

Appendix

Table of Contents

A Related works	13
B Preliminaries	14
B.1 Properties of the robust Bellman operator	14
B.2 Concentration inequalities	15
B.3 Kullback-Leibler (KL) divergence	15
C Analysis: episodic finite-horizon RMDPs	16
C.1 Proof of the upper bound: Theorem 1	16
C.2 Proof of Lemma 8	20
C.3 Proof of the lower bound: Theorem 2	24
C.4 Proof of auxiliary results	28
D Problem formulation: discounted infinite-horizon RMDPs	34
D.1 Basics about discounted infinite-horizon MDPs	35
D.2 Distributionally robust discounted infinite-horizon MDPs	35
D.3 Distributionally robust offline RL	36
E Algorithm and theory: discounted infinite-horizon RMDPs	36
E.1 Building an empirical nominal MDP	36
E.2 Algorithm: DRVI-LCB for infinite-horizon RMDPs	37
E.3 Performance guarantees: infinite-horizon RMDPs	38
F Analysis: discounted infinite-horizon RMDPs	40
F.1 Proof of the upper bound: Theorem 3	40
F.2 Proof of the lower bound: Theorem 4	45
F.3 Proof of auxiliary lemmas	50

A RELATED WORKS

We shall focus on the closely related works on offline RL and distributionally robust RL.

Offline RL. Focusing on the task of learning an optimal policy from offline data, a significant amount of prior arts sets to understand the sample complexity and efficacy of offline RL under different assumptions of the history dataset. A bulk of prior results requires the history data to cover all the state-action pairs, under assumptions such as uniformly bounded concentrability coefficients (Chen & Jiang, 2019; Munos, 2005) and uniformly lower bounded data visitation distribution (Yin & Wang, 2021; Yin et al., 2021), where the latter assumption is also related to studies of asynchronous Q-learning (Li et al., 2021). More recently, the principle of pessimism has been investigated for offline RL in both model-based (Jin et al., 2021; Xie et al., 2021; Rashidinejad et al., 2021; Li et al., 2022) and model-free algorithms (Kumar et al., 2020; Shi et al., 2022; Yan et al., 2022), without the stringent requirement of full coverage. In particular, Li et al. (2022) established the near-minimax optimality of a pessimistic variant of value iteration under the single-policy clipped concentrability of history data, which inspired our algorithm design in the distributionally robust setting.

Distributionally robust RL. While distributionally robust optimization has been mainly investigated in the context of supervised learning (Rahimian & Mehrotra, 2019; Gao, 2020; Bertsimas et al., 2018; Duchi & Namkoong, 2018; Blanchet & Murthy, 2019), distributionally robust dynamic programming has also attracted considerable amount of attention, e.g. Iyengar (2005); Nilim & Ghaoui (2003); Xu & Mannor (2012); Nilim & El Ghaoui (2005), where natural robust extensions to the standard Bellman machineries are developed under mild assumptions. Targeting robust MDPs, empirical and theoretical works have been widely explored under different forms of uncertainty sets (Iyengar, 2005; Xu & Mannor, 2012; Wolff et al., 2012; Kaufman & Schaefer, 2013; Ho et al., 2018; Smirnova et al., 2019; Ho et al., 2021; Goyal & Grand-Clement, 2022; Derman & Mannor, 2020; Tamar et al., 2014; Badrinath & Kalathil, 2021). Nonetheless, the majority of prior theoretical analyses focus on planning with an exact knowledge of the uncertainty set (Iyengar, 2005; Xu & Mannor, 2012; Tamar et al., 2014), or are asymptotic in nature (Roy et al., 2017).

A number of robust RL algorithms were proposed recently with an emphasis on finite-sample performance guarantees under different data generating mechanisms. Wang & Zou (2021) proposed a robust Q-learning algorithm with an R-contamination uncertain set for the online setting, which achieves a similar bound as its non-robust counterpart. Badrinath & Kalathil (2021) proposed a model-free algorithm for the online setting with linear function approximation to cope with large state spaces. Yang et al. (2021); Panaganti & Kalathil (2022) developed sample complexities for a model-based robust RL algorithm with a variety of uncertainty sets where the data are collected using a generative model. In addition, Zhou et al. (2021) examined the uncertainty set defined by the KL divergence for offline data with uniformly lower bounded data visitation distribution. These works all require full coverage of the state-action space, whereas ours is the first one to leverage the principle of pessimism in robust offline RL.

B PRELIMINARIES

Before starting, let’s introduce some additional notation useful throughout the theoretical analysis. Let $\text{ess inf } X$ denote the essential infimum of a function/variable X .

B.1 PROPERTIES OF THE ROBUST BELLMAN OPERATOR

To begin with, we introduce the following strong duality lemma which is widely used in distributionally robust optimization when the uncertainty set is defined with respect to the KL divergence.

Lemma 2 ((Hu & Hong, 2013), Theorem 1). *Suppose $f(x)$ has a finite moment generating function in some neighborhood around $x = 0$, then for any $\sigma > 0$ and a nominal distribution P^0 , we have*

$$\sup_{\mathcal{P} \in \mathcal{U}^\sigma(P^0)} \mathbb{E}_{X \sim \mathcal{P}}[f(X)] = \inf_{\lambda \geq 0} \left\{ \lambda \log \mathbb{E}_{X \sim P^0} \left[\exp \left(\frac{f(X)}{\lambda} \right) \right] + \lambda \sigma \right\}. \quad (30)$$

Armed with the above lemma, it is easily verified that for any positive constant M and a nominal distribution vector $P^0 \in \mathbb{R}^{1 \times S}$ supported over the state space \mathcal{S} , if $X(s) \in [0, M]$ for all $s \in \mathcal{S}$, then

$$\inf_{\mathcal{P} \in \mathcal{U}^\sigma(P^0)} \mathcal{P}X = \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(P^0 \exp \left(-\frac{X}{\lambda} \right) \right) - \lambda \sigma \right\}. \quad (31)$$

For convenience, we introduce the following lemma, paraphrased from Zhou et al. (2021, Lemma 4) and its proof, to further characterize several essential properties of the optimal dual value.

Lemma 3 ((Zhou et al., 2021)). *Let $X \sim P$ be a bounded random variable with $X \in [0, M]$. Let $\sigma > 0$ be any uncertainty level and the corresponding optimal dual variable be*

$$\lambda^* \in \arg \max_{\lambda \geq 0} f(\lambda, P), \quad \text{where } f(\lambda, P) := \left\{ -\lambda \log \mathbb{E}_{X \sim P} \left[\exp \left(-\frac{X}{\lambda} \right) \right] - \lambda \sigma \right\}. \quad (32)$$

Then the optimal value λ^ obeys*

$$\lambda^* \in \left[0, \frac{M}{\sigma} \right], \quad (33)$$

where $\lambda^* = 0$ if and only if

$$\log(\mathbb{P}(X = \text{essinf } X)) + \sigma \geq 0. \quad (34)$$

Moreover, when $\lambda^* = 0$, we have

$$\lim_{\lambda \rightarrow 0} f(\lambda, P) = \lim_{\lambda \rightarrow 0} \left\{ -\lambda \log \mathbb{E}_{X \sim P} \left[\exp \left(\frac{-X}{\lambda} \right) \right] - \lambda \sigma \right\} = \text{essinf } X. \quad (35)$$

B.2 CONCENTRATION INEQUALITIES

In light of Lemma 3 (cf. 35), we are interested in comparing the values of $\text{essinf } X$ when X is drawn from the population nominal distribution or its empirical estimate. This is supplied by the following lemma from Zhou et al. (2021).

Lemma 4 ((Zhou et al., 2021)). *Let $X \sim P$ be a discrete bounded random variable with $X \in [0, M]$. Let P_n denote the empirical distribution constructed from n independent samples X_1, X_2, \dots, X_n , and let $\hat{X} \sim P_n$. Denote $P_{\min, X}$ as the smallest positive probability $P_{\min, X} := \min\{\mathbb{P}(X = x) : x \in \text{supp}(X)\}$, where $\text{supp}(X)$ is the support of X . Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, we have*

$$\min_{i \in [n]} X_i = \text{essinf } \hat{X} = \text{essinf } X, \quad (36)$$

as long as

$$n \geq -\frac{\log(2/\delta)}{\log(1 - P_{\min, X})}. \quad (37)$$

We next gather a few elementary facts about the Binomial distribution, which will be useful throughout the proof.

Lemma 5 (Chernoff's inequality). *Suppose $N \sim \text{Binomial}(n, p)$, where $n \geq 1$ and $p \in [0, 1]$. For some universal constant $c_f > 0$, we have*

$$\mathbb{P}(|N/n - p| \geq pt) \leq \exp(-c_f n p t^2), \quad \forall t \in [0, 1]. \quad (38)$$

Lemma 6 ((Shi et al., 2022, Lemma 8)). *Suppose $N \sim \text{Binomial}(n, p)$, where $n \geq 1$ and $p \in [0, 1]$. For any $\delta \in (0, 1)$, we have*

$$N \geq \frac{np}{8 \log(\frac{1}{\delta})} \quad \text{if } np \geq 8 \log\left(\frac{1}{\delta}\right), \quad (39a)$$

$$N \leq \begin{cases} e^2 np & \text{if } np \geq \log\left(\frac{1}{\delta}\right), \\ 2e^2 \log\left(\frac{1}{\delta}\right) & \text{if } np \leq 2 \log\left(\frac{1}{\delta}\right) \end{cases} \quad (39b)$$

hold with probability at least $1 - 4\delta$.

B.3 KULLBACK-LEIBLER (KL) DIVERGENCE

We next introduce some useful facts about the Kullback-Leibler (KL) divergence for two distributions P and Q , denoted as $\text{KL}(P \parallel Q)$. Denoting $\text{Ber}(p)$ (resp. $\text{Ber}(q)$) as the Bernoulli distribution with mean p (resp. q), we introduce

$$\text{KL}(\text{Ber}(p) \parallel \text{Ber}(q)) := p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}, \quad (40)$$

which represents the KL divergence from $\text{Ber}(p)$ to $\text{Ber}(q)$. We now introduce the following lemma.

Lemma 7. *For any $p, q \in [\frac{1}{2}, 1]$ and $p > q$, it holds that*

$$\text{KL}(\text{Ber}(p) \parallel \text{Ber}(q)) \leq \text{KL}(\text{Ber}(q) \parallel \text{Ber}(p)) \leq \frac{(p-q)^2}{p(1-p)}. \quad (41)$$

Moreover, for any $0 \leq x < y < q$, it holds

$$\text{KL}(\text{Ber}(x) \parallel \text{Ber}(q)) > \text{KL}(\text{Ber}(y) \parallel \text{Ber}(q)). \quad (42)$$

Proof. The first half of this lemma is proven in (Li et al. 2022 Lemma 10). For the latter half, it follows from that the function

$$f(x, q) := \text{KL}(\text{Ber}(x) \parallel \text{Ber}(q))$$

is monotonically decreasing for all $x \in (0, q]$, since its derivative with respect to x satisfies $\frac{\partial f(x, q)}{\partial x} = \log \frac{x}{q} + \log \frac{1-q}{1-x} < 0$. \square

C ANALYSIS: EPISODIC FINITE-HORIZON RMDPs

C.1 PROOF OF THE UPPER BOUND: THEOREM 1

In this section, we outline the proof of Theorem 1. Before starting, we introduce several additional notation that will be useful in the analysis. First, we denote the state-action space covered by the behavior policy π^b in the nominal model P^0 as

$$\mathcal{C}^b = \left\{ (h, s, a) : d_h^{b, P^0}(s, a) > 0 \right\}. \quad (43)$$

Moreover, we recall the definition in (22) and define a similar one based on the exact nominal model P^0 as

$$P_{\min, h}(s, a) := \min_{s'} \left\{ P_h^0(s' | s, a) : P_h^0(s' | s, a) > 0 \right\}. \quad (44)$$

Clearly, by comparing with the definitions (23) and (24), it holds that

$$P_{\min}^* = \min_{h, s} P_{\min, h}(s, \pi_h^*(s)), \quad P_{\min}^b = \min_{(h, s, a) \in \mathcal{C}^b} P_{\min, h}(s, a). \quad (45)$$

For any time step $h \in [H]$, we denote the set of possible state occupancy distributions associated with the optimal policy π^* in a model within the uncertainty set $P \in \mathcal{U}^\sigma(P^0)$ as

$$\mathcal{D}_h^* := \left\{ \left[d_h^{*, P}(s) \right]_{s \in \mathcal{S}} : P \in \mathcal{U}^\sigma(P^0) \right\} = \left\{ \left[d_h^{*, P}(s, \pi_h^*(s)) \right]_{s \in \mathcal{S}} : P \in \mathcal{U}^\sigma(P^0) \right\}, \quad (46)$$

where the second equality is due to the fact that π^* is chosen to be deterministic.

With these in place, the proof of Theorem 1 is separated into several key steps, as outlined below.

Step 1: establishing the pessimism property. To achieve this claim, we heavily count on the following lemma whose proof can be found in Appendix C.2.

Lemma 8. *Instate the assumptions in Theorem 1. Then for all $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$, consider any vector $V \in \mathbb{R}^S$ independent of $\hat{P}_{h, s, a}^0$ obeying $\|V\|_\infty \leq H$. With probability at least $1 - \delta$, one has*

$$\left| \inf_{P \in \mathcal{U}^\sigma(\hat{P}_{h, s, a}^0)} \mathcal{P}V - \inf_{P \in \mathcal{U}^\sigma(P_{h, s, a}^0)} \mathcal{P}V \right| \leq b_h(s, a) \quad (47)$$

with $b_h(s, a)$ given in (21). Moreover, for all $(h, s, a) \in \mathcal{C}^b$, with probability at least $1 - \delta$, one has

$$\frac{P_{\min, h}(s, a)}{8 \log(KHS/\delta)} \leq \hat{P}_{\min, h}(s, a) \leq e^2 P_{\min, h}(s, a). \quad (48)$$

Armed with the above lemma, with probability at least $1 - \delta$, we shall show the following relation holds

$$\forall (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H + 1] : \quad \hat{Q}_h(s, a) \leq Q_h^{\hat{\pi}, \sigma}(s, a), \quad \hat{V}_h(s) \leq V_h^{\hat{\pi}, \sigma}(s), \quad (49)$$

which means that \hat{Q}_h (resp. \hat{V}_h) is a pessimistic estimate of $Q_h^{\hat{\pi}, \sigma}$ (resp. $V_h^{\hat{\pi}, \sigma}$). Towards this, it is easily verified that the latter assertion concerning $V_h^{\hat{\pi}, \sigma}$ is implied by the former, since

$$\hat{V}_h(s) = \max_a \hat{Q}_h(s, a) \leq \max_a Q_h^{\hat{\pi}, \sigma}(s, a) = V_h^{\hat{\pi}, \sigma}(s). \quad (50)$$

Therefore, the remainder of this step focuses on verifying the former assertion in (49) by induction.

- To begin, the claim (49) holds at the base case when $h = H + 1$, by invoking the trivial fact $\hat{Q}_{H+1}(s, a) = Q_{H+1}^{\hat{\pi}, \sigma}(s, a) = 0$.
- Then, suppose that $\hat{Q}_{h+1}(s, a) \leq Q_{h+1}^{\hat{\pi}, \sigma}(s, a)$ holds for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ at some time step $h \in [H]$, it boils down to show $\hat{Q}_h(s, a) \leq Q_h^{\hat{\pi}, \sigma}(s, a)$.

By the update rule of $\hat{Q}_h(s, a)$ in Algorithm 1 (cf. line 7), the above relation holds immediately if $\hat{Q}_h(s, a) = 0$ since $\hat{Q}_h(s, a) = 0 \leq Q_h^{\hat{\pi}, \sigma}(s, a)$. Otherwise, $\hat{Q}_h(s, a)$ is updated via

$$\begin{aligned} \hat{Q}_h(s, a) &= r_h(s, a) + \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(\hat{P}_{h,s,a}^0 \cdot \exp \left(\frac{-\hat{V}_{h+1}}{\lambda} \right) \right) - \lambda \sigma \right\} - b_h(s, a) \\ &\stackrel{(i)}{=} r_h(s, a) + \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{h,s,a}^0)} \mathcal{P} \hat{V}_{h+1} - b_h(s, a) \\ &\leq r_h(s, a) + \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,a}^0)} \mathcal{P} \hat{V}_{h+1} - b_h(s, a) \end{aligned} \quad (51)$$

$$\begin{aligned} &+ \left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{h,s,a}^0)} \mathcal{P} \hat{V}_{h+1} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,a}^0)} \mathcal{P} \hat{V}_{h+1} \right| \\ &\stackrel{(ii)}{\leq} r_h(s, a) + \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,a}^0)} \mathcal{P} V_{h+1}^{\hat{\pi}, \sigma} + 0 \stackrel{(iii)}{=} Q_h^{\hat{\pi}, \sigma}(s, a), \end{aligned} \quad (52)$$

where (i) rewrites the update rule back to its primal form (cf. (18)), (ii) holds by applying (47) with the condition (27) satisfied and the induction hypothesis $\hat{V}_{h+1} \leq V_{h+1}^{\hat{\pi}, \sigma}$, and lastly, (iii) follows by the robust Bellman consistency equation (8).

Putting them together, we have verified the claim (49) by induction.

Step 2: bounding $V_h^{*, \sigma}(s) - V_h^{\hat{\pi}, \sigma}(s)$. With the pessimism property (49) in place, we observe that the following relation holds

$$0 \leq V_h^{*, \sigma}(s) - V_h^{\hat{\pi}, \sigma}(s) \leq V_h^{*, \sigma}(s) - \hat{V}_h(s) \leq Q_h^{*, \sigma}(s, \pi_h^*(s)) - \hat{Q}_h(s, \pi_h^*(s)), \quad (53)$$

where the last inequality follows from $\hat{Q}_h(s, \pi_h^*(s)) \leq \max_a \hat{Q}_h(s, a) = \hat{V}_h(s)$. Then, by the robust Bellman optimality equation in (9) and the primal version of the update rule (cf. (18))

$$\begin{aligned} Q_h^{*, \sigma}(s, \pi_h^*(s)) &= r_h(s, \pi_h^*(s)) + \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,\pi_h^*(s)}^0)} \mathcal{P} V_{h+1}^{*, \sigma}, \\ \hat{Q}_h(s, \pi_h^*(s)) &= r_h(s, \pi_h^*(s)) + \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{h,s,\pi_h^*(s)}^0)} \mathcal{P} \hat{V}_{h+1} - b_h(s, \pi_h^*(s)), \end{aligned}$$

we arrive at

$$\begin{aligned} V_h^{*, \sigma}(s) - \hat{V}_h(s) &\leq Q_h^{*, \sigma}(s, \pi_h^*(s)) - \hat{Q}_h(s, \pi_h^*(s)) \\ &= \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,\pi_h^*(s)}^0)} \mathcal{P} V_{h+1}^{*, \sigma} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{h,s,\pi_h^*(s)}^0)} \mathcal{P} \hat{V}_{h+1} + b_h(s, \pi_h^*(s)) \\ &\leq \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,\pi_h^*(s)}^0)} \mathcal{P} V_{h+1}^{*, \sigma} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,\pi_h^*(s)}^0)} \mathcal{P} \hat{V}_{h+1} \\ &\quad + \left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{h,s,\pi_h^*(s)}^0)} \mathcal{P} \hat{V}_{h+1} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,\pi_h^*(s)}^0)} \mathcal{P} \hat{V}_{h+1} \right| + b_h(s, \pi_h^*(s)) \\ &\stackrel{(i)}{\leq} \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,\pi_h^*(s)}^0)} \mathcal{P} V_{h+1}^{*, \sigma} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,\pi_h^*(s)}^0)} \mathcal{P} \hat{V}_{h+1} + 2b_h(s, \pi_h^*(s)) \\ &\stackrel{(ii)}{\leq} \hat{P}_{h,s,\pi_h^*(s)}^{\inf} (V_{h+1}^{*, \sigma} - \hat{V}_{h+1}) + 2b_h(s, \pi_h^*(s)), \end{aligned} \quad (54)$$

where (i) holds by applying Lemma 2 (cf. (47)) since \widehat{V}_{h+1} is independent of $P_{h,s,\pi_h^*}^0(s)$ by construction, and (ii) arises from introducing the notation

$$\widehat{P}_{h,s,\pi_h^*}^{\text{inf}} := \operatorname{argmin}_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,\pi_h^*}^0(s))} \mathcal{P}\widehat{V}_{h+1} \quad (55)$$

and consequently,

$$\inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,\pi_h^*}^0(s))} \mathcal{P}V_{h+1}^{\star,\sigma} \leq \widehat{P}_{h,s,\pi_h^*}^{\text{inf}} V_{h+1}^{\star,\sigma}, \quad \text{and} \quad \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,\pi_h^*}^0(s))} \mathcal{P}\widehat{V}_{h+1} = \widehat{P}_{h,s,\pi_h^*}^{\text{inf}} \widehat{V}_{h+1}.$$

To continue, let us introduce some additional notation for convenience. Define a sequence of matrices $\widehat{P}_h^{\text{inf}} \in \mathbb{R}^{S \times S}$ and vectors $b_h^* \in \mathbb{R}^S$ for $h \in [H]$, where their s -th rows (resp. entries) are given by

$$\left[\widehat{P}_h^{\text{inf}}\right]_{s,\cdot} = \widehat{P}_{h,s,\pi_h^*}^{\text{inf}}, \quad \text{and} \quad b_h^*(s) = b_h(s, \pi_h^*(s)). \quad (56)$$

Applying (54) recursively over the time steps $h, h+1, \dots, H$ using the above notation gives

$$\begin{aligned} 0 \leq V_h^{\star,\sigma} - \widehat{V}_h &\leq \widehat{P}_h^{\text{inf}} (V_{h+1}^{\star,\sigma} - \widehat{V}_{h+1}) + 2b_h^* \\ &\leq \widehat{P}_h^{\text{inf}} \widehat{P}_{h+1}^{\text{inf}} (V_{h+2}^{\star,\sigma} - \widehat{V}_{h+2}) + 2\widehat{P}_h^{\text{inf}} b_{h+1}^* + 2b_h^* \leq \dots \leq 2 \sum_{i=h}^H \left(\prod_{j=h}^{i-1} \widehat{P}_j^{\text{inf}} \right) b_i^*, \end{aligned} \quad (57)$$

where we let $\left(\prod_{j=i}^{i-1} \widehat{P}_j^{\text{inf}}\right) = I$ for convenience.

