j' = j\} \geq \mathbb{P}_{(i,j)} \{\hat{\pi}(j | x_i) > \frac{1}{A}\} \geq \frac{3}{4}. \quad (236)$$

1212 **Step 2: developing the probability of error in testing multiple hypotheses.** Before proceeding,
 1213 we discuss the dataset consisting of in total N_{all} independent samples. Observing that each RMDP
 1214 inside the set $\mathcal{M}_{\ell_\infty}$ are constructed symmetrically associated with one pair of states (x_i, y_i) for all
 1215 $i \in \mathcal{S}$ and another action $j \in \mathcal{A} \times \{0\}$, respectively. Therefore, it is obvious that the dataset is
 1216 supposed to be generated uniformly on each (x_i, y_i, j) to maximize the information gain, leading to
 1217 $\frac{N_{\text{all}}}{S(A-1)}$ samples for any states-action (x_i, y_i, j) with $i \in \mathcal{S}, j \in \mathcal{A} \setminus \{0\}$.

1218 Then we are ready to turn to the hypothesis testing problem over $(i, j) \in \Theta$. Towards this, we
 1219 consider the minimax probability of error defined as follows:

$$p_e := \inf_{\phi} \max_{(i,j) \in \Theta} \{\mathbb{P}_{(i,j)}(\phi \neq (i, j))\}, \quad (237)$$

1220 where the infimum is taken over all possible tests ϕ constructed from the dataset introduced above.

1221 To continue, armed with the above dataset with N_{all} independent samples, we denote $\mu^{i,j}$
 1222 (resp. $\mu^{i,j}(s, a)$) as the distribution vector (resp. distribution) of each sample tuple (s, a, s') un-
 1223 der the nominal transition kernel $P^{(i,j)}$ associated with $\mathcal{M}_{(i,j)}$. With this in mind, combined with
 1224 Fano's inequality from Tsybakov [2009, Theorem 2.2] and the additivity of the KL divergence
 1225 (cf. Tsybakov [2009, Page 85]), we obtain

$$\begin{aligned} p_e &\geq 1 - N_{\text{all}} \frac{\max_{(i,j), (i',j') \in \Theta, (i,j) \neq (i',j')} \text{KL}(\mu^{i,j} | \mu^{i',j'}) + \log 2}{\log |\Theta|} \\ &\stackrel{(i)}{\geq} 1 - N_{\text{all}} \max_{(i,j), (i',j') \in \Theta, (i,j) \neq (i',j')} \text{KL}(\mu^{i,j} | \mu^{i',j'}) - \frac{1}{2} \\ &= \frac{1}{2} - N_{\text{all}} \max_{(i,j), (i',j') \in \Theta, (i,j) \neq (i',j')} \text{KL}(\mu^{i,j} | \mu^{i',j'}) \end{aligned} \quad (238)$$

1226 where (i) holds by $\log |\Theta| \geq 2 \log 2$ as long as $S(A-1)$ are large enough.

1227 Then following the same proof pipeline of Shi et al. [2023, Section C.2], we can arrive at

$$p_e \geq \frac{1}{2} - \frac{N_{\text{all}}}{S(A-1)} \frac{4096}{c_1} (1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2 \geq \frac{1}{4}, \quad (239)$$

1228 if the sample size is selected as

$$N_{\text{all}} \leq \frac{c_1 S(A-1)}{16396(1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2}. \quad (240)$$

1229 **Step 3: summing up the results together.** Finally, we suppose that there exists an estimator $\hat{\pi}$
 1230 such that

$$\max_{(i,j) \in \Theta} \mathbb{P}_{(i,j)} \left[\max_{s \in \mathcal{X} \cup \mathcal{Y}} \left(V_{(i,j)}^{*,\sigma}(s) - V_{(i,j)}^{\hat{\pi},\sigma}(s) \right) \geq \varepsilon \right] < \frac{1}{4}, \quad (241)$$