For any $d_h^* \in \mathcal{D}_h^*$ (cf. (46)), taking inner product with (57) leads to

$$\left\langle d_h^*, V_h^{\star,\sigma} - \widehat{V}_h \right\rangle \leq \left\langle d_h^*, 2 \sum_{i=h}^H \left(\prod_{j=h}^{i-1} \widehat{P}_j^{\text{inf}} \right) b_i^* \right\rangle = 2 \sum_{i=h}^H \langle d_i^*, b_i^* \rangle, \quad (58)$$

where

$$d_i^* := \left[(d_h^*)^\top \left(\prod_{j=h}^{i-1} \widehat{P}_j^{\text{inf}} \right) \right]^\top \in \mathcal{D}_i^* \quad (59)$$

by the definition of \mathcal{D}_i^* (cf. (46)) for all $i = h+1, \dots, H$.

Step 3: controlling $\langle d_i^*, b_i^* \rangle$ using concentrability. Since $\langle d_i^*, b_i^* \rangle = \sum_{s \in S} d_i^*(s) b_i^*(s)$, we shall divide the discussion in two different cases.

- For $s \in S$ where $\max_{P \in \mathcal{U}^\sigma(P^0)} d_i^{\star,P}(s, \pi_i^*(s)) = \max_{P \in \mathcal{U}^\sigma(P^0)} d_i^{\star,P}(s) = 0$, it follows from the definition (cf. (46)) that for any $d_i^* \in \mathcal{D}_i^*$, it satisfies that

$$d_i^*(s) = 0. \quad (60)$$

- For $s \in S$ where $\max_{P \in \mathcal{U}^\sigma(P^0)} d_i^{\star,P}(s, \pi_i^*(s)) = \max_{P \in \mathcal{U}^\sigma(P^0)} d_i^{\star,P}(s) > 0$, by the assumption in (12)

$$\max_{P \in \mathcal{U}^\sigma(P^0)} \frac{\min \{d_i^{\star,P}(s, \pi_i^*(s)), \frac{1}{S}\}}{d_i^{\text{b},P^0}(s, \pi_i^*(s))} = \max_{P \in \mathcal{U}^\sigma(P^0)} \frac{\min \{d_i^{\star,P}(s), \frac{1}{S}\}}{d_i^{\text{b},P^0}(s, \pi_i^*(s))} \leq C_{\text{rob}}^* < \infty,$$

it implies that

$$d_i^{\text{b},P^0}(s, \pi_i^*(s)) > 0 \quad \text{and} \quad (i, s, \pi_i^*(s)) \in \mathcal{C}^{\text{b}}. \quad (61)$$

Lemma 1 tells that with probability at least $1 - 8\delta$,

$$\begin{aligned} N_i(s, \pi_i^*(s)) &\geq \frac{K d_i^{b, P^0}(s, \pi_i^*(s))}{8} - 5 \sqrt{K d_i^{b, P^0}(s, \pi_i^*(s)) \log \frac{KH}{\delta}} \stackrel{(i)}{\geq} \frac{K d_i^{b, P^0}(s, \pi_i^*(s))}{16} \\ &\stackrel{(ii)}{\geq} \frac{K \max_{P \in \mathcal{U}^\sigma(P^0)} \min \left\{ d_i^{*, P}(s, \pi_i^*(s)), \frac{1}{S} \right\}}{16 C_{\text{rob}}^*} \geq \frac{K \min \left\{ d_i^*(s), \frac{1}{S} \right\}}{16 C_{\text{rob}}^*}, \end{aligned} \quad (62)$$

where (i) holds due to

$$K d_i^{b, P^0}(s, \pi_i^*(s)) \geq c_1 \frac{d_i^{b, P^0}(s, \pi_i^*(s)) \log(KHS/\delta)}{d_{\min}^b P_{\min}^b} \geq \frac{c_1 \log \frac{KH}{\delta}}{P_{\min}^b} \geq c_1 \log \frac{KH}{\delta} \quad (63)$$

for some sufficiently large c_1 , where the first inequality follows from Condition (27), the second inequality follows from

$$d_{\min}^b = \min_{h, s, a} \left\{ d_h^{b, P^0}(s, a) : d_h^{b, P^0}(s, a) > 0 \right\} \leq d_i^{b, P^0}(s, \pi_i^*(s)) \quad (64)$$

and the last inequality follows from $P_{\min}^b \leq 1$. In addition, (ii) follows from Assumption 1.

With this in place, we observe that the pessimistic penalty (see (21)) obeys

$$\begin{aligned} b_i^*(s) &\leq c_b \frac{H}{\sigma} \sqrt{\frac{\log(\frac{KHS}{\delta})}{\widehat{P}_{\min, i}(s, \pi_i^*(s)) N_i(s, \pi_i^*(s))}} \stackrel{(i)}{\leq} 4c_b \frac{H}{\sigma} \sqrt{\frac{\log^2(\frac{KHS}{\delta})}{P_{\min, i}(s, \pi_i^*(s)) N_i(s, \pi_i^*(s))}} \\ &\leq 16c_b \frac{H}{\sigma} \sqrt{\frac{C_{\text{rob}}^* \log^2 \frac{KHS}{\delta}}{P_{\min, i}(s, \pi_i^*(s)) K \min \left\{ d_i^*(s), \frac{1}{S} \right\}}}, \end{aligned} \quad (65)$$

where (i) holds by applying (48) in view of the fact that $(i, s, \pi_i^*(s)) \in \mathcal{C}^b$ by (61), and the last inequality holds by (62).

Combining the results in the above two cases leads to

$$\begin{aligned} \sum_{s \in \mathcal{S}} d_i^*(s) b_i^*(s) &\leq \sum_{s \in \mathcal{S}} 16 d_i^*(s) c_b \frac{H}{\sigma} \sqrt{\frac{C_{\text{rob}}^* \log^2 \frac{KHS}{\delta}}{P_{\min, i}(s, \pi_i^*(s)) K \min \left\{ d_i^*(s), \frac{1}{S} \right\}}} \\ &\stackrel{(i)}{\leq} 16c_b \frac{H}{\sigma} \sqrt{\sum_{s \in \mathcal{S}} d_i^*(s) \frac{C_{\text{rob}}^* \log^2 \frac{KHS}{\delta}}{P_{\min, i}(s, \pi_i^*(s)) K \min \left\{ d_i^*(s), \frac{1}{S} \right\}}} \sqrt{\sum_{s \in \mathcal{S}} d_i^*(s)} \\ &\leq 32c_b \frac{H}{\sigma} \sqrt{\frac{S C_{\text{rob}}^* \log^2 \frac{KHS}{\delta}}{P_{\min, i}(s, \pi_i^*(s)) K}}, \end{aligned} \quad (66)$$

where (i) follows from the Cauchy-Schwarz inequality and the last inequality hold by the trivial fact

$$\sum_{s \in \mathcal{S}} \frac{d_i^*(s)}{\min \left\{ d_i^*(s), \frac{1}{S} \right\}} \leq \sum_{s \in \mathcal{S}} d_i^*(s) \left(\frac{1}{d_i^*(s)} + \frac{1}{1/S} \right) = \sum_{s \in \mathcal{S}} 1 + \frac{1}{S} \sum_{s \in \mathcal{S}} d_i^*(s) \leq 2S. \quad (67)$$

Step 4: finishing up the proof. Then, inserting (66) back into (58) with $h = 1$ shows

$$\left\langle d_1^*, V_1^{\star, \sigma} - \widehat{V}_1 \right\rangle \leq 2 \sum_{i=1}^H \langle d_i^*, b_i^* \rangle \leq \sum_{i=1}^H 64c_b \frac{H}{\sigma} \sqrt{\frac{S C_{\text{rob}}^* \log^2 \frac{KH}{\delta}}{P_{\min, i}(s, \pi_i^*(s)) K}} \leq c_2 \frac{H^2}{\sigma} \sqrt{\frac{S C_{\text{rob}}^* \log^2 \frac{KH}{\delta}}{P_{\min}^* K}}, \quad (68)$$

where the last inequality holds by plugging in the relation $P_{\min}^* \leq P_{\min, i}(s, \pi_i^*(s))$ for $i = 1, \dots, H$ by the definition in (23) (see also (45)), and choosing c_2 to be large enough. The proof is completed.

C.2 PROOF OF LEMMA 8

To begin, we shall introduce the following fact that

$$\forall (h, s, a) \in \mathcal{C}^b : \quad N_h(s, a) \geq \frac{c_1 \log \frac{KHS}{\delta}}{16P_{\min,h}(s, a)} \geq -\frac{\log \frac{2KHS}{\delta}}{\log(1 - P_{\min,h}(s, a))}, \quad (69)$$

as long as Condition (27) holds. The proof is postponed to Appendix C.2.3. With this in mind, we shall first establish the simpler bound (48) and then move on to show (47).

C.2.1 PROOF OF (48)

To begin, recall that (69) is satisfied for all $(h, s, a) \in \mathcal{C}^b$. By Lemma 6 and the union bound, it holds that with probability at least $1 - \delta$ that for all $(h, s, a) \in \mathcal{C}^b$:

$$\forall s' \in \mathcal{S} : \quad P_h^0(s' | s, a) \geq \frac{\hat{P}_h^0(s' | s, a)}{e^2} \geq \frac{P_h^0(s' | s, a)}{8e^2 \log(\frac{KHS}{\delta})}. \quad (70)$$

To characterize the relation between $P_{\min,h}(s, a)$ and $\hat{P}_{\min,h}(s, a)$ for any $(h, s, a) \in \mathcal{C}^b$, we suppose—without loss of generality—that $P_{\min,h}(s, a) = P_h^0(s_1 | s, a)$ and $\hat{P}_{\min,h}(s, a) = \hat{P}_h^0(s_2 | s, a)$ for some $s_1, s_2 \in \mathcal{S}$. Then, it follows that

$$\begin{aligned} P_{\min,h}(s, a) &= P_h^0(s_1 | s, a) \stackrel{(i)}{\geq} \frac{\hat{P}_h^0(s_1 | s, a)}{e^2} \geq \frac{\hat{P}_{\min,h}(s, a)}{e^2} = \frac{\hat{P}_h^0(s_2 | s, a)}{e^2} \\ &\stackrel{(ii)}{\geq} \frac{P_h^0(s_2 | s, a)}{8e^2 \log(\frac{KHS}{\delta})} \geq \frac{P_{\min,h}(s, a)}{8e^2 \log(\frac{KHS}{\delta})}, \end{aligned}$$

where (i) and (ii) follow from (70).

C.2.2 PROOF OF (47)

The main goal of (47) is to control the gap between robust Bellman operations based on the nominal transition kernel $P_{h,s,a}^0$ and the estimated kernel $\hat{P}_{h,s,a}^0$ by the constructed penalty term. Towards this, first consider $(h, s, a) \notin \mathcal{C}^b$, which corresponds to the state-action pairs (s, a) that haven't been visited at step h by the behavior policy. In other words, $N_h(s, a) = 0$. In this case, (47) can be easily verified that

$$\left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{h,s,a}^0)} \mathcal{P}V - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,a}^0)} \mathcal{P}V \right| \stackrel{(i)}{=} \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,a}^0)} \mathcal{P}V \leq \|V\|_\infty \stackrel{(ii)}{\leq} H \stackrel{(iii)}{=} b_h(s, a), \quad (71)$$

where (i) follows from the fact $\hat{P}_{h,s,a}^0 = 0$ when $N_h(s, a) = 0$ (see (15)), (ii) arises from the assumption $\|V\|_\infty \leq H$, and (iii) holds by the definition of $b_h(s, a)$ in (21). Therefore, the remainder of the proof will focus on verifying (47) for $(h, s, a) \in \mathcal{C}^b$. Rewriting the term of interest via duality (cf. Lemma 2) yields

$$\begin{aligned} &\left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{h,s,a}^0)} \mathcal{P}V - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,a}^0)} \mathcal{P}V \right| \\ &= \left| \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(\hat{P}_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\} - \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(P_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\} \right|. \end{aligned} \quad (72)$$

Denoting

$$\hat{\lambda}_{h,s,a}^* := \arg \max_{\lambda \geq 0} \left\{ -\lambda \log \left(\hat{P}_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\}, \quad (73a)$$

$$\lambda_{h,s,a}^* := \arg \max_{\lambda \geq 0} \left\{ -\lambda \log \left(P_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\}, \quad (73b)$$

Lemma 3 (cf. (33)) then gives that

$$\lambda_{h,s,a}^* \in \left[0, \frac{H}{\sigma}\right], \quad \hat{\lambda}_{h,s,a}^* \in \left[0, \frac{H}{\sigma}\right], \quad (74)$$

due to $\|V\|_\infty \leq H$. We shall control (72) in three different cases separately: (a) $\lambda_{h,s,a}^* = 0$ and $\hat{\lambda}_{h,s,a}^* = 0$; (b) $\lambda_{h,s,a}^* > 0$ and $\hat{\lambda}_{h,s,a}^* = 0$ or $\lambda_{h,s,a}^* = 0$ and $\hat{\lambda}_{h,s,a}^* > 0$; and (c) $\lambda_{h,s,a}^* \neq 0$ or $\hat{\lambda}_{h,s,a}^* \neq 0$.

Case (a): $\lambda_{h,s,a}^* = 0$ and $\hat{\lambda}_{h,s,a}^* = 0$. Applying Lemma 3 and Lemma 4 to (72) gives that, with probability at least $1 - \frac{\delta}{KH}$,

$$\begin{aligned} \left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{h,s,a}^0)} \mathcal{P}V - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,a}^0)} \mathcal{P}V \right| &\stackrel{(i)}{=} \left| \text{essinf}_{s \sim \hat{P}_{h,s,a}^0} V(s) - \text{essinf}_{s \sim P_{h,s,a}^0} V(s) \right| \\ &\stackrel{(ii)}{=} \left| \text{essinf}_{s \sim P_{h,s,a}^0} V(s) - \text{essinf}_{s \sim P_{h,s,a}^0} V(s) \right| \\ &= 0 \leq b_h(s, a). \end{aligned} \quad (75)$$

where (i) holds by Lemma 3 (cf. (35)) and (ii) arises from Lemma 4 (cf. (36)) given (69).

Case (b): $\lambda_{h,s,a}^* > 0$ and $\hat{\lambda}_{h,s,a}^* = 0$ or $\lambda_{h,s,a}^* = 0$ and $\hat{\lambda}_{h,s,a}^* > 0$. Towards this, note that two trivial facts are implied by the definition (73):

$$\sup_{\lambda \geq 0} \left\{ -\lambda \log \left(P_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\} \geq -\hat{\lambda}_{h,s,a}^* \log \left(P_{h,s,a}^0 \cdot \exp \left(\frac{-V}{\hat{\lambda}_{h,s,a}^*} \right) \right) - \hat{\lambda}_{h,s,a}^* \sigma, \quad (76a)$$

$$\sup_{\lambda \geq 0} \left\{ -\lambda \log \left(\hat{P}_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\} \geq -\lambda_{h,s,a}^* \log \left(\hat{P}_{h,s,a}^0 \cdot \exp \left(\frac{-V}{\lambda_{h,s,a}^*} \right) \right) - \lambda_{h,s,a}^* \sigma. \quad (76b)$$

To continue, first, we consider a subcase when $\lambda_{h,s,a}^* = 0$ and $\hat{\lambda}_{h,s,a}^* > 0$. With probability at least $1 - \frac{\delta}{KH}$, it follows from Lemma 3 (cf. (35)) and Lemma 4 (cf. (36)) that

$$\begin{aligned} \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(\hat{P}_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\} &\geq \lim_{\lambda \rightarrow 0} \left\{ -\lambda \log \left(\hat{P}_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\} \\ &= \text{essinf}_{s \sim \hat{P}_{h,s,a}^0} V(s) = \text{essinf}_{s \sim P_{h,s,a}^0} V(s) \\ &= \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(P_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\}, \end{aligned} \quad (77)$$

leading to

$$\begin{aligned} &\left| \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(\hat{P}_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\} - \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(P_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\} \right| \\ &\stackrel{(i)}{\leq} \left(-\hat{\lambda}_{h,s,a}^* \log \left(\hat{P}_{h,s,a}^0 \cdot \exp \left(\frac{-V}{\hat{\lambda}_{h,s,a}^*} \right) \right) - \hat{\lambda}_{h,s,a}^* \sigma \right) \\ &\quad - \left(-\lambda_{h,s,a}^* \log \left(P_{h,s,a}^0 \cdot \exp \left(\frac{-V}{\lambda_{h,s,a}^*} \right) \right) - \lambda_{h,s,a}^* \sigma \right) \\ &\leq \hat{\lambda}_{h,s,a}^* \left| \log \left(\hat{P}_{h,s,a}^0 \cdot \exp \left(\frac{-V}{\hat{\lambda}_{h,s,a}^*} \right) \right) - \log \left(P_{h,s,a}^0 \cdot \exp \left(\frac{-V}{\lambda_{h,s,a}^*} \right) \right) \right|, \end{aligned} \quad (78)$$

where (i) follows from the definition of $\hat{\lambda}_{h,s,a}^*$ in (73) and the fact in (76a).

We pause to claim that with probability at least $1 - \delta$, the following bound holds

$$\forall (h, s, a) \in \mathcal{C}^b, V \in \mathbb{R}^S : \frac{\left| \left(\hat{P}_{h,s,a}^0 - P_{h,s,a}^0 \right) \cdot \exp \left(\frac{-V}{\lambda} \right) \right|}{P_{h,s,a}^0 \cdot \exp \left(\frac{-V}{\lambda} \right)} \leq \sqrt{\frac{\log \left(\frac{KHS}{\delta} \right)}{c_f N_h(s, a) P_{\min, h}(s, a)}} \leq \frac{1}{2}. \quad (79)$$

The proof is postponed to Appendix C.2.4. With (79) in place, we can further bound (78) (which is plugged into (72)) as

$$\begin{aligned} \left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{h,s,a}^0)} \mathcal{P}V - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,a}^0)} \mathcal{P}V \right| &\leq \hat{\lambda}_{h,s,a}^* \left| \log \left(1 + \frac{\left(\hat{P}_{h,s,a}^0 - P_{h,s,a}^0 \right) \cdot \exp \left(\frac{-V}{\lambda} \right)}{P_{h,s,a}^0 \cdot \exp \left(\frac{-V}{\lambda} \right)} \right) \right| \\ &\stackrel{(i)}{\leq} 2\hat{\lambda}_{h,s,a}^* \frac{\left| \left(\hat{P}_{h,s,a}^0 - P_{h,s,a}^0 \right) \cdot \exp \left(\frac{-V}{\lambda} \right) \right|}{P_{h,s,a}^0 \cdot \exp \left(\frac{-V}{\lambda} \right)} \\ &\stackrel{(ii)}{\leq} \frac{2H}{\sigma} \sqrt{\frac{\log \left(\frac{KHS}{\delta} \right)}{c_f N_h(s, a) P_{\min, h}(s, a)}} \\ &\leq \frac{2eH}{\sigma} \sqrt{\frac{\log \left(\frac{KHS}{\delta} \right)}{c_f N_h(s, a) \hat{P}_{\min, h}(s, a)}} \\ &\leq c_b \frac{H}{\sigma} \sqrt{\frac{\log \left(\frac{KHS}{\delta} \right)}{\hat{P}_{\min, h}(s, a) N_h(s, a)}}, \end{aligned} \quad (80)$$

where (i) follows from $\log(1+x) \leq 2|x|$ for any $|x| \leq \frac{1}{2}$ in view of (79), (ii) follows from (74) as well as (79), and the last line follows from (48) and choosing c_b to be sufficiently large.

Moreover, note that it can be easily verified that

$$\left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{h,s,a}^0)} \mathcal{P}V - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,a}^0)} \mathcal{P}V \right| \leq H$$

due to the assumption $\|V\|_\infty \leq H$. Plugging in the definition of $b_h(s, a)$ in (21), combined with the above bounds, we have that with probability at least $1 - \delta$,

$$\left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{h,s,a}^0)} \mathcal{P}V - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,a}^0)} \mathcal{P}V \right| \leq \min \left\{ c_b \frac{H}{\sigma} \sqrt{\frac{\log \left(\frac{KHS}{\delta} \right)}{N_h(s, a) \hat{P}_{\min, h}(s, a)}}, H \right\} =: b_h(s, a). \quad (81)$$

The other subcase when $\lambda_{h,s,a}^* > 0$ and $\hat{\lambda}_{h,s,a}^* = 0$ follows similarly from the bound

$$\begin{aligned} &\left| \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(\hat{P}_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\} - \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(P_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\} \right| \\ &\leq \lambda_{h,s,a}^* \left| \log \left(\hat{P}_{h,s,a}^0 \cdot \exp \left(\frac{-V}{\lambda_{h,s,a}^*} \right) \right) - \log \left(P_{h,s,a}^0 \cdot \exp \left(\frac{-V}{\lambda_{h,s,a}^*} \right) \right) \right|, \end{aligned} \quad (82)$$

and therefore, will be omitted for simplicity.

Case (c): $\lambda_{h,s,a}^* > 0$ and $\hat{\lambda}_{h,s,a}^* > 0$. It follows that

$$\begin{aligned} &\left| \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(\hat{P}_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\} - \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(P_{h,s,a}^0 \exp \left(\frac{-V}{\lambda} \right) \right) - \lambda \sigma \right\} \right| \\ &\stackrel{(i)}{\leq} \max \left\{ \left(-\hat{\lambda}_{h,s,a}^* \log \left(\hat{P}_{h,s,a}^0 \cdot e^{\frac{-V}{\hat{\lambda}_{h,s,a}^*}} \right) - \hat{\lambda}_{h,s,a}^* \sigma \right) \right. \\ &\quad \left. - \left(-\hat{\lambda}_{h,s,a}^* \log \left(P_{h,s,a}^0 \cdot e^{\frac{-V}{\hat{\lambda}_{h,s,a}^*}} \right) - \hat{\lambda}_{h,s,a}^* \sigma \right) \right\}, \end{aligned}$$

$$\begin{aligned}
& \left(-\lambda_{h,s,a}^* \log \left(P_{h,s,a}^0 \cdot e^{\frac{-V}{\lambda_{h,s,a}^*}} \right) - \lambda_{h,s,a}^* \sigma \right) - \left(-\lambda_{h,s,a}^* \log \left(\hat{P}_{h,s,a}^0 \cdot e^{\frac{-V}{\lambda_{h,s,a}^*}} \right) - \lambda_{h,s,a}^* \sigma \right) \Big\} \\
& \leq \max_{\lambda \in \{\lambda_{h,s,a}^*, \hat{\lambda}_{h,s,a}^*\}} \lambda \left| \log \left(\hat{P}_{h,s,a}^0 \cdot \exp \left(\frac{-V}{\lambda} \right) \right) - \log \left(P_{h,s,a}^0 \cdot \exp \left(\frac{-V}{\lambda} \right) \right) \right|, \quad (83)
\end{aligned}$$

where (i) can be verified by applying the facts in (76). Hence, the above term (83) can be controlled again in a similar manner as (78); we omit the details for simplicity.

Summing up. Combining the previous results in different cases by the union bound, with probability at least $1 - 10\delta$, it is satisfied that for all $(h, s, a) \in \mathcal{C}^b$:

$$\left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{h,s,a}^0)} \mathcal{P}V - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,s,a}^0)} \mathcal{P}V \right| \leq b_h(s, a),$$

which concludes the proof.

C.2.3 PROOF OF (69)

Observe that for all $(h, s, a) \in \mathcal{C}^b$:

$$K d_h^{b, P^0}(s, a) \stackrel{(i)}{\geq} \frac{c_1 d_h^{b, P^0}(s, a) \log(KHS/\delta)}{d_{\min}^b P_{\min}^b} \stackrel{(ii)}{\geq} \frac{c_1 \log(KHS/\delta)}{P_{\min}^b} \stackrel{(iii)}{\geq} \frac{c_1 \log(KHS/\delta)}{P_{\min, h}(s, a)}, \quad (84)$$

where (i) follows from Condition (27), (ii) follows from the definition that $d_{\min}^b \leq d_h^{b, P^0}(s, a)$ for $(h, s, a) \in \mathcal{C}^b$, and (iii) comes from (45).

Lemma 1 then tells that with probability at least $1 - 8\delta$,

$$\begin{aligned}
N_h(s, a) & \geq \frac{K d_h^{b, P^0}(s, a)}{8} - 5 \sqrt{K d_h^{b, P^0}(s, a) \log \frac{KH}{\delta}} \\
& \geq \frac{K d_i^{b, P^0}(s, a)}{16} \geq \frac{c_1 \log \frac{KH}{\delta}}{16 P_{\min, h}(s, a)}, \quad (85)
\end{aligned}$$

where the second line follows from the above relation as long as c_1 is sufficiently large. The last inequality of (69) then follows from

$$\frac{c_1 \log \frac{KHS}{\delta}}{16 P_{\min, h}(s, a)} \geq -\frac{\log \frac{2KHS}{\delta}}{\log(1 - P_{\min, h}(s, a))}, \quad (86)$$

since $x \leq -\log(1 - x)$ for all $x \in [0, 1]$.

C.2.4 PROOF OF (79)

Denoting

$$\text{supp}(P_{h,s,a}^0) := \{s' \in \mathcal{S} : P_h^0(s' | s, a) > 0\}$$

as the support of $P_{h,s,a}^0$, we observe that

$$\begin{aligned}
& \left| \frac{(\hat{P}_{h,s,a}^0 - P_{h,s,a}^0) \cdot \exp \left(\frac{-V}{\lambda} \right)}{P_{h,s,a}^0 \cdot \exp \left(\frac{-V}{\lambda} \right)} \right| \leq \frac{\sum_{s' \in \text{supp}(P_{h,s,a}^0)} |\hat{P}_h^0(s' | s, a) - P_h^0(s' | s, a)| \exp \left(\frac{-V(s')}{\lambda} \right)}{\sum_{s' \in \text{supp}(P_{h,s,a}^0)} P_h^0(s' | s, a) \exp \left(\frac{-V(s')}{\lambda} \right)} \\
& \leq \max_{s' \in \text{supp}(P_{h,s,a}^0)} \frac{|\hat{P}_h^0(s' | s, a) - P_h^0(s' | s, a)|}{P_h^0(s' | s, a)}, \quad (87)
\end{aligned}$$

where the second line follows from $\sum_i a_i = \sum_i b_i \frac{a_i}{b_i} \leq (\max_i \frac{a_i}{b_i}) \sum_i b_i$ for any positive sequences $\{a_i, b_i\}_i$ obeying $a_i, b_i > 0$.

To continue, note that for any $(h, s, a) \in \mathcal{C}^b$ and $s' \in \text{supp}(P_{h,s,a}^0)$, $N_h(s, a)\hat{P}_h^0(s' | s, a)$ follows the binomial distribution $\text{Binomial}(N_h(s, a), P_h^0(s' | s, a))$. Thus, applying Lemma 5 with $t =$

$\sqrt{\frac{\log(\frac{KHS}{\delta})}{c_f N_h(s, a) P_h^0(s' | s, a)}}$ yields

$$\mathbb{P}\left(\left|\hat{P}_h^0(s' | s, a) - P_h^0(s' | s, a)\right| \geq P_h^0(s' | s, a)t\right) \leq \exp(-c_f N_h(s, a) P_h^0(s' | s, a)t^2) \leq \frac{\delta}{KHS}, \quad (88)$$

as soon as $t \leq \frac{1}{2}$, which can be verified by the fact (69) and $P_{\min, h}(s, a) \leq P_h^0(s' | s, a)$ (cf. (44)), namely,

$$N_h(s, a) \geq \frac{c_1 \log \frac{KHS}{\delta}}{16 P_{\min, h}(s, a)} \geq \frac{\log(\frac{KHS}{\delta})}{4c_f P_{\min, h}(s, a)} \geq \frac{\log(\frac{KHS}{\delta})}{4c_f P_h^0(s' | s, a)} \quad (89)$$

as long as c_1 is sufficiently large.

Applying (88) and taking the union bound over $s \in \text{supp}(P_{h,s,a}^0)$ lead to that with probability at least $1 - \frac{\delta}{KH}$,

$$\begin{aligned} \max_{s' \in \text{supp}(P_{h,s,a}^0)} \frac{|\hat{P}_h^0(s' | s, a) - P_h^0(s' | s, a)|}{P_h^0(s' | s, a)} &\leq \max_{s' \in \text{supp}(P_{h,s,a}^0)} \frac{P_h^0(s' | s, a) \sqrt{\frac{\log(\frac{KHS}{\delta})}{c_f N_h(s, a) P_h^0(s' | s, a)}}}{P_h^0(s' | s, a)} \\ &= \max_{s' \in \text{supp}(P_{h,s,a}^0)} \sqrt{\frac{\log(\frac{KHS}{\delta})}{c_f N_h(s, a) P_h^0(s' | s, a)}} \\ &\leq \sqrt{\frac{\log(\frac{KHS}{\delta})}{c_f N_h(s, a) P_{\min, h}(s, a)}} \leq \frac{1}{2}, \end{aligned}$$

where the last line uses again (89). Plugging this back into (87) and applying the union bound over $(h, s, a) \in \mathcal{C}^b$ then completes the proof.

C.3 PROOF OF THE LOWER BOUND: THEOREM 2

The proof of Theorem 2 is inspired by the construction in Li et al. (2022) for standard MDPs, but is considerably more involved to handle the uncertainty set unique in robust MDPs. We shall first construct some hard instances and then characterize the sample complexity requirements over these instances.

C.3.1 CONSTRUCTION OF HARD PROBLEM INSTANCES

Construction of a collection of hard MDPs. Let us introduce two MDPs

$$\left\{ \mathcal{M}_\phi = \left(\mathcal{S}, \mathcal{A}, P^\phi = \{P_h^\phi\}_{h=1}^H, \{r_h\}_{h=1}^H, H \right) \mid \phi = \{0, 1\} \right\}, \quad (90)$$

where the state space is $\mathcal{S} = \{0, 1, \dots, S-1\}$, and the action space is $\mathcal{A} = \{0, 1\}$. The transition kernel P^ϕ of the constructed MDP \mathcal{M}_ϕ is defined as

$$P_1^\phi(s' | s, a) = \begin{cases} p\mathbb{1}(s' = 0) + (1-p)\mathbb{1}(s' = 1) & \text{if } (s, a) = (0, \phi) \\ q\mathbb{1}(s' = 0) + (1-q)\mathbb{1}(s' = 1) & \text{if } (s, a) = (0, 1-\phi) \\ \mathbb{1}(s' = 1) & \text{if } s = 1 \\ q\mathbb{1}(s' = s) + (1-q)\mathbb{1}(s' = 1) & \text{if } s > 1 \end{cases} \quad (91a)$$

and

$$P_h^\phi(s' | s, a) = \mathbb{1}(s' = s), \quad \forall (h, s, a) \in \{2, \dots, H\} \times \mathcal{S} \times \mathcal{A}. \quad (91b)$$

In words, except at step $h = 1$, the MDP always stays in the same state. Additionally, the MDP will always stay in the state subset $\{0, 1\}$ if the initial distribution is supported only on $\{0, 1\}$, in view of (91). Here, p and q are set to be

$$p = 1 - \frac{1}{H} + \Delta \quad \text{and} \quad q = 1 - \frac{1}{H} \quad (92)$$

for some $H \geq e^8$ and Δ (whose value will be specified later) obeying

$$\frac{1}{H} \leq \frac{1}{H^{1-3/\beta}} \leq \frac{1}{2} \quad \text{and} \quad \Delta \leq \frac{1}{2H}, \quad (93)$$

where β is set as

$$\beta := \frac{\log H}{2} \geq 4. \quad (94)$$

The assumption (93) immediately indicates the facts

$$1 > p > q \geq \frac{1}{2}. \quad (95)$$

Moreover, for any $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$, the reward function is defined as

$$r_h(s, a) = \begin{cases} 1 & \text{if } s = 0 \\ 0 & \text{otherwise} \end{cases}. \quad (96)$$

Construction of the history/batch dataset. In the nominal environment \mathcal{M}_ϕ , a batch dataset is generated consisting of K independent sample trajectories each of length H , where each trajectory is generated according to (10), based on the following initial state distribution ρ^b and behavior policy $\pi^b = \{\pi_h^b\}_{h=1}^H$:

$$\rho^b(s) = \mu(s) \quad \text{and} \quad \pi_h^b(a | s) = \frac{1}{2}, \quad \forall (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]. \quad (97)$$

Here, $\mu(s)$ is defined as the following state distribution supported on the state subset $\{0, 1\}$:

$$\mu(s) = \frac{1}{CS} \mathbb{1}(s = 0) + \left(1 - \frac{1}{CS}\right) \mathbb{1}(s = 1), \quad (98)$$

where $\mathbb{1}(\cdot)$ is the indicator function, and $C > 0$ is some constant that will determine the concentration coefficient C_{rob}^* (as we shall detail momentarily) and obeys

$$\frac{1}{CS} \leq \frac{1}{4}. \quad (99)$$

As it turns out, for any MDP \mathcal{M}_ϕ , the occupancy distributions of the above batch dataset are the same (due to symmetry) and admit the following simple characterization:

$$d_1^{b, P^\phi}(0, a) = \frac{1}{2} \mu(0), \quad \forall a \in \mathcal{A}, \quad (100a)$$

$$\frac{\mu(s)}{2} \leq d_h^{b, P^\phi}(s) \leq 2\mu(s), \quad \frac{\mu(s)}{4} \leq d_h^{b, P^\phi}(s, a) \leq \mu(s), \quad \forall (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]. \quad (100b)$$

In addition, we choose the following initial state distribution

$$\rho(s) = \begin{cases} 1, & \text{if } s = 0 \\ 0, & \text{if } s > 0 \end{cases}. \quad (101)$$

The proof of the claim (100) is postponed to Appendix C.4.1

Uncertainty set of the transition kernels. Denote the transition kernel vector as

$$P_{h,s,a}^\phi := P_h^\phi(\cdot | s, a) \in [0, 1]^{1 \times S}. \quad (102)$$

For any $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$, the perturbation of the transition kernels in \mathcal{M}_ϕ is restricted to the following uncertainty set

$$\mathcal{U}^\sigma(P^\phi) := \otimes \mathcal{U}^\sigma(P_{h,s,a}^\phi), \quad \mathcal{U}^\sigma(P_{h,s,a}^\phi) := \left\{ P_{h,s,a} \in \Delta(\mathcal{S}) : \text{KL}(P_{h,s,a} \parallel P_{h,s,a}^\phi) \leq \sigma \right\}, \quad (103)$$

where the radius of the uncertainty set σ obeys

$$\left(1 - \frac{3}{\beta}\right) \log(H) \leq \sigma \leq \left(1 - \frac{2}{\beta}\right) \log(H). \quad (104)$$

Before continuing, we shall introduce some notation for convenience. For any $P_h^\phi(\cdot | s, a)$ in (91), we define the limit of the perturbed kernel transiting to the next state s' from the current state-action pair (s, a) by

$$\underline{P}_h^\phi(s' | s, a) := \inf_{P_{h,s,a} \in \mathcal{U}^\sigma(P_{h,s,a}^\phi)} P_h(s' | s, a), \quad (105)$$

and in particular, denote

$$\underline{p} := \underline{P}_1^\phi(0 | 0, \phi), \quad \underline{q} := \underline{P}_1^\phi(0 | 0, 1 - \phi). \quad (106)$$

Armed with the above definitions, we introduce the following lemma which implies some useful properties of the uncertainty set.

Lemma 9. *When β satisfies (94) and the uncertainty level σ satisfies (234), the perturbed transition kernels obey*

$$\underline{p} \geq \underline{q} \geq \frac{1}{\beta}. \quad (107)$$

Proof. See Appendix C.4.2. \square

Value functions and optimal policies. We take a moment to derive the corresponding value functions and identify the optimal policies. With some abuse of the notation, for any MDP \mathcal{M}_ϕ , we denote $\pi^{*,\phi} = \{\pi_h^{*,\phi}\}_{h=1}^H$ as the optimal policy, and let $V_h^{\pi,\sigma,\phi}$ (resp. $V_h^{*,\sigma,\phi}$) represent the robust value function of policy π (resp. $\pi^{*,\phi}$) at step h with uncertainty radius σ . Armed with these notation, we introduce the following lemma which collects the properties concerning the value functions and optimal policies.

Lemma 10. *For any $\phi = \{0, 1\}$ and any policy π , defining*

$$z_\phi^\pi := \underline{p}\pi_1(\phi | 0) + \underline{q}\pi_1(1 - \phi | 0), \quad (108)$$

it holds that

$$V_1^{\pi,\sigma,\phi}(0) = 1 + z_\phi^\pi(H - 1). \quad (109)$$

In addition, the optimal policies and the optimal value functions obey

$$V_1^{*,\sigma,\phi}(0) = 1 + \underline{p}(H - 1), \quad (110a)$$

$$\forall h \in [H] \setminus \{1\} : \quad V_h^{*,\sigma,\phi}(0) = H - h + 1, \quad (110b)$$

$$\forall h \in [H] : \quad \pi_h^{*,\phi}(\phi | 0) = 1, \quad \pi_h^{*,\phi}(\phi | 1) = 1, \quad V_h^{*,\sigma,\phi}(1) = 0. \quad (110c)$$

The robust single-policy clipped concentrability coefficient C_{rob}^ obeys*

$$2C \leq C_{\text{rob}}^* \leq 4C. \quad (111)$$

Proof. See Appendix C.4.3. \square

In view of Lemma 10, we note that the smallest positive state transition probability of the optimal policy π^* under any MDP \mathcal{M}_ϕ with $\phi \in \{0, 1\}$ thus can be given by

$$P_{\min}^* := \min_{h,s,s'} \left\{ P_h^\phi(s' | s, \pi_h^{*,\phi}(s)) : P_h^\phi(s' | s, \pi_h^{*,\phi}(s)) > 0 \right\} = P_1^\phi(1 | 0, \phi) = 1 - p. \quad (112)$$

C.3.2 ESTABLISHING THE MINIMAX LOWER BOUND

We are now ready to establish the sample complexity lower bound. With the choice of the initial distribution ρ in (101), for any policy estimator $\hat{\pi}$ computed based on the batch dataset, we plan to control the quantity

$$\langle \rho, V_1^{*,\sigma,\phi} - V_1^{\hat{\pi},\sigma,\phi} \rangle = V_1^{*,\sigma,\phi}(0) - V_1^{\hat{\pi},\sigma,\phi}(0).$$

Step 1: converting the goal to estimate ϕ . We make the following claim which shall be verified in Appendix C.4.4: given $\varepsilon \leq \frac{H}{256e^6 \log H}$, choosing

$$\Delta = \frac{128e^6 \sigma(1-q)\varepsilon}{H} \leq \frac{\sigma}{2H \log H} \leq \frac{1}{2H}, \quad (113)$$

which satisfies (93) with the aid of (234) and (92), it holds that for any policy $\hat{\pi}$,

$$\langle \rho, V_1^{\star, \sigma, \phi} - V_1^{\hat{\pi}, \sigma, \phi} \rangle \geq 2\varepsilon(1 - \hat{\pi}_1(\phi | 0)). \quad (114)$$

Armed with this relation between the policy $\hat{\pi}$ and its sub-optimality gap, we are positioned to construct an estimate of ϕ . We denote \mathbb{P}_ϕ as the probability distribution when the MDP is \mathcal{M}_ϕ , for any $\phi \in \{0, 1\}$.

Suppose for the moment that a policy estimate $\hat{\pi}$ achieves

$$\mathbb{P}_\phi \left\{ \langle \rho, V_1^{\star, \sigma, \phi} - V_1^{\hat{\pi}, \sigma, \phi} \rangle \leq \varepsilon \right\} \geq \frac{7}{8}, \quad (115)$$

then in view of (114), we necessarily have $\hat{\pi}_1(\phi | 0) \geq \frac{1}{2}$ with probability at least $\frac{7}{8}$. With this in mind, we are motivated to construct the following estimate $\hat{\phi}$ for $\phi \in \{0, 1\}$:

$$\hat{\phi} = \arg \max_{a \in \{0, 1\}} \hat{\pi}_1(a | 0), \quad (116)$$

which obeys

$$\mathbb{P}_\phi \{ \hat{\phi} = \phi \} \geq \mathbb{P}_\phi \{ \hat{\pi}_1(\phi | 0) > 1/2 \} \geq \frac{7}{8}. \quad (117)$$

In what follows, we would like to show (117) cannot happen without enough samples, which would in turn contradict (114).

Step 2: probability of error in testing two hypotheses. Armed with the above preparation, we shall focus on differentiating the two hypotheses $\phi \in \{0, 1\}$. Towards this, consider the minimax probability of error defined as follows:

$$p_e := \inf_{\psi} \max \{ \mathbb{P}_0(\psi \neq 0), \mathbb{P}_1(\psi \neq 1) \}, \quad (118)$$

where the infimum is taken over all possible tests ψ constructed from the batch dataset.