1231 then according to (229), we necessarily have

$$\forall s \in \mathcal{X} : \max_{(i,j) \in \Theta} \mathbb{P}_{(i,j)} \left[V_{(i,j)}^{*,\sigma}(s) - V_{(i,j)}^{\hat{\pi},\sigma}(s) \geq \varepsilon \right] < \frac{1}{4}, \quad (242)$$

1232 which indicates

$$\forall s \in \mathcal{X} : \max_{(i,j) \in \Theta} \mathbb{P}_{(i,j)} \left[V_{(i,j)}^{*,\sigma}(s) - V_{(i,j)}^{\hat{\pi},\sigma}(s) < \varepsilon \right] \geq \frac{3}{4}. \quad (243)$$

1233 As a consequence, (236) shows we must have

$$\forall (i,j) \in \Theta : \mathbb{P}_{(i,j)} \left[\hat{\theta} = (i,j) \right] \geq \frac{3}{4} \quad (244)$$

1234 to achieve (241). However, this would contract with (239) if the sample size condition in (240) is
1235 satisfied. Thus, we complete the proof.

1236 11.3 Proof of Lemma 12

1237 Without loss of generality, we first consider any $\mathcal{M}_{(i,j)}$ with $(i,j) \in \mathcal{S} \times \mathcal{A} \setminus \{0\}$. Following the
1238 same routine of Shi et al. [2023, Section C.3.1], we can verify that the order of the robust value
1239 function $V_{(i,j)}^{\pi,\sigma}$ over different states satisfies

$$\forall k \in \mathcal{S} : V_{(i,j)}^{\pi,\sigma}(x_k) \leq V_{(i,j)}^{\pi,\sigma}(y_k), \quad (245)$$

1240 which means the robust value function of the states inside \mathcal{X} are always not larger than the corre-
1241 sponding states inside \mathcal{Y} .

1242 Then we denote the minimum of the robust value function over states as below:

$$V_{(i,j),\min}^{\pi,\sigma} := \min_{s \in \mathcal{S}} V_{(i,j)}^{\pi,\sigma}(s). \quad (246)$$

1243 In the following arguments, we first take a moment to assume $V_{(i,j),\min}^{\pi,\sigma} = V_{(i,j)}^{\pi,\sigma}(x_i)$. With this in
1244 mind, we arrive at

$$V_{(i,j)}^{\pi,\sigma}(y_i) = 1 + \gamma(1 - \sigma) V_{(i,j)}^{\pi,\sigma}(y_i) + \gamma\sigma V_{(i,j),\min}^{\pi,\sigma} = \frac{1 + \gamma\sigma V_{(i,j)}^{\pi,\sigma}(x_i)}{1 - \gamma(1 - \sigma)}. \quad (247)$$

1245 Then, when we move on to the characterization of the robust value function at state x_i . To do so, we
1246 notice two important facts:

- 1247 1) The nominal transition probability $P_{x_i,a}^{(i,j)}$ at state-action pair (x_i, a) for any $a \in \mathcal{A}$ is a
1248 Bernoulli distribution (see (218) and (216)). The TV distance and the ℓ_∞ norm between
1249 two Bernoulli distribution are the same.
- 1250 2) Invoking the definitions of the nominal transition probability in (218) and (216), we have

$$\begin{aligned} P_{x_i,j}^{(i,j)} &= p\mathbb{1}(s' = y_i) + (1 - p)\mathbb{1}(s' = x_i) \\ P_{x_i,a}^{(i,j)} &= q\mathbb{1}(s' = y_i) + (1 - q)\mathbb{1}(s' = x_i) \quad \forall a \in \mathcal{A} \setminus \{j\}. \end{aligned} \quad (248)$$

1251 With the above two facts in hand, our problem setting is reduced to the same one in Shi et al. [2023]
1252 and can reuse the results in Shi et al. [2023, Section C.3.1] to achieve

$$V_{(i,j)}^{\pi,\sigma}(x_i) \leq \frac{\frac{\gamma(z_{(i,j)}^\pi - \sigma)}{1 - \gamma(1 - \sigma)}}{(1 - \gamma) \left(1 + \frac{\gamma(z_{(i,j)}^\pi - \sigma)}{1 - \gamma(1 - \sigma)} \right)}. \quad (249)$$