Let $\mu^{b, \phi}$ (resp. $\mu_h^{b, \phi}(s_h)$) be the distribution of a sample trajectory $\{s_h, a_h\}_{h=1}^H$ (resp. a sample (a_h, s_{h+1}) conditional on s_h) for the MDP \mathcal{M}_ϕ . Following standard results from [Tsybakov & Zaiats (2009), Theorem 2.2] and the additivity of the KL divergence (cf. [Tsybakov & Zaiats (2009, Page 85)]), we obtain

$$\begin{aligned} p_e &\geq \frac{1}{4} \exp \left(-K \text{KL}(\mu^{b, 0} \parallel \mu^{b, 1}) \right) \\ &\geq \frac{1}{4} \exp \left\{ -\frac{1}{2} K \mu(0) \left(\text{KL}(P_1^0(\cdot | 0, 0) \parallel P_1^1(\cdot | 0, 0)) + \text{KL}(P_1^0(\cdot | 0, 1) \parallel P_1^1(\cdot | 0, 1)) \right) \right\}, \end{aligned} \quad (119)$$

where we also use the independence of the K trajectories in the batch dataset in the first line. Here, the second line arises from the chain rule of the KL divergence ([Duchi, 2018, Lemma 5.2.8]) and the Markov property of the sample trajectories (recall that $d_h^{b, P^0} = d_h^{b, P^1}$) according to

$$\begin{aligned} \text{KL}(\mu^{b, 0} \parallel \mu^{b, 1}) &= \sum_{h=1}^H \mathbb{E}_{s_h \sim d_h^{b, P^0}} \left[\text{KL}(\mu_h^{b, 0}(s_h) \parallel \mu_h^{b, 1}(s_h)) \right] \\ &= \sum_{a \in \{0, 1\}} d_1^{b, P^0}(0, a) \text{KL}(P_1^0(\cdot | 0, a) \parallel P_1^1(\cdot | 0, a)) \\ &= \frac{1}{2} \mu(0) \sum_{a \in \{0, 1\}} \text{KL}(P_1^0(\cdot | 0, a) \parallel P_1^1(\cdot | 0, a)), \end{aligned}$$

where the penultimate equality holds by the fact that $P_h^0(\cdot | s, a)$ and $P_h^1(\cdot | s, a)$ only differ when $h = 1$ and $s = 0$, and the last equality follows from (100).

It remains to control the KL divergence terms in (119). Given $p \geq q \geq 1/2$ (cf. (95)), applying Lemma 7 (cf. (41)) yields

$$\begin{aligned} \text{KL}(P_1^0(\cdot | 0, 0) \parallel P_1^1(\cdot | 0, 0)) &= \text{KL}(p \parallel q) \leq \frac{(p-q)^2}{(1-p)p} \stackrel{(i)}{=} \frac{\Delta^2}{p(1-p)} \\ &\stackrel{(ii)}{=} \frac{128^2 e^{12} \sigma^2 (1-q)^2 \varepsilon^2}{H^2 p(1-p)} \\ &\stackrel{(iii)}{\leq} \frac{c_1 \sigma^2 P_{\min}^* \varepsilon^2}{H^2}, \end{aligned} \quad (120)$$

where (i) follows from the definition (92), (ii) holds by plugging in the expression of Δ in (113), (iii) arises from $1 - q \leq 2(1 - p) = 2P_{\min}^*$ (see (93) and (112)), $p > \frac{1}{2}$, as long as c_1 is a large enough constant. It can be shown that $\text{KL}(P_1^0(\cdot | 0, 1) \parallel P_1^1(\cdot | 0, 1))$ can be upper bounded in the same way. Substituting (120) back into (119) demonstrates that: if the sample size is chosen as

$$KH \leq \frac{H^3 SC_{\text{rob}}^* \log 2}{4c_1 P_{\min}^* \sigma^2 \varepsilon^2}, \quad (121)$$

then one necessarily has

$$\begin{aligned} p_e &\geq \frac{1}{4} \exp \left\{ -\frac{1}{2} K \mu(0) \cdot 2 \frac{c_1 \sigma^2 P_{\min}^* \varepsilon^2}{H^2} \right\} \stackrel{(i)}{=} \frac{1}{4} \exp \left\{ -K \frac{c_1 \sigma^2 P_{\min}^* \varepsilon^2}{SCH^2} \right\} \\ &\stackrel{(ii)}{\geq} \frac{1}{4} \exp \left\{ -K \frac{4c_1 \sigma^2 P_{\min}^* \varepsilon^2}{SC_{\text{rob}}^* H^2} \right\} \geq \frac{1}{8}, \end{aligned} \quad (122)$$

where (i) follows from (98) and (ii) holds by (111).

Step 3: putting things together. Finally, suppose that there exists an estimator $\hat{\pi}$ such that

$$\mathbb{P}_0 \{ \langle \rho, V_1^{*, \sigma, 0} - V_1^{\hat{\pi}, \sigma, 0} \rangle > \varepsilon \} < \frac{1}{8} \quad \text{and} \quad \mathbb{P}_1 \{ \langle \rho, V_1^{*, \sigma, 1} - V_1^{\hat{\pi}, \sigma, 1} \rangle > \varepsilon \} < \frac{1}{8}.$$

Then Step 1 tells us that the estimator $\hat{\phi}$ defined in (116) must satisfy

$$\mathbb{P}_0(\hat{\phi} \neq 0) < \frac{1}{8} \quad \text{and} \quad \mathbb{P}_1(\hat{\phi} \neq 1) < \frac{1}{8},$$

which cannot happen under the sample size condition (121) to avoid contradiction with (122). The proof is thus finished.

C.4 PROOF OF AUXILIARY RESULTS

C.4.1 PROOF OF (100)

With the initial state distribution and behavior policy defined in (97), we have for any MDP \mathcal{M}_ϕ with $\phi \in \{0, 1\}$,

$$d_1^{b, P^\phi}(s) = \rho^b(s) = \mu(s),$$

which leads to

$$\forall a \in \mathcal{A}: \quad d_1^{b, P^\phi}(0, a) = \frac{1}{2} \mu(0). \quad (123)$$

In view of (91a), the state occupancy distribution at step $h = 2$ obeys

$$d_2^{b, P^\phi}(0) = \mathbb{P} \left\{ s_2 = 0 \mid s_1 \sim d_1^{b, P^\phi}; \pi^b \right\} = \mu(0) [\pi_1^b(\phi | 0)p + \pi_1^b(1 - \phi | 0)q] = \frac{(p+q)\mu(0)}{2},$$

and

$$d_2^{b, P^\phi}(1) = \mathbb{P} \left\{ s_2 = 1 \mid s_1 \sim d_1^{b, P^\phi}; \pi^b \right\}$$

$$= \mu(0) [\pi_1^b(\phi | 0)(1-p) + \pi_1^b(1-\phi | 0)(1-q)] + \mu(1) = \mu(1) + \frac{(2-p-q)\mu(0)}{2}.$$

With the above result in mind and recalling the assumption in (95), we arrive at

$$\frac{\mu(0)}{2} \leq d_2^{b, P^\phi}(0) \leq \mu(0), \quad \mu(1) \leq d_2^{b, P^\phi}(1) \stackrel{(i)}{\leq} 2\mu(1), \quad (124)$$

where (i) holds by applying (95) and (99) (which implies $\mu(0) \leq \mu(1)$ by the assumption in (99))

$$d_2^{b, P^\phi}(1) = \mu(1) + \frac{(2-p-q)\mu(0)}{2} \leq \mu(1) + \mu(0) \leq 2\mu(1).$$

Finally, from the definitions of $P_h^\phi(\cdot | s, a)$ in (91b) and the Markov property, we arrive at for any $(h, s) \in [H] \times \mathcal{S}$,

$$\frac{\mu(s)}{2} \leq d_h^{b, P^\phi}(s) \leq 2\mu(s), \quad (125)$$

which directly leads to

$$\frac{\mu(s)}{4} \leq d_h^{b, P^\phi}(s, a) = \pi_1^b(a | s) d_h^{b, P^\phi}(s) \leq \mu(s). \quad (126)$$

C.4.2 PROOF OF LEMMA 9

Note that $p \geq q$ can be easily verified since $p > q$, which indicates that the first assertion is true. So we will focus on the second assertion in (107). Towards this, invoking the definition in (40), let σ' be the KL divergence from $\text{Ber}(\frac{1}{\beta})$ to $\text{Ber}(q)$, defined as follows

$$\begin{aligned} \sigma' &:= \text{KL} \left(\text{Ber} \left(\frac{1}{\beta} \right) \parallel \text{Ber}(q) \right) = \frac{1}{\beta} \log \frac{\frac{1}{\beta}}{q} + \left(1 - \frac{1}{\beta} \right) \log \frac{\left(1 - \frac{1}{\beta} \right)}{1-q} \\ &= \left(\frac{1}{\beta} \right) \log \left(\frac{1}{\beta} \right) - \left(\frac{1}{\beta} \right) \log(q) + \left(1 - \frac{1}{\beta} \right) \log(H) + \left(1 - \frac{1}{\beta} \right) \log \left(1 - \frac{1}{\beta} \right), \end{aligned} \quad (127)$$

where the second line uses the definition of q in (92). We claim that σ' satisfies the following relation with σ , which will be proven at the end of this proof:

$$\left(1 - \frac{3}{\beta} \right) \log(H) \leq \sigma \leq \left(1 - \frac{2}{\beta} \right) \log(H) \leq \sigma' \leq \left(1 - \frac{1}{\beta} \right) \log(H). \quad (128)$$

Recalling the definition of the transition kernel in (91a)

$$P_1^\phi(0 | 0, 1-\phi) = q, \quad P_1^\phi(1 | 0, 1-\phi) = 1-q, \quad P_1^\phi(s | 0, 1-\phi) = 0, \quad \forall s \in \mathcal{S} \setminus \{0, 1\},$$

the uncertainty set of the transition kernel with radius σ is thus given as

$$\begin{aligned} \mathcal{U}^\sigma(P_{1,0,1-\phi}^\phi) \\ = \{P_{1,0,1-\phi} \in \Delta(\mathcal{S}) : P(0 | 0, 1-\phi) = q', P(1 | 0, 1-\phi) = 1-q', \text{KL}(\text{Ber}(q') \parallel \text{Ber}(q)) \leq \sigma\}. \end{aligned} \quad (129)$$

Recalling the definition of \underline{q} in (106), we can bound

$$\begin{aligned} \underline{q} &= \inf_{P_{1,0,1-\phi} \in \mathcal{U}^\sigma(P_{1,0,1-\phi}^\phi)} P(0 | 0, 1-\phi) = \inf_{q' : \text{KL}(\text{Ber}(q') \parallel \text{Ber}(q)) \leq \sigma} q' \\ &\stackrel{(i)}{\geq} \inf_{q' : \text{KL}(\text{Ber}(q') \parallel \text{Ber}(q)) \leq \sigma'} q' = \frac{1}{\beta}, \end{aligned}$$

where (i) holds by $\sigma \leq \sigma'$ (cf. (128)) and the last equality follows from applying Lemma 7 (cf. (42)) and (127) to arrive at

$$\forall 0 \leq q' < \frac{1}{\beta} : \quad \text{KL}(\text{Ber}(q') \parallel \text{Ber}(q)) > \text{KL} \left(\text{Ber} \left(\frac{1}{\beta} \right) \parallel \text{Ber}(q) \right) = \sigma'.$$

Proof of (128). To control σ' , we plug in the assumptions in (95) and $\beta \geq 4$ and arrive at the trivial facts

$$\left(\frac{1}{\beta}\right) \log\left(\frac{1}{\beta}\right) - \left(\frac{1}{\beta}\right) \log(q) < 0, \quad \left(1 - \frac{1}{\beta}\right) \log\left(1 - \frac{1}{\beta}\right) < 0.$$

The above facts directly lead to

$$\sigma' \leq \left(1 - \frac{1}{\beta}\right) \log(H). \quad (130)$$

Similarly, observing

$$-1 \leq \left(\frac{1}{\beta}\right) \log\left(\frac{1}{\beta}\right) + \left(1 - \frac{1}{\beta}\right) \log\left(1 - \frac{1}{\beta}\right) \leq 0, \quad -\left(\frac{1}{\beta}\right) \log(q) \geq 0,$$

we arrive at

$$\sigma' \geq -1 + \left(1 - \frac{1}{\beta}\right) \log(H) \geq \left(1 - \frac{2}{\beta}\right) \log(H) \quad (131)$$

as long as $\log H \geq \beta$ (cf. (94)). With (130) and (131) in hand, it is straightforward to see that the choice of the uncertainty radius σ in (234) obeys the advertised bound (128).

C.4.3 PROOF OF LEMMA 10

For notational conciseness, we shall drop the superscript ϕ and use the shorthand $V_h^{\pi, \sigma} = V_h^{\pi, \sigma, \phi}$ and $V_h^{*, \sigma} = V_h^{*, \sigma, \phi}$ whenever it is clear from the context. We begin by deriving the robust value function for any policy π . Starting with state 1, at any step $h \in [H]$, it obeys

$$V_h^{\pi, \sigma}(1) = \mathbb{E}_{a \sim \pi_h(\cdot | 1)} \left[r_h(1, a) + \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,1,a}^\phi)} \mathcal{P} V_{h+1}^{\pi, \sigma} \right] = 0 + V_{h+1}^{\pi, \sigma}(1),$$

where the first equality follows from the robust Bellman consistency equation (cf. (8)), and the second equality follows from the observation that the distribution $P_{h,1,a}^\phi$ is supported solely on state 1 in view of (91a), therefore $\mathcal{U}^\sigma(P_{h,1,a}^\phi) = P_{h,1,a}^\phi$. Leveraging the terminal condition $V_{H+1}^{\pi, \sigma}(1) = 0$, and recursively applying the previous relation, we have

$$V_h^{*, \sigma}(1) = V_h^{\pi, \sigma}(1) = 0, \quad \forall h \in [H]. \quad (132)$$

Similarly, turning to state 0, at any step $h > 1$, the robust value function satisfies

$$V_h^{\pi, \sigma}(0) = \mathbb{E}_{a \sim \pi_h(\cdot | 0)} \left[r_h(0, a) + \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{h,0,a}^\phi)} \mathcal{P} V_{h+1}^{\pi, \sigma} \right] = 1 + V_{h+1}^{\pi, \sigma}(0),$$

which again uses the fact that the distribution $P_{h,0,a}^\phi$ is supported solely on state 0 in view of (91b), therefore $\mathcal{U}^\sigma(P_{h,0,a}^\phi) = P_{h,0,a}^\phi$. Leveraging the terminal condition $V_{H+1}^{\pi, \sigma}(0) = 0$, and recursively applying the previous relation, we have

$$V_h^{*, \sigma}(0) = V_h^{\pi, \sigma}(0) = H - h + 1, \quad 2 \leq h \leq H. \quad (133)$$

Taking (132) and (133) together, it follows that

$$\forall 2 \leq h \leq H : \quad V_h^{\pi, \sigma}(0) > V_h^{\pi, \sigma}(1). \quad (134)$$

Consequently, the robust value function of state 0 at step $h = 1$ satisfies

$$\begin{aligned} V_1^{\pi, \sigma}(0) &= \mathbb{E}_{a \sim \pi_1(\cdot | 0)} \left[r_1(0, a) + \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{1,0,a}^\phi)} \mathcal{P} V_2^{\pi, \sigma} \right] \\ &\stackrel{(i)}{=} 1 + \pi_1(\phi | 0) \left(\inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{1,0,\phi}^\phi)} \mathcal{P} V_2^{\pi, \sigma} \right) + \pi_1(1 - \phi | 0) \left(\inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{1,0,1-\phi}^\phi)} \mathcal{P} V_2^{\pi, \sigma} \right) \end{aligned}$$

$$\begin{aligned}
&\stackrel{(ii)}{=} 1 + \pi_1(\phi|0) \left[\underline{p} V_2^{\pi, \sigma}(0) + (1 - \underline{p}) V_2^{\pi, \sigma}(1) \right] \\
&\quad + \pi_1(1 - \phi|0) \left[\underline{q} V_2^{\pi, \sigma}(0) + (1 - \underline{q}) V_2^{\pi, \sigma}(1) \right] \\
&\stackrel{(iii)}{=} 1 + V_2^{\pi, \sigma}(1) + z_\phi^\pi [V_2^{\pi, \sigma}(0) - V_2^{\pi, \sigma}(1)] \\
&= 1 + z_\phi^\pi V_2^{\pi, \sigma}(0)
\end{aligned} \tag{135}$$

where (i) uses the definition of the reward function in (96), (ii) uses (134) so that the infimum is attained by picking the choice specified in (106) with a smallest probability mass imposed on the transition to state 0. Finally, we plug in the definition (108) of z_ϕ^π in (iii), and the last line follows from (132).

Therefore, taking $\pi = \pi^{\star, \phi}$ in the previous relation directly leads to

$$V_1^{\star, \sigma}(0) = 1 + z_\phi^{\pi^{\star, \phi}} V_2^{\star, \sigma}(0) = 1 + z_\phi^{\pi^{\star, \phi}} (H - 1), \tag{136}$$

where the second equality follows from (133). Observing that the function $(H - 1)z$ is increasing in z and that z_ϕ^π is increasing in $\pi_1(\phi|0)$ (due to the fact $\underline{p} \geq \underline{q}$ in (107)). As a result, the optimal policy obeys

$$\pi_1^{\star, \phi}(\phi|0) = 1 \tag{137}$$

at state 0, and plugging back to (136) gives

$$V_1^{\star, \sigma}(0) = 1 + z_\phi^{\pi^{\star, \phi}} (H - 1) = 1 + \underline{p} (H - 1),$$

where $z_\phi^{\pi^{\star, \phi}} = \underline{p} \pi_1^{\star, \phi}(\phi|0) + \underline{q} \pi_1^{\star, \phi}(1 - \phi|0) = \underline{p}$. For the rest of the states, without loss of generality, we choose the optimal policy obeying

$$\forall h \in [H] : \quad \pi_h^{\star, \phi}(\phi|0) = 1, \quad \pi_h^{\star, \phi}(\phi|1) = 1. \tag{138}$$

Proof of claim (111). Given that $\pi_h^{\star, \phi}(\phi|0) = 1$ for all $h \in [H]$ and $\rho(0) = 1$, for any $P \in \mathcal{U}^\sigma(P^\phi)$, we have

$$\begin{aligned}
d_2^{\star, P}(0, \phi) &= d_2^{\star, P}(0) \pi_2^{\star, \phi}(\phi|0) = d_2^{\star, P}(0) = \mathbb{P}_{s_2 \sim P(\cdot | s_1, \pi_1^{\star, \phi}(s_1))} \{s_2 = 0 | s_1 \sim \rho; \pi^{\star, \phi}\} \\
&= P_1(0|0, \phi) \stackrel{(i)}{\geq} \underline{P}_1^\phi(0|0, \phi) \stackrel{(ii)}{=} \underline{p} \geq \frac{1}{\beta},
\end{aligned} \tag{139}$$

which (i) holds by plugging in the definition (105), (ii) follows from the definition (106), and the final inequality arises from Lemma 9. Hence, for all $2 \leq h \leq H$, by the Markov property and $P_h^\phi(0|0, \phi) = 1$, we have

$$d_h^{\star, P}(0, \phi) = d_2^{\star, P}(0, \phi) \geq \frac{1}{\beta}. \tag{140}$$

Examining the definition of C_{rob}^* in (12), we make the following observations.

- For $h = 1$, we have

$$\begin{aligned}
\max_{(s, a, P) \in \mathcal{S} \times \mathcal{A} \times \mathcal{U}^\sigma(P^\phi)} \frac{\min \{d_1^{\star, P}(s, a), \frac{1}{S}\}}{d_1^{\text{b}, P^\phi}(s, a)} &\stackrel{(i)}{=} \max_{P \in \mathcal{U}^\sigma(P^\phi)} \frac{\min \{d_1^{\star, P}(0, \phi), \frac{1}{S}\}}{d_1^{\text{b}, P^\phi}(0, \phi)} \\
&\stackrel{(ii)}{=} \max_{P \in \mathcal{U}^\sigma(P^\phi)} \frac{1}{S d_1^{\text{b}, P^\phi}(0, \phi)} \stackrel{(iii)}{=} \frac{2}{S \mu(0)} = 2C,
\end{aligned} \tag{141}$$

where (i) holds by $d_1^{\star, P}(s) = \rho(s) = 0$ for all $s \in \mathcal{S} \setminus \{0\}$ (see (101)) and $\pi_h^{\star, \phi}(\phi|0) = 1$ for all $h \in [H]$, (ii) follows from the fact $d_1^{\star, P}(0, \phi) = 1$, (iii) is verified in (100), and the last equality arises from the definition in (98).

- Similarly, for $h = 2$, we arrive at

$$\begin{aligned} \max_{(s,a,P) \in \mathcal{S} \times \mathcal{A} \times \mathcal{U}^\sigma(P^\phi)} \frac{\min \{d_2^{*,P}(s,a), \frac{1}{S}\}}{d_2^{b,P^\phi}(s,a)} &\stackrel{(i)}{=} \max_{s \in \{0,1\}, P \in \mathcal{U}^\sigma(P^\phi)} \frac{\min \{d_2^{*,P}(s,\phi), \frac{1}{S}\}}{d_2^{b,P^\phi}(s,\phi)} \\ &\leq \max_{s \in \{0,1\}, P \in \mathcal{U}^\sigma(P^\phi)} \frac{1}{S d_2^{b,P^\phi}(s,\phi)} \stackrel{(ii)}{\leq} \frac{4}{S\mu(0)} = 4C, \end{aligned} \quad (142)$$

where (i) holds by the optimal policy in (110) and the trivial fact that $d_2^{*,P}(s) = 0$ for all $s \in \mathcal{S} \setminus \{0,1\}$ (see (101) and (91a)), (ii) arises from (100), and the last equality comes from (98).

- For all other steps $h = 3, \dots, H$, observing from the deterministic transition kernels in (91b), it can be easily verified that

$$\max_{(s,a,P) \in \mathcal{S} \times \mathcal{A} \times \mathcal{U}^\sigma(P^\phi)} \frac{\min \{d_h^{*,P}(s,a), \frac{1}{S}\}}{d_h^{b,P^\phi}(s,a)} = \max_{(s,a,P) \in \mathcal{S} \times \mathcal{A} \times \mathcal{U}^\sigma(P^\phi)} \frac{\min \{d_2^{*,P}(s,a), \frac{1}{S}\}}{d_2^{b,P^\phi}(s,a)} \leq 4C. \quad (143)$$

Combining the above cases, we complete the proof by

$$2C \leq C_{\text{rob}}^* = \max_{(h,s,a,P) \in [H] \times \mathcal{S} \times \mathcal{A} \times \mathcal{U}^\sigma(P^\phi)} \frac{\min \{d_h^{*,P}(s,a), \frac{1}{S}\}}{d_h^{b,P^\phi}(s,a)} \leq 4C.$$

C.4.4 PROOF OF THE CLAIM (114)

Recall that by virtue of (108) and (110), we arrive at

$$z_\phi^* := z_{\phi}^{\pi^*,\phi} = \underline{p}\pi_1^{\pi,\phi}(\phi|0) + \underline{q}\pi_1^{\pi,\phi}(1-\phi|0) = \underline{p}.$$

Applying (109) yields

$$\langle \rho, V_1^{*,\sigma,\phi} - V_1^{\pi,\sigma,\phi} \rangle = V_h^{*,\sigma,\phi}(0) - V_h^{\pi,\sigma,\phi}(0) = (\underline{p} - z_\phi^*) (H-1) = (\underline{p} - \underline{q}) (H-1) (1 - \pi_1(\phi|0)), \quad (144)$$

where the last equality uses the definition (108). Therefore, it boils down to control $\underline{p} - \underline{q}$.

To continue, we define an auxiliary value function vector $\bar{V} \in \mathbb{R}^{S \times 1}$ obeying

$$\bar{V}(0) = H-1 \quad \text{and} \quad \bar{V}(s) = 0, \quad \forall s \in \mathcal{S} \setminus \{0\}. \quad (145)$$

With this in hand, applying Lemma 2 gives

$$\begin{aligned} (H-1)(\underline{p} - \underline{q}) &\stackrel{(i)}{=} \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{1,0,\phi}^\phi)} \mathcal{P}\bar{V} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{1,0,1-\phi}^\phi)} \mathcal{P}\bar{V} \\ &= \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(P_{1,0,\phi}^\phi \cdot \exp \left(\frac{-\bar{V}}{\lambda} \right) \right) - \lambda \sigma \right\} - \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(P_{1,0,1-\phi}^\phi \cdot \exp \left(\frac{-\bar{V}}{\lambda} \right) \right) - \lambda \sigma \right\} \\ &\stackrel{(ii)}{\geq} \left\{ -\lambda^* \log \left(P_{1,0,\phi}^\phi \cdot \exp \left(\frac{-\bar{V}}{\lambda^*} \right) \right) - \lambda^* \sigma \right\} - \left\{ -\lambda^* \log \left(P_{1,0,1-\phi}^\phi \cdot \exp \left(\frac{-\bar{V}}{\lambda^*} \right) \right) - \lambda^* \sigma \right\} \\ &= -\lambda^* \left[\log \left(P_{1,0,\phi}^\phi \cdot \exp \left(\frac{-\bar{V}}{\lambda^*} \right) \right) - \log \left(P_{1,0,1-\phi}^\phi \cdot \exp \left(\frac{-\bar{V}}{\lambda^*} \right) \right) \right], \end{aligned} \quad (146)$$

where (i) follows from (see the definition of \underline{p} in (106))

$$\begin{aligned} \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{1,0,\phi}^\phi)} \mathcal{P}\bar{V} &= \underline{P}_1^\phi(0|0, \phi) \bar{V}(0) = (H-1)\underline{p}, \\ \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{1,0,1-\phi}^\phi)} \mathcal{P}\bar{V} &= \underline{P}_1^\phi(0|0, 1-\phi) \bar{V}(0) = (H-1)\underline{q}. \end{aligned}$$

Here, (ii) holds by letting

$$\lambda^* := \arg \max_{\lambda \geq 0} f(\lambda) := \arg \max_{\lambda \geq 0} \left\{ -\lambda \log \left(P_{1,0,1-\phi}^\phi \cdot \exp \left(\frac{-\bar{V}}{\lambda} \right) \right) - \lambda \sigma \right\}. \quad (147)$$

The rest of the proof is then to control (146). We start with the observation that $\lambda^* > 0$; this is because in view of Lemma 3 (cf. (34)), it suffices to verify that

$$\log(1-q) + \sigma \stackrel{(i)}{\leq} \log\left(\frac{1}{H}\right) + \left(1 - \frac{2}{\beta}\right) \log H = -\frac{2}{\beta} \log H < 0, \quad (148)$$

where (i) holds by (234). We now claim the following bound for λ^* holds, whose proof is postponed to the end:

$$\frac{H}{16\sigma} \leq \frac{H-1}{\log(\beta H)} \leq \lambda^* \leq \frac{H-1}{\left(1 - \frac{3}{\beta}\right) \log(H)}, \quad (149)$$

which immediately implies the following by taking exponential maps given $\lambda^* > 0$:

$$\frac{1}{\beta H} \leq e^{-(H-1)/\lambda^*} \leq \frac{1}{H^{1-3/\beta}}. \quad (150)$$

Moving to the second term of (146), it follows that

$$\begin{aligned} & \log\left(P_{1,0,\phi}^\phi \cdot \exp\left(\frac{-\bar{V}}{\lambda^*}\right)\right) - \log\left(P_{1,0,1-\phi}^\phi \cdot \exp\left(\frac{-\bar{V}}{\lambda^*}\right)\right) \\ & \stackrel{(i)}{=} \log \frac{pe^{-(H-1)/\lambda^*} + (1-p)}{qe^{-(H-1)/\lambda^*} + (1-q)} \\ & = \log\left(1 + \frac{(p-q)(e^{-(H-1)/\lambda^*} - 1)}{qe^{-(H-1)/\lambda^*} + (1-q)}\right) \\ & \stackrel{(ii)}{<} -\frac{\Delta(1 - e^{-(H-1)/\lambda^*})}{qe^{-(H-1)/\lambda^*} + (1-q)} \\ & \stackrel{(iii)}{\leq} -\frac{1}{2} \frac{\Delta}{H^{3/\beta}(1-q) + (1-q)} \\ & \leq -\frac{\Delta}{4e^6(1-q)}, \end{aligned} \quad (151)$$

where (i) follows from the definitions in (91) and (145), (ii) holds by $\log(1+x) < x$ for $x \in (-1, \infty)$, (iii) can be verified by (150) and (93):

$$1 - e^{-(H-1)/\lambda^*} \geq 1 - \frac{1}{H^{1-3/\beta}} \geq \frac{1}{2},$$

and the last line uses $H^{3/\beta} = H^{6/\log H} = e^6$ by the definition of β in (94). Plugging (149) and (151) back into (146) and (144), we arrive at

$$\begin{aligned} \langle \rho, V_1^{\star,\sigma,\phi} - V_1^{\pi,\sigma,\phi} \rangle &= (H-1)(p-q)(1 - \pi_1(\phi|0)) \\ &\geq \frac{H\Delta}{64e^6\sigma(1-q)}(1 - \pi_1(\phi|0)) \geq 2\varepsilon(1 - \pi_1(\phi|0)), \end{aligned}$$

where (i) holds by the definition of β in (94) and the last inequality follows directly from the choice of Δ in (113).

Proof of inequality (149). Applying (33) in Lemma 3 to λ^* in (147) leads to the upper bound in (149):

$$\lambda^* \leq \frac{H-1}{\sigma} \leq \frac{H-1}{\left(1 - \frac{3}{\beta}\right) \log(H)}, \quad (152)$$

where the last inequality holds by (234). As a result, we shall focus on showing the lower bounds in (149) in the remainder of the proof.

Recalling the definition of q in (92), we can reparameterize $1-q$ using two positive variables c_q and λ_q (whose choices will be made clearer soon) as follows:

$$1-q = \frac{1}{H} = c_q e^{-(H-1)/\lambda_q}. \quad (153)$$

Deriving the first derivative of the function of interest $f(\lambda)$ in (147) as follows:

$$\begin{aligned}\nabla_{\lambda} f(\lambda) &= \nabla_{\lambda} \left(-\lambda \log \left(P_{1,0,1-\phi}^{\phi} \cdot \exp \left(\frac{-\bar{V}}{\lambda} \right) \right) - \lambda \sigma \right) \\ &\stackrel{(i)}{=} \nabla_{\lambda} \left(-\lambda \log \left(qe^{-(H-1)/\lambda} + 1 - q \right) - \lambda \sigma \right) \\ &= -\sigma - \log \left(qe^{-(H-1)/\lambda} + 1 - q \right) - \frac{1}{\lambda} \cdot \frac{q(H-1)e^{-(H-1)/\lambda}}{qe^{-(H-1)/\lambda} + 1 - q},\end{aligned}\quad (154)$$

where (i) holds by the chosen transition kernels in (91) and the last line arises from basic calculus. To continue, when $\lambda = \lambda_q$, the derivative of the function $f(\lambda)$ can be expressed as

$$\begin{aligned}\nabla_{\lambda} f(\lambda) |_{\lambda=\lambda_q} &= -\sigma - \log \left((1-q) \frac{q}{c_q} + 1 - q \right) + \frac{(1-q) \frac{q}{c_q} \log \frac{1-q}{c_q}}{(1-q) \frac{q}{c_q} + 1 - q} \\ &= -\sigma - \log(1-q) - \log \left(1 + \frac{q}{c_q} \right) + \frac{\frac{q}{c_q} \log \frac{1-q}{c_q}}{\frac{q}{c_q} + 1} \\ &= -\sigma - \log(1-q) \left(1 - \frac{q/c_q}{q/c_q + 1} \right) - \log \left(1 + \frac{q}{c_q} \right) - \frac{\frac{q}{c_q} \log(c_q)}{1 + q/c_q} \\ &\stackrel{(i)}{=} -\sigma + \log H \left(1 - \frac{q/c_q}{q/c_q + 1} \right) - \log \left(1 + \frac{q}{c_q} \right) - \frac{\frac{q}{c_q} \log(c_q)}{1 + q/c_q} \\ &\stackrel{(ii)}{\geq} \log H \left(\frac{2}{\beta} - \frac{q/c_q}{q/c_q + 1} \right) - \log \left(1 + \frac{q}{c_q} \right) - \frac{\frac{q}{c_q} \log(c_q)}{1 + q/c_q} \\ &\stackrel{(iii)}{\geq} \frac{1}{\beta} \log H - \log(1 + \frac{1}{\beta}) - 1 \\ &\geq \frac{1}{\beta} \log H - 2 = 0,\end{aligned}\quad (155)$$

where (i) holds by (153), (ii) follows from the bound of σ in (234), (iii) arises from letting $c_q = \beta \geq 4$ and noting the fact $1/2 \leq q < 1$ (see (95)), leading to

$$\frac{1}{2\beta} \leq \frac{q}{c_q} < \frac{1}{\beta}, \quad \frac{q/c_q}{q/c_q + 1} \leq \frac{1}{\beta}, \quad \frac{\frac{q}{c_q} \log(c_q)}{1 + q/c_q} < 1. \quad (157)$$

Finally, the last line holds by $1/\beta \leq \frac{1}{4}$ and $\log H = 2\beta$ (see (94)).

To proceed, note that the function $f(\lambda)$ is concave with respect to λ . Therefore, observing $\nabla_{\lambda} f(\lambda) |_{\lambda=\lambda_q} \geq 0$ with $c_q = \beta$, we have $\lambda_q \leq \lambda^*$, which implies (see (153))

$$1 - q = \frac{1}{H} = \beta e^{-(H-1)/\lambda_q} \leq \beta e^{-(H-1)/\lambda^*}. \quad (158)$$

The above assertion directly gives

$$\lambda^* \geq \frac{H-1}{\log(\beta H)}.$$

The proof is completed by noticing

$$\frac{H-1}{\log(\beta H)} = \frac{H-1}{\log(H) + \log \beta} \stackrel{(i)}{\geq} \frac{H-1}{2 \log H} \geq \frac{\left(1 - \frac{3}{\beta}\right)(H-1)}{2\sigma} \geq \frac{H}{16\sigma},$$

where (i) follows from (94), the penultimate inequality follows from (234), and the last inequality follows from $\beta \in [4, \infty)$.

D PROBLEM FORMULATION: DISCOUNTED INFINITE-HORIZON RMDPs

In this section, we turn to the studies of distributionally robust offline RL for discounted infinite-horizon MDPs.

D.1 BASICS ABOUT DISCOUNTED INFINITE-HORIZON MDPs

A discounted infinite-horizon MDP can be denoted by $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \gamma, P, r\}$. Here, $\mathcal{S} = \{1, 2, \dots, S\}$ is the state space, $\mathcal{A} = \{1, 2, \dots, A\}$ is the action space, $\gamma \in [0, 1]$ is the discounted factor, $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ represent the transition kernel of the MDP, and $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is the intermediate reward function.

Policy, value/Q function and occupancy distribution. A (possibly random) stationary policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ represents the selection rule of the agent, namely, $\pi(a | s)$ denote the probability of choosing a in state s . With some abuse of notation, let $\pi(s)$ represent the action chosen by π when π is a deterministic policy. With this in mind, let ρ be some initial state distribution. We denote $d^{\pi, P}(s | \rho)$ and $d^{\pi, P}(s, a | \rho)$ respectively as the state occupancy distribution and the state-action occupancy distribution induced by π , namely

$$\forall s \in \mathcal{S} : \quad d^{\pi, P}(s) := (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s | s_0 \sim \rho, \pi, P), \quad (159a)$$

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad d^{\pi, P}(s, a) := \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s | s_0 \sim \rho, \pi, P) \pi(a | s). \quad (159b)$$

Here, the occupancy distributions are conditioned on $s_0 \sim \rho$ and the sequence of actions and states are generated based on policy π and transition kernel P .

In addition, the value function $V^{\pi, P}$ and Q-function $Q^{\pi, P}$ w.r.t. policy π and transition kernel P are defined respectively by

$$\forall s \in \mathcal{S} : \quad V^{\pi, P}(s) := \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right], \quad (160)$$

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad Q^{\pi, P}(s, a) := \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right], \quad (161)$$

where the expectation is taken over the randomness of the trajectory.

D.2 DISTRIBUTIONALLY ROBUST DISCOUNTED INFINITE-HORIZON MDPs

Before continuing, we introduce the vector of a transition kernel P at (s, a) as

$$P_{s,a} := P(\cdot | s, a) \in \mathbb{R}^{1 \times S}, \quad (162)$$

which is used throughout Appendix [D](#) to Appendix [F](#).

In this work, instead of the standard MDP introduced above, we consider the discounted infinite-horizon robust MDPs (RMDPs) represented by $\mathcal{M}_{\text{rob}} = \{\mathcal{S}, \mathcal{A}, \gamma, \mathcal{U}^\sigma(P^0), r\}$. Here, $\mathcal{U}^\sigma(P^0)$ denote the set of possible transition kernels within an uncertainty set centered around a nominal kernel P^0 using the distance measured in terms of the KL divergence. In particular, given an uncertainty level $\sigma > 0$, the uncertainty set around P^0 is specified as

$$\mathcal{U}^\sigma(P^0) := \otimes \mathcal{U}^\sigma(P_{s,a}^0), \quad \mathcal{U}^\sigma(P_{s,a}^0) := \{P_{s,a} \in \Delta(\mathcal{S}) : \text{KL}(P_{s,a} \parallel P_{s,a}^0) \leq \sigma\}. \quad (163)$$

Armed with this, we define the *robust value function* $V^{\pi, \sigma}$ and the *robust Q-function* $Q^{\pi, \sigma}$ in the discounted infinite-horizon setting respectively as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad V^{\pi, \sigma}(s) := \inf_{P \in \mathcal{U}^\sigma(P^0)} V^{\pi, P}(s), \quad Q^{\pi, \sigma}(s, a) := \inf_{P \in \mathcal{U}^\sigma(P^0)} Q^{\pi, P}(s, a).$$

In words, the robust value/Q functions characterize the worst case over all the instances in the uncertainty set.

Optimal policy and robust Bellman equation It is well-known that there exists at least one deterministic policy that maximizes the robust value function and Q-function simultaneously in the infinite-horizon setting as well ([Iyengar, 2005](#); [Nilim & El Ghaoui, 2005](#)). With this in mind, we

denote the optimal policy as π^* and the corresponding *optimal robust value function* (resp. *optimal robust Q function*) as $V^{*,\sigma}$ (resp. $Q^{*,\sigma}$), namely

$$\forall s \in \mathcal{S} : V^{*,\sigma}(s) := V^{\pi^*,\sigma}(s) = \max_{\pi} V^{\pi,\sigma}(s), \quad (164a)$$

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : Q^{*,\sigma}(s, a) := Q^{\pi^*,\sigma}(s, a) = \max_{\pi} Q^{\pi,\sigma}(s, a). \quad (164b)$$

Next, applying (159) with $\pi = \pi^*$, we adopt the the following short-hand notation for the occupancy distributions associated with the optimal policy:

$$\forall s \in \mathcal{S} : d^{*,P}(s) := d^{\pi^*,P}(s), \quad (165a)$$

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : d^{*,P}(s, a) := d^{\pi^*,P}(s, a) = d^{*,P}(s) \mathbb{1}\{a = \pi^*(s)\}. \quad (165b)$$

In addition, the Bellman’s optimality principle can also be extended to the infinite-horizon robust MDPs, which is essential. Specifically, for any policy π (resp. optimal policy π^*), the robust value function and Q-function obey the following *robust Bellman consistency equation* (resp. *robust Bellman optimality equation*):

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : Q^{\pi,\sigma}(s, a) = r(s, a) + \gamma \inf_{P \in \mathcal{U}^\sigma(P^0)} V^{\pi,\sigma}, \quad (166)$$

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : Q^{*,\sigma}(s, a) = r(s, a) + \gamma \inf_{P \in \mathcal{U}^\sigma(P^0)} V^{*,\sigma}. \quad (167)$$

D.3 DISTRIBUTIONALLY ROBUST OFFLINE RL

We observe a batch dataset $\mathcal{D} = \{(s_i, a_i, s'_i)\}_{1 \leq i \leq N}$ consisting of N sample transitions. These transitions are independently generated based on some behavior distribution $d^b \in \Delta(\mathcal{S} \times \mathcal{A})$ over the nominal transition kernel P^0 , i.e.,

$$(s_i, a_i) \stackrel{\text{i.i.d.}}{\sim} d^b \quad \text{and} \quad s'_i \stackrel{\text{i.i.d.}}{\sim} P^0(\cdot | s_i, a_i), \quad 1 \leq i \leq N. \quad (168)$$

Similar to Assumption 1, we design the following *robust single-policy clipped concentrability* for infinite-horizon RMDPs to characterize the quality of the history dataset.

Assumption 2 (Robust single-policy clipped concentrability for infinite-horizon MDPs). *The behavior policy of the history dataset \mathcal{D} satisfies*

$$\max_{(s,a,P) \in \mathcal{S} \times \mathcal{A} \times \mathcal{U}^\sigma(P^0)} \frac{\min \{d^{*,P}(s, a), \frac{1}{S}\}}{d^{b,P^0}(s, a)} \leq C_{\text{rob}}^* \quad (169)$$

for some finite quantity $C_{\text{rob}}^* \in [\frac{1}{S}, \infty)$. Following the convention $0/0 = 0$, we denote C_{rob}^* to be the smallest quantity satisfying (169), and refer to it as the *robust single-policy clipped concentrability coefficient*.

Armed with these, we are ready to introduce the goal in the infinite-horizon case. Given the history dataset \mathcal{D} obeying Assumption 2 for some target accuracy $\varepsilon > 0$, we aim to find a near-optimal robust policy $\hat{\pi}$, which satisfies

$$V^{\hat{\pi},\sigma}(\rho) \geq V^{*,\sigma}(\rho) - \varepsilon \quad (170)$$

in a sample-efficient manner.

E ALGORITHM AND THEORY: DISCOUNTED INFINITE-HORIZON RMDPs

In this section, we present both the model-based algorithm —DRVI-LCB— for robust offline RL and its performance guarantees for the discounted infinite-horizon setting.

E.1 BUILDING AN EMPIRICAL NOMINAL MDP

Recalling that we have N independent samples in the dataset $\mathcal{D} = \{(s_i, a_i, s'_i)\}_{1 \leq i \leq N}$. First, we denote $N(s, a)$ as the total number of sample transitions from (s, a) as

$$N(s, a) := \sum_{i=1}^N \mathbb{1}\{(s_i, a_i) = (s, a)\}. \quad (171)$$

Algorithm 2: Robust value iteration with LCB (DRVI-LCB) for infinite-horizon RMDPs.