1253 and

$$\begin{aligned} \pi_{(i,j)}^*(j | x_i) &= 1 \\ V_{(i,j)}^{*,\sigma}(x_i) &= \frac{\frac{\gamma(z_{(i,j)}^{\pi^*} - \sigma)}{1 - \gamma(1 - \sigma)}}{(1 - \gamma) \left(1 + \frac{\gamma(z_{(i,j)}^{\pi^*} - \sigma)}{1 - \gamma(1 - \sigma)} \right)} = \frac{\frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)}}{(1 - \gamma) \left(1 + \frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)} \right)}. \end{aligned} \quad (250)$$

1254 Analogously, we can verify that for other $x_k \in \mathcal{X} \setminus \{x_i\}$,

$$\begin{aligned} \pi_{(i,j)}^*(0 | x_k) &= 1 \\ V_{(i,j)}^{*,\sigma}(x_k) &= \frac{\frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)}}{(1 - \gamma) \left(1 + \frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)} \right)}. \end{aligned} \quad (251)$$

1255 **12 DRVI for sa - rectangular algorithm for arbitrary norm**

1256 In order to compute the fixed point of $\widehat{\mathcal{T}}^\sigma$, distributionally robust value iteration (DRVI), is defined
 1257 in Algorithm 1. For sa -rectangularity, starting from an initialization $\widehat{Q}_0 = 0$, the update rule at the
 1258 t -th ($t \geq 1$) iteration is the following $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$:

$$\widehat{Q}_t^\pi(s, a) = \widehat{\mathcal{T}}^\sigma \widehat{Q}_{t-1}^\pi(s, a) = r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}_{\|\cdot\|}^{\sigma, \sigma}(\widehat{P}_{s,a}^0)} \mathcal{P} \widehat{V}_{t-1}, \quad (252)$$

1259 where $\widehat{V}_{t-1}(s) = \max_\pi \widehat{Q}_{t-1}^\pi(s, a)$ for all $s \in \mathcal{S}$.

1260 Directly solving (252) is computationally expensive since it involves optimization over a S -
 1261 dimensional probability simplex at each iteration, especially when the dimension of the state space \mathcal{S}
 1262 is large. Fortunately, given strong duality (252) can be equivalently solved using its dual problem,
 1263 which concerns optimizing a two variable (λ and ω) and thus can be solved efficiently. The specific
 1264 form of the dual problem depends on the choice of the norm $\|\cdot\|$, which we shall discuss separately in
 1265 Appendix 8.3. To complete the description, we output the greedy policy of the final Q-estimate \widehat{Q}_T
 1266 as the final policy $\widehat{\pi}$, namely,

$$\forall s \in \mathcal{S} : \quad \widehat{\pi}(s) = \arg \max_a \widehat{Q}_T(s, a). \quad (253)$$

1267 Encouragingly, the iterates $\{\widehat{Q}_t\}_{t \geq 0}$ of DRVI converge linearly to the fixed point $\widehat{Q}^{*, \sigma}$, owing to
 1268 the appealing γ -contraction property of $\widehat{\mathcal{T}}^\sigma$.

input: empirical nominal transition kernel \widehat{P}^0 ; reward function r ; uncertainty level σ ; number of iterations T .

initialization: $\widehat{Q}_0(s, a) = 0, \widehat{V}_0(s) = 0$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.

```

for  $t = 1, 2, \dots, T$  do
  for  $s \in \mathcal{S}, a \in \mathcal{A}$  do
    | Set  $\widehat{Q}_t(s, a)$  according to (252);
  end
  for  $s \in \mathcal{S}$  do
    | Set  $\widehat{V}_t(s) = \max_a \widehat{Q}_t(s, a)$ ;
  end
end

```

output: $\widehat{Q}_T, \widehat{V}_T$ and $\widehat{\pi}$ obeying $\widehat{\pi}(s) := \arg \max_a \widehat{Q}_T(s, a)$.