-
- 1 **input:** a dataset \mathcal{D} ; reward function r ; uncertainty level σ .
 - 2 **initialization:** $\hat{Q}_0(s, a) = 0$, $\hat{V}_0(s) = 0$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.
 - 3 Compute the empirical nominal transition kernel \hat{P}^0 according to (172);
 - 4 Compute the penalty term $b(s, a)$ according to (176);
 - 5 **for** $m = 1, 2, \dots, M$ **do**
 - 6 **for** $s \in \mathcal{S}, a \in \mathcal{A}$ **do**
 - 7 Set $\hat{Q}_m(s, a)$ according to (179);
 - 8 **for** $s \in \mathcal{S}$ **do**
 - 9 Set $\hat{V}_m(s) = \max_a \hat{Q}_m(s, a)$;
 - 10 **output:** $\hat{\pi}$ s.t. $\hat{\pi}(s) = \arg \max_a \hat{Q}_M(s, a)$ for all $s \in \mathcal{S}$.
-

Armed with $N(s, a)$, we construct the empirical estimate \hat{P}^0 of the nominal kernel P^0 by the visiting frequencies of state-action pairs as follows:

$$\hat{P}^0(s' | s, a) := \begin{cases} \frac{1}{N(s, a)} \sum_{i=1}^N \mathbb{1}\{(s_i, a_i, s'_i) = (s, a, s')\}, & \text{if } N(s, a) > 0 \\ 0, & \text{else} \end{cases} \quad (172)$$

for any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$.

E.2 ALGORITHM: DRVI-LCB FOR INFINITE-HORIZON RMDPs

With the estimate \hat{P}^0 of the nominal transition kernel P^0 in hand, we are positioned to introduce our algorithm DRVI-LCB for infinite-horizon RMDPs, which taking the uncertainty into consideration by incorporating some penalty term inside the value estimation, summarized in Algorithm 2.

The pessimistic robust Bellman operator. To begin with, recall the classical distributionally robust Bellman operator (Zhou et al., 2021; Iyengar, 2005; Nilim & El Ghaoui, 2005) in the infinite-horizon case $\mathcal{T}^\sigma(\cdot) : \mathbb{R}^{\mathcal{S}\mathcal{A}} \rightarrow \mathbb{R}^{\mathcal{S}\mathcal{A}}$,

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad \mathcal{T}^\sigma(Q)(s, a) := r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^0)} \mathcal{P}V, \quad (173)$$

where $V = [V(s)]_{s \in \mathcal{S}}$, $V(s) := \max_a Q(s, a)$.

The contraction property of the above robust Bellman operator plays a fundamental role in the convergence and performance of prior works, e.g. Yang et al. (2021); Zhou et al. (2021); Panaganti & Kalathil (2022). In view of this, to apply the pessimistic principle in RMDPs, we define a pessimistic variant of the robust Bellman operator $\hat{\mathcal{T}}_{\text{pe}}^\sigma(\cdot)$ w.r.t. the empirical nominal kernel \hat{P}^0 as follows:

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad \hat{\mathcal{T}}_{\text{pe}}^\sigma(Q)(s, a) = \max \left\{ r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P}V - b(s, a), 0 \right\}, \quad (174)$$

where $b(s, a)$ denotes the new penalty term that measures the data-dependent uncertainty of the value estimates.

To specify the tailored penalty term $b(s, a)$ in (174), we first introduce an additional term

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad \hat{P}_{\min}(s, a) := \min_{s'} \left\{ \hat{P}^0(s' | s, a) : \hat{P}^0(s' | s, a) > 0 \right\}, \quad (175)$$

which in words represents the smallest positive transition probability of the estimated nominal kernel $\hat{P}^0(s' | s, a)$. Then for some $\delta \in (0, 1)$, some universal constant $c_b > 0$ and any $(s, a) \in \mathcal{S} \times \mathcal{A}$, $b(s, a)$ is defined as

$$b(s, a) = \begin{cases} \min \left\{ \frac{c_b}{\sigma(1-\gamma)} \sqrt{\frac{\log \left(\frac{2N^3 S}{(1-\gamma)\delta} \right)}{\hat{P}_{\min}(s, a) N(s, a)}} + \frac{4}{\sigma N(1-\gamma)}, \frac{1}{1-\gamma} \right\} + \frac{2}{\sigma N} & \text{if } N(s, a) > 0, \\ \frac{1}{1-\gamma} + \frac{2}{\sigma N} & \text{otherwise.} \end{cases} \quad (176)$$

Our proposed pessimistic robust Bellman operator $\hat{T}_{\text{pe}}^\sigma(\cdot)$ (cf. (174)) will play an important role in DRVI-LCB. Encouragingly, although the pessimistic robust Bellman operator $\hat{T}_{\text{pe}}^\sigma(\cdot)$ involves an additional penalty term $b(s, a)$ compared to the classical robust Bellman operator $\mathcal{T}^\sigma(\cdot)$ used in DRVI (Zhou et al., 2021), it still enjoys the celebrated γ -contractive property. Before continuing, we summarize the contraction property below, whose proof is postponed to Appendix F.3.1

Lemma 11 (γ -Contraction). *For any $\gamma \in [\frac{1}{2}, 1)$, the operator $\hat{T}_{\text{pe}}^\sigma(\cdot)$ (cf. (174)) is a γ -contraction w.r.t. $\|\cdot\|_\infty$. Namely, for any $Q_1, Q_2 \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ s.t. $Q_1(s, a), Q_2(s, a) \in [0, \frac{1}{1-\gamma}]$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, one has*

$$\|\hat{T}_{\text{pe}}^\sigma(Q_1) - \hat{T}_{\text{pe}}^\sigma(Q_2)\|_\infty \leq \gamma \|Q_1 - Q_2\|_\infty. \quad (177)$$

Additionally, there exists a unique fixed point $\hat{Q}_{\text{pe}}^{*,\sigma}$ of the operator $\hat{T}_{\text{pe}}^\sigma(\cdot)$ obeying $0 \leq \hat{Q}_{\text{pe}}^{*,\sigma}(s, a) \leq \frac{1}{1-\gamma}$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.

Our algorithm DRVI-LCB for infinite-horizon robust offline RL. Armed with the γ -contraction of the pessimistic robust Bellman operator $\hat{T}_{\text{pe}}^\sigma(\cdot)$, we are positioned to introduce DRVI-LCB for infinite-horizon RMDPs, summarized in Algorithm 2. Specifically, DRVI-LCB can be seen as a value iteration algorithm w.r.t. $\hat{T}_{\text{pe}}^\sigma(\cdot)$ (cf. (174)), whose update rule at the m -th iteration can be formulated as

$$\hat{Q}_m(s, a) = \hat{T}_{\text{pe}}^\sigma(\hat{Q}_{m-1})(s, a) = \max \left\{ r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P}\hat{V}_{m-1} - b(s, a), 0 \right\} \quad (178)$$

for all $m = 1, 2, \dots, M$. In view of strong duality (Hu & Hong, 2013), the above convex problem can be translated into a dual formulation, leading to the following equivalent update rule:

$$\hat{Q}_m(s, a) = \max \left\{ r(s, a) + \sup_{\lambda \geq 0} \left\{ -\lambda \log \left(\hat{P}_{s,a}^0 \cdot \exp \left(\frac{-\hat{V}_{m-1}}{\lambda} \right) \right) - \lambda \sigma \right\} - b(s, a), 0 \right\}, \quad (179)$$

which can be solved efficiently (Iyengar, 2005; Yang et al., 2021; Panaganti & Kalathil, 2022). The computational efficiency is due to the fact that the variable to be optimized is a scalar and the dimension of it won't explode as the size of the state space \mathcal{S} increases (independent w.r.t. the size of the state-action space).

To continue, we initialize the estimates of Q-function (\hat{Q}_0) and value function (\hat{V}_0) to be zero and output the greedy policy of the final Q-estimates (\hat{Q}_M) as the final policy $\hat{\pi}$, namely,

$$\hat{\pi}(s) = \arg \max_a \hat{Q}_M(s, a) \quad \text{for all } s \in \mathcal{S}. \quad (180)$$

Finally, we introduce a useful fact that the iterates $\{\hat{Q}_m\}_{m \geq 0}$ of our algorithm DRVI-LCB converge linearly to the fixed point $\hat{Q}_{\text{pe}}^{*,\sigma}$, summarized in the following lemma; its proof is postponed to Appendix F.3.2

Lemma 12. *Let $\hat{Q}_0 = 0$. The iterates of Algorithm 2 obey*

$$\forall m \geq 0: \quad \hat{Q}_m \leq \hat{Q}_{\text{pe}}^{*,\sigma} \quad \text{and} \quad \|\hat{Q}_m - \hat{Q}_{\text{pe}}^{*,\sigma}\|_\infty \leq \frac{\gamma^m}{1-\gamma}. \quad (181)$$

In addition, choosing $M \geq \frac{\log \frac{\sigma N}{1-\gamma}}{\log \frac{1}{\gamma}}$, one has

$$\|\hat{Q}_M - \hat{Q}_{\text{pe}}^{*,\sigma}\|_\infty \leq \frac{1}{\sigma N}. \quad (182)$$

E.3 PERFORMANCE GUARANTEES: INFINITE-HORIZON RMDPs

Before introducing the main theorems, we first define several essential metrics.

- d_{\min}^b : the smallest positive state-action occupancy distribution of the dataset \mathcal{D} under the nominal model P^0 , i.e.,

$$d_{\min}^b := \min_{s,a} \left\{ d^{b,P^0}(s,a) : d^{b,P^0}(s,a) > 0 \right\}. \quad (183)$$

- P_{\min}^b : the smallest positive state transition probability under the nominal kernel P^0 in the region covered by dataset \mathcal{D} , i.e.,

$$P_{\min}^b := \min_{s,a,s'} \left\{ P^0(s' | s, a) : d^{b,P^0}(s,a) > 0, P^0(s' | s, a) > 0 \right\}. \quad (184)$$

Note that P_{\min}^b is determined only by the state-action pairs covered by the batch dataset \mathcal{D} .

- P_{\min}^* : the smallest positive state transition probability of the optimal robust policy π^* under the nominal kernel P^0 , namely

$$P_{\min}^* := \min_{s,s'} \left\{ P^0(s' | s, \pi^*(s)) : P^0(s' | s, \pi^*(s)) > 0 \right\}. \quad (185)$$

We also note that P_{\min}^* is determined only by the state-action pairs covered by the optimal robust policy π^* under the nominal model P^0 .

We are now positioned to introduce the sample complexity upper bound of DRVI-LCB, together with the minimax lower bound of solving infinite-horizon RMDPs. First, we present the performance guarantees of DRVI-LCB for robust offline RL in the infinite-horizon case, with the proof deferred to Appendix F.1

Theorem 3. *Let c_0 and c_1 be some sufficiently large universal constants. Given an uncertainty level $\sigma > 0$, suppose that the penalty terms in Algorithm 2 are chosen as (176) for sufficiently large c_b . With probability at least $1 - \delta$, the output $\hat{\pi}$ of Algorithm 2 obeys*

$$V^{\star,\sigma}(\rho) - V^{\hat{\pi},\sigma}(\rho) \leq \frac{c_0}{\sigma(1-\gamma)^2} \sqrt{\frac{SC_{\text{rob}}^* \log^2 \left(\frac{(1+\sigma)N^3S}{(1-\gamma)\delta} \right)}{P_{\min}^* N}}, \quad (186)$$

as long as the number of samples N satisfies

$$N \geq \frac{c_1 \log(NS/\delta)}{d_{\min}^b P_{\min}^b}. \quad (187)$$

The result directly indicates that DRVI-LCB can find an ε -optimal policy as long as the sample size in dataset \mathcal{D} exceeds the order of (ignoring the logarithmic factor)

$$\underbrace{\frac{SC_{\text{rob}}^*}{P_{\min}^* (1-\gamma)^4 \sigma^2 \varepsilon^2}}_{\varepsilon\text{-dependent}} + \underbrace{\frac{\log(NS/\delta)}{d_{\min}^b P_{\min}^b}}_{\text{burn-in cost}}. \quad (188)$$

Note that the burn-in cost is independent with the accuracy level ε , which tells us that the sample complexity is on the order of

$$\tilde{O} \left(\frac{SC_{\text{rob}}^*}{P_{\min}^* (1-\gamma)^4 \sigma^2 \varepsilon^2} \right) \quad (189)$$

as long as ε is small enough.

The sample complexity of DRVI-LCB dramatically outperforms prior works, which has been compared in detail in Section 1.2 (cf. Table 1). We highlight that our sample complexity depends linearly on the state space S , in sharp contrast to prior works that all suffer from a quadratic dependency S^2 . To achieve this, we resort to a delicate technique called leave-one-out analysis (Agarwal et al., 2020; Li et al., 2020; 2022) to decouple the statistical dependency introduced across the iterates of robust value iteration, which has potential to be used in deriving tighter sample complexity of other RMDP problems.

In addition, we develop an information-theoretic lower bound for robust offline RL as provided in the following theorem whose proof can be found in Appendix F.2

Theorem 4. Suppose $(S, C_{\text{rob}}^*, \gamma, \sigma, \varepsilon)$ obeying $\frac{1}{1-\gamma} \geq e^8$, $S \geq \log(\frac{1}{1-\gamma})/2$, $C_{\text{rob}}^* \geq 8/S$, $\varepsilon \leq \frac{1}{256e^6(1-\gamma) \log \frac{1}{1-\gamma}}$, and $\log \frac{1}{1-\gamma} - 6 \leq \sigma \leq \log \frac{1}{1-\gamma} - 4$, we consider two infinite-horizon robust MDPs $\mathcal{M}_0, \mathcal{M}_1$, an initial state distribution ρ , and a batch dataset with N independent samples. Consequently, denoting \mathbb{P}_0 (resp. \mathbb{P}_1) as the probability when the MDP is \mathcal{M}_0 (resp. \mathcal{M}_1), one has

$$\inf_{\hat{\pi}} \max \left\{ \mathbb{P}_0(V^{*,\sigma}(\rho) - V^{\hat{\pi},\sigma}(\rho) > \varepsilon), \mathbb{P}_1(V^{*,\sigma}(\rho) - V^{\hat{\pi},\sigma}(\rho) > \varepsilon) \right\} \geq \frac{1}{8},$$

as long as

$$N \leq \frac{c_1 S C_{\text{rob}}^*}{P_{\min}^* (1-\gamma)^2 \sigma^2 \varepsilon^2}.$$

Here, the infimum is taken over all estimator $\hat{\pi}$ and $c_1 > 0$ is some universal constant.

The sample complexity minimax lower bound shown above indicates that no any algorithm can find an ε -optimal policy if the sample complexity is below the order of $\Omega\left(\frac{S C_{\text{rob}}^*}{P_{\min}^* (1-\gamma)^2 \sigma^2 \varepsilon^2}\right)$. It directly confirms that DRVI-LCB is near-optimal up to a polynomial factor of the effective horizon length $\frac{1}{(1-\gamma)}$ (cf. (188)), which is the first provable algorithm with near-optimal sample complexity for infinite-horizon robust offline RL. Moreover, the requirement of the offline history dataset is also much weaker than prior literature on robust offline RL (Yang et al., 2021; Zhou et al., 2021), without the need of full coverage of the state-action space.

F ANALYSIS: DISCOUNTED INFINITE-HORIZON RMDPS

F.1 PROOF OF THE UPPER BOUND: THEOREM 3

In this section, we outline the proof of Theorem 3. Before moving to the main proof, we introduce some notations and important properties that are useful throughout the analysis.

F.1.1 NOTATION AND PRELIMINARY FACTS

Notation. To begin with, we denote the state-action space covered by the batch dataset \mathcal{D} generated in the nominal model P^0 as

$$\mathcal{C}^b = \left\{ (s, a) : d^{b, P^0}(s, a) > 0 \right\}. \quad (190)$$

Armed with it, in view of the definition in (175), we define a similar one based on the exact nominal model P^0 as

$$P_{\min}(s, a) := \min_{s'} \left\{ P^0(s' | s, a) : P^0(s' | s, a) > 0 \right\}, \quad (191)$$

which combined with (184) and (185) directly indicates that

$$P_{\min}^b = \min_{(s,a) \in \mathcal{C}^b} P_{\min}(s, a), \quad P_{\min}^* = \min_s P_{\min}(s, \pi^*(s)). \quad (192)$$

Finally, we denote the robust Q-value and robust value function outputted from Algorithm 2 as

$$\hat{Q} = \hat{Q}_M \quad \text{and} \quad \hat{V} = \hat{V}_M. \quad (193)$$

Properties of $N(s, a)$. We recall a key property of the visiting times $N(s, a)$ over the state-action pair (s, a) which has been established as follows:

Lemma 13 ((Li et al., 2022), Lemma 7). For any $\delta \in (0, 1)$, with probability at least $1 - \delta$, the quantities $\{N(s, a)\}$ in (171) obey

$$\max \left\{ N(s, a), \frac{2}{3} \log \frac{NS}{\delta} \right\} \geq \frac{N d^{b, P^0}(s, a)}{12} \quad (194)$$

simultaneously for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.

Armed with Lemma 13 we introduce another fact that is useful throughout the proof. With probability at least $1 - \delta$, we have

$$\forall (s, a) \in \mathcal{C}^b : \quad N(s, a) \geq \frac{Nd^{b, P^0}(s, a)}{12} \geq \frac{c_1 \log(NS/\delta)}{12P_{\min}(s, a)} \geq -\frac{\log \frac{2NS}{\delta}}{\log(1 - P_{\min}(s, a))}, \quad (195)$$

as long as c_1 is some sufficient large universal constant and (187) holds. The proof is postponed to Appendix F.3.4.

Properties of \hat{Q} and \hat{V} . Invoking Lemma 12 directly leads to

$$\left\| \hat{Q} - \hat{Q}_{\text{pe}}^{*, \sigma} \right\|_{\infty} \leq \frac{1}{\sigma N} \quad (196)$$

and therefore

$$\left\| \hat{V} - \hat{V}_{\text{pe}}^{*, \sigma} \right\|_{\infty} = \max_s \left| \max_a \hat{Q}(s, a) - \max_a \hat{Q}_{\text{pe}}^{*, \sigma}(s, a) \right| \leq \left\| \hat{Q} - \hat{Q}_{\text{pe}}^{*, \sigma} \right\|_{\infty} \leq \frac{1}{\sigma N}, \quad (197)$$

where the penultimate inequality arises from that the maximum operator is a 1-contraction.

F.1.2 PROOF OF THEOREM 3

We are now positioned to outline the proof of Theorem 3 which is separated into several key steps.

Step 1: controlling the estimation uncertainty. In view of the access to only finite and partial coverage samples for estimating the nominal transition kernel P^0 , we need to efficiently control the uncertainty term $\inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\hat{P}_{s,a}^0)} \mathcal{P}\hat{V} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s,a}^0)} \mathcal{P}\hat{V}$. However, the statistical dependency between the estimated value \hat{V} and the kernel estimation $\hat{P}_{s,a}^0$ (since $\hat{P}_{s,a}^0$ will be reused in the update rule (cf. (179)) for all the iterations) adds daunting challenges of controlling it tightly. This is also why the expensive quadratical dependency w.r.t. state space S appears in the sample complexity of the prior works (Zhou et al., 2021; Panaganti & Kalathil, 2022; Yang et al., 2021), which was addressed by the covering theory. To overcome this challenge, we count on a leave-one-out argument motivated by (Agarwal et al., 2020; Li et al., 2020, 2022) to decouple the dependency. The results are summarized as the following lemma with the proof deferred to Appendix F.3.3.

Lemma 14. Suppose the assumptions in Theorem 3 are satisfied and $\gamma \in [\frac{1}{2}, 1)$. Then for all vector \tilde{V} obeying $\left\| \tilde{V} - \hat{V}_{\text{pe}}^{*, \sigma} \right\|_{\infty} \leq \frac{1}{N}$ and $\left\| \tilde{V} \right\|_{\infty} \leq \frac{1}{1-\gamma}$, with probability at least $1 - \delta$, one has

$$\begin{aligned} & \left| \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\hat{P}_{s,a}^0)} \mathcal{P}\tilde{V} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s,a}^0)} \mathcal{P}\tilde{V} \right| \\ & \leq \min \left\{ \frac{c_b}{\sigma(1-\gamma)} \sqrt{\frac{\log(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta})}{\hat{P}_{\min}(s, a)N(s, a)} + \frac{4}{N\sigma(1-\gamma)}}, \frac{1}{1-\gamma} \right\} \end{aligned} \quad (198)$$

simultaneously for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. In addition, with probability at least $1 - \delta$, we have

$$\forall (s, a) \in \mathcal{C}^b : \quad \frac{P_{\min}(s, a)}{8 \log(NS/\delta)} \leq \hat{P}_{\min}(s, a) \leq e^2 P_{\min}(s, a). \quad (199)$$

Step 2: establishing the pessimism property. Armed with above lemma, we shall show that the estimated $\hat{Q}(s, a)$ is a lower bound of $Q^{\hat{\pi}, \sigma}(s, a)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. Towards this, we first recall that the fixed-point $\hat{Q}_{\text{pe}}^{*, \sigma}$ of the pessimistic robust Bellman operator $\hat{\mathcal{T}}_{\text{pe}}^{\sigma}(\cdot)$ (cf. (174)) obeys

$$\hat{Q}_{\text{pe}}^{*, \sigma} = \hat{\mathcal{T}}_{\text{pe}}^{\sigma}(\hat{Q}_{\text{pe}}^{*, \sigma}) = \max \left\{ r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\hat{P}_{s,a}^0)} \mathcal{P}\hat{V}_{\text{pe}}^{*, \sigma} - b(s, a), 0 \right\}. \quad (200)$$

With this in mind, we shall show $\hat{Q} \leq Q^{\hat{\pi}, \sigma}$ in two different conditions of $\hat{Q}_{\text{pe}}^{*, \sigma}$. In the state-action pairs obeying $\hat{Q}_{\text{pe}}^{*, \sigma}(s, a) = 0$, recalling that the initial $\hat{Q}_0 = 0$ and the definition in (193) gives

$$\hat{Q}(s, a) = \hat{Q}_M \leq \hat{Q}_{\text{pe}}^{*, \sigma}(s, a) = 0, \quad (201)$$

where the inequality holds by applying Lemma 12. As a result, $Q^{\hat{\pi},\sigma}(s, a) \geq 0$ directly indicates $\hat{Q}(s, a) \leq Q^{\hat{\pi},\sigma}(s, a)$.

Consequently, we shall focus on the second case when $\hat{Q}_{\text{pe}}^{\star,\sigma}(s, a) > 0$. To continue, we observe that

$$\begin{aligned}
\hat{Q}(s, a) &\stackrel{(i)}{\leq} \hat{Q}_{\text{pe}}^{\star,\sigma}(s, a) + \frac{1}{\sigma N} = r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P} \hat{V}_{\text{pe}}^{\star,\sigma} - b(s, a) + \frac{1}{\sigma N} \\
&\leq r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P} \hat{V} - b(s, a) + \frac{1}{\sigma N} + \gamma \left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P} \hat{V}_{\text{pe}}^{\star,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P} \hat{V} \right| \\
&\stackrel{(ii)}{\leq} r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P} \hat{V} - b(s, a) + \frac{2}{\sigma N} \\
&\leq r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^0)} \mathcal{P} \hat{V} - b(s, a) + \frac{2}{\sigma N} + \gamma \left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P} \hat{V} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^0)} \mathcal{P} \hat{V} \right| \\
&\leq r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^0)} \mathcal{P} \hat{V}, \tag{202}
\end{aligned}$$

where (i) follows from (196), (ii) arises from (197) and the basic fact that infimum operator is a 1-contraction w.r.t. $\|\cdot\|_\infty$, and the final inequality holds by the definition of $b(s, a)$ (see (176)) and Lemma 14.

With above result in mind and invoking the robust Bellman equation $Q^{\hat{\pi},\sigma}(s, a) = r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^0)} \mathcal{P} V^{\hat{\pi},\sigma}$ (see (166)), we arrive at

$$\begin{aligned}
Q^{\hat{\pi},\sigma}(s, a) - \hat{Q}(s, a) &\geq \gamma \left[\inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^0)} \mathcal{P} V^{\hat{\pi},\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^0)} \mathcal{P} \hat{V} \right] \\
&\stackrel{(i)}{=} \gamma \left[\tilde{P}_{s,a} V^{\hat{\pi},\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^0)} \mathcal{P} \hat{V} \right] \geq \gamma \tilde{P}_{s,a} (V^{\hat{\pi},\sigma} - \hat{V}), \tag{203}
\end{aligned}$$

where (i) holds by letting $\tilde{P}_{s,a} := \arg \min_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^0)} \mathcal{P} V^{\hat{\pi},\sigma}$. Consequently, one has

$$\begin{aligned}
\min_{s,a} [Q^{\hat{\pi},\sigma}(s, a) - \hat{Q}(s, a)] &\geq \min_{s,a} [\gamma \tilde{P}_{s,a} (V^{\hat{\pi},\sigma} - \hat{V})] \stackrel{(i)}{\geq} \gamma \min_s [V^{\hat{\pi},\sigma}(s) - \hat{V}(s)] \\
&= \gamma \min_s [Q^{\hat{\pi},\sigma}(s, \hat{\pi}(s)) - \hat{Q}(s, \hat{\pi}(s))] \\
&\geq \gamma \min_{s,a} [Q^{\hat{\pi},\sigma}(s, a) - \hat{Q}(s, a)] \tag{204}
\end{aligned}$$

where (i) follows from $\tilde{P}_{s,a} \in \Delta(\mathcal{S})$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. Noting that $0 < \gamma < 1$, we conclude $Q^{\hat{\pi},\sigma}(s, a) - \hat{Q}(s, a) \geq 0$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ in this case, otherwise (204) won't happen.

Summing up these two cases, we arrive at

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad Q^{\hat{\pi},\sigma}(s, a) \geq \hat{Q}(s, a) \tag{205}$$

and consequently

$$\forall s \in \mathcal{S}: \quad V^{\star,\sigma}(s) \geq V^{\hat{\pi},\sigma}(s) = \max_a Q^{\hat{\pi},\sigma}(s, a) \geq \max_a \hat{Q}(s, a) = \hat{V}(s). \tag{206}$$

Step 3: bounding $V^{\star,\sigma}(s) - V^{\hat{\pi},\sigma}(s)$. First, armed with above pessimistic property (cf. (206)), we convert the goal to control another term in view of

$$V^{\star,\sigma}(s) - V^{\hat{\pi},\sigma}(s) \leq V^{\star,\sigma}(s) - \hat{V}(s). \tag{207}$$

Towards this, we observe that

$$\hat{V}(s) = \max_a \hat{Q}(s, a) \geq \hat{Q}(s, \pi^\star(s)) \stackrel{(i)}{\geq} \hat{Q}_{\text{pe}}^{\star,\sigma}(s, \pi^\star(s)) - \frac{1}{\sigma N}$$

$$\begin{aligned}
&\stackrel{(ii)}{\geq} r(s, \pi^*(s)) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s, \pi^*(s)}^0)} \mathcal{P} \hat{V}_{pe}^{*, \sigma} - b(s, \pi^*(s)) - \frac{1}{\sigma N} \\
&\geq r(s, \pi^*(s)) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s, \pi^*(s)}^0)} \mathcal{P} \hat{V} - b(s, \pi^*(s)) - \frac{1}{\sigma N} \\
&\quad - \gamma \left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s, \pi^*(s)}^0)} \mathcal{P} \hat{V}_{pe}^{*, \sigma} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s, \pi^*(s)}^0)} \mathcal{P} \hat{V} \right| \\
&\stackrel{(iii)}{\geq} r(s, \pi^*(s)) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s, \pi^*(s)}^0)} \mathcal{P} \hat{V} - b(s, \pi^*(s)) - \frac{2}{\sigma N} \\
&\geq r(s, \pi^*(s)) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s, \pi^*(s)}^0)} \mathcal{P} \hat{V} - b(s, \pi^*(s)) - \frac{2}{\sigma N} \\
&\quad - \gamma \left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s, \pi^*(s)}^0)} \mathcal{P} \hat{V} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s, \pi^*(s)}^0)} \mathcal{P} \hat{V} \right| \\
&\geq r(s, \pi^*(s)) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s, \pi^*(s)}^0)} \mathcal{P} \hat{V} - 2b(s, \pi^*(s)), \tag{208}
\end{aligned}$$

where (i) follows from (196), (ii) holds by applying (200), (iii) arises from (197) and the basic fact that the infimum operator is a 1-contraction w.r.t. $\|\cdot\|_\infty$, and the final inequality holds by the definition of $b(s, a)$ (see (176)) and Lemma 14.

To continue, invoking the robust Bellman optimality equation in (167) gives

$$V^{*, \sigma}(s) = Q^{*, \sigma}(s, \pi^*(s)) = r(s, \pi^*(s)) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s, \pi^*(s)}^0)} \mathcal{P} V^{*, \sigma}. \tag{209}$$

Armed with above results and (208), we arrive at

$$\begin{aligned}
V^{*, \sigma}(s) - \hat{V}(s) &\leq Q^{*, \sigma}(s, \pi^*(s)) - \left(r(s, \pi^*(s)) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s, \pi^*(s)}^0)} \mathcal{P} \hat{V} - 2b(s, \pi^*(s)) \right) \\
&\leq \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s, \pi^*(s)}^0)} \mathcal{P} V^{*, \sigma} - \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s, \pi^*(s)}^0)} \mathcal{P} \hat{V} + 2b(s, \pi^*(s)) \\
&\leq \gamma \hat{P}_{s, \pi^*(s)}^{\text{inf}} (V^{*, \sigma} - \hat{V}) + 2b(s, \pi^*(s)) \tag{210}
\end{aligned}$$

where the final inequality holds by introducing the additional notation

$$\hat{P}_{s, \pi^*(s)}^{\text{inf}} := \operatorname{argmin}_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s, \pi^*(s)}^0)} \mathcal{P} \hat{V} \tag{211}$$

and evidently,

$$\inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s, \pi^*(s)}^0)} \mathcal{P} V^{*, \sigma} \leq \hat{P}_{s, \pi^*(s)}^{\text{inf}} V^{*, \sigma}, \quad \text{and} \quad \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s, \pi^*(s)}^0)} \mathcal{P} \hat{V} = \hat{P}_{s, \pi^*(s)}^{\text{inf}} \hat{V}.$$

Before continuing, for convenience, we introduce a matrix $\hat{P}^{\text{inf}} \in \mathbb{R}^{S \times S}$ and a vector $b^* \in \mathbb{R}^S$, where their s -th rows (resp. entries) are defined as

$$[\hat{P}^{\text{inf}}]_{s, \cdot} = \hat{P}_{s, \pi^*(s)}^{\text{inf}}, \quad \text{and} \quad b^*(s) = b(s, \pi^*(s)). \tag{212}$$

With these notation in mind, applying (210) leads to

$$\begin{aligned}
\langle \rho, V^{*, \sigma} - \hat{V} \rangle &= \sum_{s \in \mathcal{S}} \rho(s) (V^{*, \sigma}(s) - \hat{V}(s)) \\
&\leq \gamma \sum_{s \in \mathcal{S}} \rho(s) \hat{P}_{s, \pi^*(s)}^{\text{inf}} (V^{*, \sigma} - \hat{V}) + 2 \sum_{s \in \mathcal{S}} \rho(s) b(s, \pi^*(s))
\end{aligned}$$

$$= \gamma \rho^\top \hat{P}^{\text{inf}} (V^{*,\sigma} - \hat{V}) + 2\rho^\top b^*. \quad (213)$$

Applying the above result recursively gives

$$\begin{aligned} \langle \rho, V^{*,\sigma} - \hat{V} \rangle &\leq \gamma \rho^\top \hat{P}^{\text{inf}} (V^{*,\sigma} - \hat{V}) + 2\rho^\top b^* \\ &\leq \gamma \left(\gamma \rho^\top \hat{P}^{\text{inf}} \right) \hat{P}^{\text{inf}} (V^{*,\sigma} - \hat{V}) + 2 \left(\gamma \rho^\top \hat{P}^{\text{inf}} \right) b^* + 2\rho^\top b^* \\ &\leq \dots \leq \left\{ \lim_{i \rightarrow \infty} \gamma^i \rho^\top \left(\hat{P}^{\text{inf}} \right)^i (V^{*,\sigma} - \hat{V}) \right\} + 2\rho^\top \sum_{i=0}^{\infty} \gamma^i \left(\hat{P}^{\text{inf}} \right)^i b^* \\ &\stackrel{(i)}{\leq} 2\rho^\top \sum_{i=0}^{\infty} \gamma^i \left(\hat{P}^{\text{inf}} \right)^i b^* = 2\rho^\top \left(I - \gamma \hat{P}^{\text{inf}} \right)^{-1} b^*, \end{aligned} \quad (214)$$

where (i) holds by $|\rho^\top \left(\hat{P}^{\text{inf}} \right)^i (V^{*,\sigma} - \hat{V})| \leq \frac{1}{1-\gamma}$ for all $i \geq 0$ and the basic fact that $\lim_{i \rightarrow \infty} \gamma^i \rho^\top \left(\hat{P}^{\text{inf}} \right)^i (V^{*,\sigma} - \hat{V}) = 0$ since $\lim_{i \rightarrow \infty} \gamma^i = 0$ for all $0 \leq \gamma < 1$.

To further characterize the above performance gap, invoking the definition of $d^{*,P}$ (cf. (159) and (165a)), we arrive at

$$\left(d^{*,\hat{P}^{\text{inf}}} \right)^\top = (1-\gamma) \rho^\top \sum_{t=0}^{\infty} \gamma^t \left(\hat{P}^{\text{inf}} \right)^t = (1-\gamma) \rho^\top \left(I - \gamma \hat{P}^{\text{inf}} \right)^{-1}. \quad (215)$$

Combined with (207), plugging in above result back into (214) yields

$$\langle \rho, V^{*,\sigma} - V^{\pi^*,\sigma} \rangle \leq \langle \rho, V^{*,\sigma} - \hat{V} \rangle \leq \frac{2}{1-\gamma} \langle d^{*,\hat{P}^{\text{inf}}}, b^* \rangle. \quad (216)$$

Step 4: controlling $\langle d^{*,\hat{P}^{\text{inf}}}, b^* \rangle$ using concentrability. Before continuing, note that $\hat{P}^{\text{inf}} \in \mathcal{U}^\sigma(P^0)$ (see (211) and (212)), which in words means \hat{P}^{inf} is some transition kernel inside $\mathcal{U}^\sigma(P^0)$ — the uncertainty set around the nominal kernel P^0 .

Observing that we can express $\langle d^{*,\hat{P}^{\text{inf}}}, b^* \rangle = \sum_{s \in \mathcal{S}} d^{*,\hat{P}^{\text{inf}}}(s) b^*(s)$, we divide the set of state into two types and control them separately.

- **For $s \in \mathcal{S}$ where $\max_{P \in \mathcal{U}^\sigma(P^0)} d^{*,P}(s, \pi^*(s)) = 0$.** In this case, one has for any $\tilde{P} \in \mathcal{U}^\sigma(P^0)$,

$$0 \leq d^{*,\tilde{P}}(s) = d^{*,\tilde{P}}(s, \pi^*(s)) \leq \max_{P \in \mathcal{U}^\sigma(P^0)} d^{*,P}(s, \pi^*(s)) = 0, \quad (217)$$

which consequently indicates

$$d^{*,\hat{P}^{\text{inf}}}(s) = 0. \quad (218)$$

- **For $s \in \mathcal{S}$ where $\max_{P \in \mathcal{U}^\sigma(P^0)} d^{*,P}(s, \pi^*(s)) > 0$.** For such state s , we claim that

$$d^{\mathbf{b},P^0}(s, \pi^*(s)) > 0 \quad \text{and} \quad (s, \pi^*(s)) \in \mathcal{C}^{\mathbf{b}}, \quad (219)$$

which can be verified by recalling Assumption 2 which requires that the history dataset \mathcal{D} obeys

$$\max_{P \in \mathcal{U}^\sigma(P^0)} \frac{\min \{d^{*,P}(s, \pi^*(s)), \frac{1}{S}\}}{d^{\mathbf{b},P^0}(s, \pi^*(s))} = \max_{P \in \mathcal{U}^\sigma(P^0)} \frac{\min \{d^{*,P}(s), \frac{1}{S}\}}{d^{\mathbf{b},P^0}(s)} \leq C_{\text{rob}}^* < \infty. \quad (220)$$

To continue, invoking the fact in (195) with $(s, \pi^*(s)) \in \mathcal{C}^{\mathbf{b}}$ gives

$$N(s, \pi^*(s)) \geq \frac{N d^{\mathbf{b},P^0}(s, \pi^*(s))}{12}$$

$$\stackrel{(i)}{\geq} \frac{N \max_{P \in \mathcal{U}^\sigma(P^0)} \min \{d^{*,P}(s, \pi^*(s)), \frac{1}{S}\}}{12C_{\text{rob}}^*} \geq \frac{N \min \{d^{*,\hat{P}^{\text{inf}}}(s), \frac{1}{S}\}}{12C_{\text{rob}}^*}, \quad (221)$$

where (i) holds by Assumption 2 and the last inequality holds by $\hat{P}^{\text{inf}} \in \mathcal{U}^\sigma(P^0)$. With this in mind, we can control the pessimistic penalty $b^*(s)$ (cf. (176)) by

$$\begin{aligned} b^*(s) &\leq \frac{c_b}{\sigma(1-\gamma)} \sqrt{\frac{\log\left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta}\right)}{\hat{P}_{\min}(s, \pi^*(s))N(s, \pi^*(s))}} + \frac{4}{\sigma N(1-\gamma)} + \frac{2}{\sigma N} \\ &\stackrel{(i)}{\leq} \frac{4c_b}{\sigma(1-\gamma)} \sqrt{\frac{\log^2\left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta}\right)}{P_{\min}(s, \pi^*(s))N(s, \pi^*(s))}} + \frac{4}{\sigma N(1-\gamma)} + \frac{2}{\sigma N} \\ &\leq \frac{16c_b}{\sigma(1-\gamma)} \sqrt{\frac{C_{\text{rob}}^* \log^2\left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta}\right)}{P_{\min}(s, \pi^*(s))N \min \{d^{*,\hat{P}^{\text{inf}}}(s), \frac{1}{S}\}}} + \frac{6}{\sigma N(1-\gamma)} \\ &\leq \frac{20c_b}{\sigma(1-\gamma)} \sqrt{\frac{C_{\text{rob}}^* \log^2\left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta}\right)}{P_{\min}(s, \pi^*(s))N \min \{d^{*,\hat{P}^{\text{inf}}}(s), \frac{1}{S}\}}}, \end{aligned}$$

where (i) arises from (199), the penultimate inequality follows from (221), and the last inequality holds as long as c_b is large enough.

Summing up the results in above two cases, we arrive at

$$\begin{aligned} \langle d^{*,\hat{P}^{\text{inf}}}, b^* \rangle &= \sum_{s \in \mathcal{S}} d^{*,\hat{P}^{\text{inf}}}(s) b^*(s) \\ &\leq \sum_{s \in \mathcal{S}} d^{*,\hat{P}^{\text{inf}}}(s) \frac{20c_b}{\sigma(1-\gamma)} \sqrt{\frac{C_{\text{rob}}^* \log^2\left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta}\right)}{P_{\min}(s, \pi^*(s))N \min \{d^{*,\hat{P}^{\text{inf}}}(s), \frac{1}{S}\}}} \\ &\stackrel{(i)}{\leq} \frac{20c_b}{\sigma(1-\gamma)} \sqrt{\sum_{s \in \mathcal{S}} d^{*,\hat{P}^{\text{inf}}}(s) \frac{C_{\text{rob}}^* \log^2\left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta}\right)}{P_{\min}(s, \pi^*(s))N \min \{d^{*,\hat{P}^{\text{inf}}}(s), \frac{1}{S}\}}} \sqrt{\sum_{s \in \mathcal{S}} d^{*,\hat{P}^{\text{inf}}}(s)} \\ &\leq \frac{40c_b}{\sigma(1-\gamma)} \sqrt{\frac{SC_{\text{rob}}^* \log^2\left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta}\right)}{P_{\min}^* N}}, \end{aligned} \quad (222)$$

where (i) arises from applying the Cauchy-Schwarz inequality, and the last inequality holds by $P_{\min}(s, \pi^*(s)) \geq P_{\min}^*$ for all $s \in \mathcal{S}$ (see (192)) and the following fact that has been established in (67):

$$\sum_{s \in \mathcal{S}} \frac{d^{*,\hat{P}^{\text{inf}}}(s)}{\min \{d^{*,\hat{P}^{\text{inf}}}(s), \frac{1}{S}\}} \leq 2S. \quad (223)$$

Step 5: finishing up the proof. Finally, inserting (222) back into (216), we complete the proof: with probability at least $1 - 2\delta$, one has

$$\langle \rho, V^{*,\sigma} - V^{\hat{\pi},\sigma} \rangle \leq \frac{2}{1-\gamma} \langle d^{*,\hat{P}^{\text{inf}}}, b^* \rangle \leq \frac{80c_b}{\sigma(1-\gamma)^2} \sqrt{\frac{SC_{\text{rob}}^* \log^2\left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta}\right)}{P_{\min}^* N}}. \quad (224)$$

F.2 PROOF OF THE LOWER BOUND: THEOREM 4

We shall first construct some hard discounted infinite-horizon MDP instances and then characterize the sample complexity requirements over these instances.