Algorithm 1: Distributionally robust value iteration (DRVI) for infinite-horizon RMDPs for sa -rectangular for arbitrary norm

1269 Using Algorithm 1, it allows getting an ϵ_{opt} error in the empirical MDP in the sa -rectangular case. In
 1270 the s -rectangular case, finding an algorithm to get ϵ_{opt} is more difficult to use, as the policy is not
 1271 deterministic anymore and 1 cannot anymore be applied. For L_p norms, Clavier et al. [2023] derived
 1272 an algorithm but for arbitrary norm we need to consider a more general problem for arbitrary norm in
 1273 Appendix 12

1274 **NeurIPS Paper Checklist**

1275 **1. Claims**

1276 Question: Do the main claims made in the abstract and introduction accurately reflect the
1277 paper's contributions and scope?

1278 Answer: [Yes]

1279 Justification: Yes

1280 Guidelines:

- 1281 • The answer NA means that the abstract and introduction do not include the claims
1282 made in the paper.
- 1283 • The abstract and/or introduction should clearly state the claims made, including the
1284 contributions made in the paper and important assumptions and limitations. A No or
1285 NA answer to this question will not be perceived well by the reviewers.
- 1286 • The claims made should match theoretical and experimental results, and reflect how
1287 much the results can be expected to generalize to other settings.
- 1288 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
1289 are not attained by the paper.

1290 **2. Limitations**

1291 Question: Does the paper discuss the limitations of the work performed by the authors?

1292 Answer: [Yes]

1293 Justification: See conclusion.

1294 Guidelines:

- 1295 • The answer NA means that the paper has no limitation while the answer No means that
1296 the paper has limitations, but those are not discussed in the paper.
- 1297 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1298 • The paper should point out any strong assumptions and how robust the results are to
1299 violations of these assumptions (e.g., independence assumptions, noiseless settings,
1300 model well-specification, asymptotic approximations only holding locally). The authors
1301 should reflect on how these assumptions might be violated in practice and what the
1302 implications would be.
- 1303 • The authors should reflect on the scope of the claims made, e.g., if the approach was
1304 only tested on a few datasets or with a few runs. In general, empirical results often
1305 depend on implicit assumptions, which should be articulated.
- 1306 • The authors should reflect on the factors that influence the performance of the approach.
1307 For example, a facial recognition algorithm may perform poorly when image resolution
1308 is low or images are taken in low lighting. Or a speech-to-text system might not be
1309 used reliably to provide closed captions for online lectures because it fails to handle
1310 technical jargon.
- 1311 • The authors should discuss the computational efficiency of the proposed algorithms
1312 and how they scale with dataset size.
- 1313 • If applicable, the authors should discuss possible limitations of their approach to
1314 address problems of privacy and fairness.
- 1315 • While the authors might fear that complete honesty about limitations might be used by
1316 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
1317 limitations that aren't acknowledged in the paper. The authors should use their best
1318 judgment and recognize that individual actions in favor of transparency play an impor-
1319 tant role in developing norms that preserve the integrity of the community. Reviewers
1320 will be specifically instructed to not penalize honesty concerning limitations.

1321 **3. Theory Assumptions and Proofs**

1322 Question: For each theoretical result, does the paper provide the full set of assumptions and
1323 a complete (and correct) proof?

1324 Answer: [Yes]

1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378

Justification: Assumptions are stated in lemmas and Theorems.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [NA]

Justification: Theoretical paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430

Answer: [NA]

Justification: Theoretical paper.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [NA]

Justification: Theoretical paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: Theoretical paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- 1431
- 1432
- 1433
- 1434
- 1435
- 1436
- 1437
- 1438
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
 - If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

1439 8. Experiments Compute Resources

1440 Question: For each experiment, does the paper provide sufficient information on the com-
1441 puter resources (type of compute workers, memory, time of execution) needed to reproduce
1442 the experiments?

1443 Answer: [NA]

1444 Justification: Theoretical paper.

1445 Guidelines:

- 1446
- 1447
- 1448
- 1449
- 1450
- 1451
- 1452
- 1453
- The answer NA means that the paper does not include experiments.
 - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
 - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
 - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

1454 9. Code Of Ethics

1455 Question: Does the research conducted in the paper conform, in every respect, with the
1456 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1457 Answer: [Yes]

1458 Justification: Done

1459 Guidelines:

- 1460
- 1461
- 1462
- 1463
- 1464
- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

1465 10. Broader Impacts

1466 Question: Does the paper discuss both potential positive societal impacts and negative
1467 societal impacts of the work performed?

1468 Answer: [NA]

1469 Justification: Theoretical paper.

1470 Guidelines:

- 1471
- 1472
- 1473
- 1474
- 1475
- 1476
- 1477
- 1478
- 1479
- 1480
- 1481
- The answer NA means that there is no societal impact of the work performed.
 - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to

1482 generate deepfakes for disinformation. On the other hand, it is not needed to point out
1483 that a generic algorithm for optimizing neural networks could enable people to train
1484 models that generate Deepfakes faster.

- 1485 • The authors should consider possible harms that could arise when the technology is
1486 being used as intended and functioning correctly, harms that could arise when the
1487 technology is being used as intended but gives incorrect results, and harms following
1488 from (intentional or unintentional) misuse of the technology.
- 1489 • If there are negative societal impacts, the authors could also discuss possible mitigation
1490 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1491 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1492 feedback over time, improving the efficiency and accessibility of ML).

1493 11. Safeguards

1494 Question: Does the paper describe safeguards that have been put in place for responsible
1495 release of data or models that have a high risk for misuse (e.g., pretrained language models,
1496 image generators, or scraped datasets)?

1497 Answer: [NA]

1498 Justification: Theoretical paper.

1499 Guidelines:

- 1500 • The answer NA means that the paper poses no such risks.
- 1501 • Released models that have a high risk for misuse or dual-use should be released with
1502 necessary safeguards to allow for controlled use of the model, for example by requiring
1503 that users adhere to usage guidelines or restrictions to access the model or implementing
1504 safety filters.
- 1505 • Datasets that have been scraped from the Internet could pose safety risks. The authors
1506 should describe how they avoided releasing unsafe images.
- 1507 • We recognize that providing effective safeguards is challenging, and many papers do
1508 not require this, but we encourage authors to take this into account and make a best
1509 faith effort.

1510 12. Licenses for existing assets

1511 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1512 the paper, properly credited and are the license and terms of use explicitly mentioned and
1513 properly respected?

1514 Answer: [NA]

1515 Justification: Theoretical paper.

1516 Guidelines:

- 1517 • The answer NA means that the paper does not use existing assets.
- 1518 • The authors should cite the original paper that produced the code package or dataset.
- 1519 • The authors should state which version of the asset is used and, if possible, include a
1520 URL.
- 1521 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1522 • For scraped data from a particular source (e.g., website), the copyright and terms of
1523 service of that source should be provided.
- 1524 • If assets are released, the license, copyright information, and terms of use in the package
1525 should be provided. For popular datasets, paperswithcode.com/datasets has
1526 curated licenses for some datasets. Their licensing guide can help determine the license
1527 of a dataset.
- 1528 • For existing datasets that are re-packaged, both the original license and the license of
1529 the derived asset (if it has changed) should be provided.
- 1530 • If this information is not available online, the authors are encouraged to reach out to
1531 the asset's creators.

1532 13. New Assets

1533 Question: Are new assets introduced in the paper well documented and is the documentation
1534 provided alongside the assets?

1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580

Answer: [NA]

Justification: Theoretical paper.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: Theoretical paper.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: Theoretical paper.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.