F.2.1 CONSTRUCTION OF HARD PROBLEM INSTANCES

Construction of a collection of hard MDPs. Suppose there are two MDPs

$$\{\mathcal{M}_\phi = (\mathcal{S}, \mathcal{A}, P^\phi, r, \gamma) \mid \phi = \{0, 1\}\}. \quad (225)$$

Here, γ is the discount parameter, the state space is $\mathcal{S} = \{0, 1, \dots, S-1\}$, and the action space is $\mathcal{A} = \{0, 1\}$. The transition kernel P^ϕ of any constructed MDP \mathcal{M}_ϕ is defined as

$$P^\phi(s' \mid s, a) = \begin{cases} p\mathbb{1}(s' = 2) + (1-p)\mathbb{1}(s' = 1) & \text{if } (s, a) = (0, \phi) \\ q\mathbb{1}(s' = 2) + (1-q)\mathbb{1}(s' = 1) & \text{if } (s, a) = (0, 1-\phi) \\ \mathbb{1}(s' = s) & \text{if } s = 1 \text{ or } s = 2 \\ q\mathbb{1}(s' = s) + (1-q)\mathbb{1}(s' = 1) & \text{if } s > 2 \end{cases} \quad (226)$$

where p and q are set as

$$p = \gamma + \Delta \quad \text{and} \quad q = \gamma \quad (227)$$

for some γ and Δ obeying

$$1 - \gamma \leq 1/e^8 \leq \frac{1}{2} \quad \text{and} \quad \Delta \leq \frac{1}{2}(1 - \gamma). \quad (228)$$

Here, Δ is some value that will be introduced later. Consequently, applying (227) directly leads to

$$1 \geq p \geq q \geq \gamma \geq \frac{1}{2}. \quad (229)$$

Regarding the introduced transition kernel, in words, state 1 and 2 are absorbing states, and states $s > 2$ will stay without moving or go to state 1. The action will only influence the transitions in state 0. In addition, if the initial distribution is supported on states $\{0, 1, 2\}$, the MDP will always stay in the state $\{1, 2\}$ after the first transition.

Finally, we define the reward function as

$$r(s, a) = \begin{cases} 1 & \text{if } s = 0 \text{ or } s = 2 \\ 0 & \text{otherwise} \end{cases}. \quad (230)$$

Uncertainty set of the transition kernels. Then we introduce the considered robust MDPs with some tailored radius σ of the uncertainty set, along with some useful properties.

To begin with, we introduce an important constant β defined as

$$\beta := \frac{\log \frac{1}{1-\gamma}}{2} \geq 4. \quad (231)$$

Then, for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, we denote the transition kernel vector of \mathcal{M}_ϕ as

$$P_{s,a}^\phi := P^\phi(\cdot \mid s, a) \in [0, 1]^{1 \times S}. \quad (232)$$

Armed with these, the perturbed transition kernels in \mathcal{M}_ϕ is limited to the following uncertainty set

$$\mathcal{U}^\sigma(P^\phi) := \otimes \mathcal{U}^\sigma(P_{s,a}^\phi), \quad \mathcal{U}^\sigma(P_{s,a}^\phi) := \{P_{s,a} \in \Delta(\mathcal{S}) : \text{KL}(P_{s,a} \parallel P_{s,a}^\phi) \leq \sigma\}, \quad (233)$$

where the radius of the uncertainty set σ obeys

$$\left(1 - \frac{3}{\beta}\right) \log \frac{1}{1-\gamma} \leq \sigma \leq \left(1 - \frac{2}{\beta}\right) \log \frac{1}{1-\gamma}. \quad (234)$$

Next, to introduce the properties, for any $P^\phi(\cdot \mid s, a)$ in (226), we denote the minimum limit of the perturbed distribution transiting from the current state-action pair (s, a) to the next state s' as

$$\underline{P}^\phi(s' \mid s, a) := \inf_{P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^\phi)} P(s' \mid s, a). \quad (235)$$

As the transition from state 0 to state 2 plays an important role in the analysis, in particular, we denote

$$\underline{p} := \underline{P}^\phi(2 \mid 0, \phi), \quad \underline{q} := \underline{P}^\phi(2 \mid 0, 1 - \phi). \quad (236)$$

With these definitions in mind, we summarize some useful properties of the uncertainty set in the following lemma.

Lemma 15. Suppose β satisfies (231) and the uncertainty level σ satisfies (234). The perturbed transition kernels obey

$$\underline{p} \geq \underline{q} \geq \frac{1}{\beta}. \quad (237)$$

Proof. The proof follows from exactly the pipeline in Appendix C.4.2 except replacing H with $\frac{1}{1-\gamma}$. We omit the details for brevity. \square

Construction of the history/batch dataset. Before continuing, we define a useful state distribution (only supported on the state subset $\{0, 1, 2\}$):

$$\mu(s) = \frac{1}{CS} \mathbb{1}(s=0) + \frac{1}{CS} \mathbb{1}(s=2) + \left(1 - \frac{2}{CS}\right) \mathbb{1}(s=1), \quad (238)$$

where $C > 0$ is some constant that determines the robust concentrability coefficient C_{rob}^* (will be introduced momentarily) and obeys

$$\frac{1}{CS} \leq \frac{1}{4}. \quad (239)$$

Generated over the nominal environment \mathcal{M}_ϕ , a batch dataset consists of N i.i.d samples $\{(s_i, a_i, s'_i)\}_{1 \leq i \leq N}$ according to (168), with the occupancy state (resp. state action) distribution chosen to be:

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad d^{\text{b}, P^\phi}(s) = \mu(s) \quad \text{and} \quad d^{\text{b}, P^\phi}(s, a) = \frac{\mu(s)}{2}. \quad (240)$$

Additionally, we choose the following initial state distribution:

$$\rho(s) = \begin{cases} 1, & \text{if } s=0 \\ 0, & \text{otherwise} \end{cases}. \quad (241)$$

Value functions and optimal policies. Now we are positioned to derive the corresponding robust value functions and identify the optimal policies. To clarify the notations, for any MDP \mathcal{M}_ϕ , we denote π_ϕ^* as the optimal policy. In addition, we denote the robust value function of any policy π (resp. optimal policy π_ϕ^*) as $V_\phi^{\pi, \sigma}$ (resp. $V_\phi^{*, \sigma}$) with uncertainty radius σ . Then, we can introduce the following lemma which describes some important properties of the robust value functions and optimal policies.

Lemma 16. For any $\phi = \{0, 1\}$ and any policy π , one has

$$V_\phi^{\pi, \sigma}(0) = 1 + \gamma z_\phi^\pi \frac{1}{1-\gamma}, \quad (242)$$

where z_ϕ^π is defined as

$$z_\phi^\pi := \underline{p}\pi(\phi|0) + \underline{q}\pi(1-\phi|0). \quad (243)$$

In addition, the optimal value functions and the optimal policies obey

$$V_\phi^{*, \sigma}(0) = 1 + \gamma \underline{p} \frac{1}{1-\gamma}, \quad V_\phi^{*, \sigma}(2) = \frac{1}{1-\gamma}, \quad V_\phi^{*, \sigma}(s) = 0 \quad \text{for } s=1 \text{ or } s>2 \quad (244a)$$

$$\pi_\phi^*(\phi|0) = 1, \quad \text{for } s \in \mathcal{S}. \quad (244b)$$

Moreover, choosing $S \geq \beta$, the robust single-policy clipped concentrability coefficient C_{rob}^* obeys

$$C_{\text{rob}}^* = 2C. \quad (245)$$

Proof. See Appendix F.2.3. \square

F.2.2 ESTABLISHING THE MINIMAX LOWER BOUND

Now we are positioned to provide the sample complexity lower bound. Before starting, we introduce a useful notation representing the smallest positive state transition probability of the optimal policy π_ϕ^* under any robust MDP \mathcal{M}_ϕ with $\phi \in \{0, 1\}$:

$$P_{\min}^* := \min_{s, s'} \left\{ P^\phi(s' | s, \pi_\phi^*(s)) : P^\phi(s' | s, \pi_\phi^*(s)) > 0 \right\} = P^\phi(1|0, \phi) = 1 - p. \quad (246)$$

To continue, the goal is to control the quantity w.r.t. any policy estimator $\hat{\pi}$ based on the batch dataset and chosen initial distribution ρ in (241)

$$\left\langle \rho, V_\phi^{*, \sigma} - V_\phi^{\hat{\pi}, \sigma} \right\rangle = V_\phi^{*, \sigma}(0) - V_\phi^{\hat{\pi}, \sigma}(0). \quad (247)$$

Towards this, we first introduce the following lemma:

Lemma 17. *Given $\epsilon \leq \frac{1}{256e^6(1-\gamma)\log(\frac{1}{1-\gamma})}$, choosing*

$$\Delta = 128e^6\sigma(1-q)\epsilon(1-\gamma), \quad (248)$$

one has Δ obeys (satisfying (228))

$$\Delta \leq \frac{\sigma(1-\gamma)}{2\log\left(\frac{1}{1-\gamma}\right)} \leq \frac{1}{2}(1-\gamma), \quad (249)$$

and for any policy $\hat{\pi}$,

$$V_\phi^{*, \sigma}(0) - V_\phi^{\hat{\pi}, \sigma}(0) \geq 2\epsilon(1 - \hat{\pi}(\phi | 0)). \quad (250)$$

Proof. This lemma can be verified by following the same pipeline in Appendix C.4.4 except replacing H with $\frac{1}{1-\gamma}$ and with the additional condition $\gamma \geq \frac{1}{2}$. \square

Armed with this lemma, following the same pipeline in Appendix C.3.2 we can complete the proof by observing that: let c_1 be some sufficient large constant, as long as the sample size is chosen as

$$N \leq \frac{SC_{\text{rob}}^* \log 2}{4c_1 P_{\min}^* \sigma^2 (1-\gamma)^2 \epsilon^2}, \quad (251)$$

then we necessarily has

$$\inf_{\hat{\pi}} \max_{\phi \in \{0, 1\}} \mathbb{P}_\phi \left\{ \left\langle \rho, V_\phi^{*, \sigma} - V_\phi^{\hat{\pi}, \sigma} \right\rangle \right\} \geq \frac{1}{8}, \quad (252)$$

where \mathbb{P}_ϕ denote the probability conditioned on that the MDP is \mathcal{M}_ϕ . We omit the details for brevity and complete the proof.

F.2.3 PROOF OF LEMMA I6

To start, for any \mathcal{M}_ϕ with $\phi \in \{0, 1\}$, it is easily observed that for any policy π , the robust value function of a state $s > 0$ obey

$$V_\phi^{\pi, \sigma}(2) = \sum_{t=0}^{\infty} \gamma^t \cdot 1 = \frac{1}{1-\gamma}, \quad (253a)$$

$$V_\phi^{\pi, \sigma}(s) = \sum_{t=0}^{\infty} \gamma^t \cdot 0 = 0, \quad \text{for } s = 1 \text{ or } s > 2. \quad (253b)$$

Similarly, the robust value function of state 0 satisfies

$$V_\phi^{\pi, \sigma}(0) = \mathbb{E}_{a \sim \pi(\cdot | 0)} \left[r(0, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,a}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} \right]$$

$$\begin{aligned}
&\stackrel{(i)}{=} 1 + \gamma\pi(\phi|0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,\phi}^\phi)} \mathcal{P}V_\phi^{\pi,\sigma} + \gamma\pi(1-\phi|0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,1-\phi}^\phi)} \mathcal{P}V_\phi^{\pi,\sigma} \\
&\stackrel{(ii)}{=} 1 + \gamma\pi(\phi|0) \left[\underline{p}V_\phi^{\pi,\sigma}(2) + (1-\underline{p})V_\phi^{\pi,\sigma}(1) \right] \\
&\quad + \gamma\pi(1-\phi|0) \left[\underline{q}V_\phi^{\pi,\sigma}(2) + (1-\underline{q})V_\phi^{\pi,\sigma}(1) \right] \\
&\stackrel{(iii)}{=} 1 + \gamma V_\phi^{\pi,\sigma}(1) + \gamma z_\phi^\pi \left[V_\phi^{\pi,\sigma}(2) - V_\phi^{\pi,\sigma}(1) \right] \\
&= 1 + \gamma z_\phi^\pi V_\phi^{\pi,\sigma}(2)
\end{aligned} \tag{254}$$

where (i) holds by the reward function defined in (230), (ii) arises from (253) which indicates $V_\phi^{\pi,\sigma}(2) \geq V_\phi^{\pi,\sigma}(1)$ so that the infimum is obtained by picking the smallest possible mass on the transition to state 2 which reaches (236), (iii) follows from plugging in the definition of z_ϕ^π in (243), and the last identity is due to (253). Consequently, taking $\pi = \pi_\phi^*$, we directly arrive at

$$V_\phi^{*,\sigma}(0) = 1 + \gamma z_\phi^{\pi^*} V_\phi^{\pi^*,\sigma}(2) = 1 + \gamma z_\phi^{\pi^*} \frac{1}{1-\gamma}, \tag{255}$$

which holds by (253). Observing that the function $z \frac{\gamma}{1-\gamma}$ is increasing in z and z_ϕ^π is also increasing in $\pi(\phi|0)$ (see the fact $\underline{p} \geq \underline{q}$ in (237)), the optimal policy in state 0 obeys

$$\pi_\phi^*(\phi|0) = 1. \tag{256}$$

Finally, plugging the above fact back into (243) leads to

$$z_\phi^* := z_\phi^{\pi^*} = \underline{p}\pi_\phi^*(\phi|0) + \underline{q}\pi_\phi^*(1-\phi|0) = \underline{p}, \tag{257}$$

which combined with (255) yields

$$V_\phi^{*,\sigma}(0) = 1 + \gamma \underline{p} \frac{1}{1-\gamma}. \tag{258}$$

For the rest of states $s > 0$, since the action does not influence the transition, without loss of generality, we choose the optimal policy to obey

$$\forall s > 0: \quad \pi_\phi^*(\phi|s) = 1. \tag{259}$$

Proof of (245). To begin with, for any MDP \mathcal{M}_ϕ with $\phi \in \{0, 1\}$, recall the definition of C_{rob}^*

$$C_{\text{rob}}^* = \max_{(s,a,P) \in \mathcal{S} \times \mathcal{A} \times \mathcal{U}^\sigma(P^\phi)} \frac{\min \{d^{*,P}(s,a), \frac{1}{S}\}}{d^{\text{b},P^\phi}(s,a)}. \tag{260}$$

Regarding $\pi_\phi^*(\phi|s) = 1$ for all $s \in \mathcal{S}$ and the initial distribution $\rho(0) = 1$, for any $P \in \mathcal{U}^\sigma(P^\phi)$, we arrive at

$$d^{*,P}(0, \phi) = (1-\gamma)\rho(0)\pi_\phi^*(\phi|0) = (1-\gamma), \tag{261}$$

which holds by that the agent transits from state 0 to some other states at the first step and then will never go back to state 0. In addition, one has for any $P \in \mathcal{U}^\sigma(P^\phi)$,

$$\begin{aligned}
d^{*,P}(2, \phi) &= (1-\gamma)\rho(0)P(2|0, \phi) \sum_{t=0}^{\infty} \gamma^t (P(2|2, \phi))^t \\
&= (1-\gamma)P(2|0, \phi) \sum_{t=0}^{\infty} \gamma^t \stackrel{(i)}{\geq} \underline{p} \geq \frac{1}{\beta}
\end{aligned} \tag{262}$$

where (i) holds by (236) and the final inequality follows from (237). Armed with above facts, we observe that

$$\max_{(s,a,P) \in \mathcal{S} \times \mathcal{A} \times \mathcal{U}^\sigma(P^\phi)} \frac{\min \{d^{*,P}(s,a), \frac{1}{S}\}}{d^{\text{b},P^\phi}(s,a)} = \max_{s \in \{0,1,2\}, P \in \mathcal{U}^\sigma(P^\phi)} \frac{\min \{d^{*,P}(s, \phi), \frac{1}{S}\}}{d^{\text{b},P^\phi}(s, \phi)} \tag{263}$$

which follows from the properties of optimal policy in (259) and consequently $d^{*,P}(s) = d^{*,P}(s, \phi) = 0$ for all $s > 2$ and all $P \in \mathcal{U}^\sigma(P^\phi)$.

To continue, we control the term in states $\{0, 1, 2\}$ separately:

$$\max_{P \in \mathcal{U}^\sigma(P^\phi)} \frac{\min \{d^{*,P}(2, \phi), \frac{1}{S}\}}{d^{b,P^\phi}(2, \phi)} \stackrel{(i)}{=} \frac{1}{S d^{b,P^\phi}(2, \phi)} \stackrel{(ii)}{=} \frac{2}{S \mu(2)} = 2C, \quad (264)$$

$$\max_{P \in \mathcal{U}^\sigma(P^\phi)} \frac{\min \{d^{*,P}(0, \phi), \frac{1}{S}\}}{d^{b,P^\phi}(0, \phi)} \leq \frac{1}{S d^{b,P^\phi}(0, \phi)} \stackrel{(iii)}{=} \frac{2}{S \mu(0)} = 2C, \quad (265)$$

$$\max_{P \in \mathcal{U}^\sigma(P^\phi)} \frac{\min \{d^{*,P}(1, \phi), \frac{1}{S}\}}{d^{b,P^\phi}(1, \phi)} \leq \frac{1}{S d^{b,P^\phi}(1, \phi)} \stackrel{(iv)}{=} \frac{2}{S(1 - \frac{2}{CS})} \stackrel{(v)}{\leq} \frac{4}{S} \leq C, \quad (266)$$

where (i) holds by (262) by choosing $S \geq \beta$, (ii), (iii) and (iv) follow from the definition in (240) and (238), and (v) and the final inequality of (266) arise from the assumption in (239). Plugging in above result back into (263) directly completes the proof by

$$C_{\text{rob}}^* = \max_{(s,a,P) \in \mathcal{S} \times \mathcal{A} \times \mathcal{U}^\sigma(P^\phi)} \frac{\min \{d^{*,P}(s, a), \frac{1}{S}\}}{d^{b,P^\phi}(s, a)} = 2C. \quad (267)$$

F.3 PROOF OF AUXILIARY LEMMAS

F.3.1 PROOF OF LEMMA 11

We shall provide the proof to show that the operator $\hat{T}_{\text{pe}}^\sigma(\cdot)$ (cf. (174)) is a γ -contraction and the existence of the unique fixed point of $\hat{T}_{\text{pe}}^\sigma(\cdot)$ subsequently.

Before starting, suppose $Q, Q', Q_1, Q_2 \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ obey $Q(s, a), Q'(s, a), Q_1(s, a), Q_2(s, a) \in [0, \frac{1}{1-\gamma}]$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. Then we introduce the following notations:

$$\begin{aligned} \forall s \in \mathcal{S}: \quad V(s) &:= \max_a Q(s, a), \quad V'(s) := \max_a Q'(s, a), \\ V_1(s) &:= \max_a Q_1(s, a), \quad V_2(s) := \max_a Q_2(s, a). \end{aligned} \quad (268)$$

γ -contraction. We first show that $\hat{T}_{\text{pe}}^\sigma(\cdot)$ is a γ -contraction. Towards this, instead of $\hat{T}_{\text{pe}}^\sigma(\cdot)$, we consider a simpler operator $\tilde{T}_{\text{pe}}^\sigma(\cdot)$ firstly, defined as follows:

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \tilde{T}_{\text{pe}}^\sigma(Q)(s, a) = r(s, a) + \gamma \inf_{P \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P}V - b(s, a), \quad (269)$$

which consequently leads to

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \hat{T}_{\text{pe}}^\sigma(Q)(s, a) = \max \left\{ \tilde{T}_{\text{pe}}^\sigma(Q)(s, a), 0 \right\}. \quad (270)$$

With this in mind, we observe that

$$\begin{aligned} \left\| \tilde{T}_{\text{pe}}^\sigma(Q_1) - \tilde{T}_{\text{pe}}^\sigma(Q_2) \right\|_\infty &= \gamma \left\| \inf_{P \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P}V_1 - \inf_{P \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P}V_2 \right\|_\infty \stackrel{(i)}{\leq} \gamma \|V_1 - V_2\|_\infty \\ &\stackrel{(ii)}{=} \gamma \max_s \left| \max_a Q_1(s, a) - \max_a Q_2(s, a) \right| \\ &\leq \gamma \max_{(s,a)} |Q_1(s, a) - Q_2(s, a)| = \gamma \|Q_1 - Q_2\|_\infty \end{aligned} \quad (271)$$

where the first equality holds by applying the definition of $b(s, a)$ (cf. (176)) and (269), (i) follows from that the infimum operator is a 1-contraction w.r.t. $\|\cdot\|_\infty$ and $\|\mathcal{P}V_1 - \mathcal{P}V_2\|_\infty \leq \|V_1 - V_2\|_\infty$ for all $P \in \Delta(\mathcal{S})$, (ii) arises from the definitions in (268), and the last inequality is due to the maximum operator is also a 1-contraction w.r.t. $\|\cdot\|_\infty$.

Taking the above result with (270), we verify the desired assertion by

$$\left\| \hat{T}_{\text{pe}}^\sigma(Q_1) - \hat{T}_{\text{pe}}^\sigma(Q_2) \right\|_\infty \leq \left\| \tilde{T}_{\text{pe}}^\sigma(Q_1) - \tilde{T}_{\text{pe}}^\sigma(Q_2) \right\|_\infty \leq \gamma \|Q_1 - Q_2\|_\infty, \quad (272)$$

where the first inequality follows from the basic fact that the maximum operator is a 1-contraction w.r.t. $\|\cdot\|_\infty$.

Existence of the unique fixed-point. To continue, we shall firstly show that there exist at least one fixed-point of $\hat{T}_{\text{pe}}^\sigma(\cdot)$. Recalling the definition of $\hat{T}_{\text{pe}}^\sigma(\cdot)$ (cf. (174))

$$\hat{T}_{\text{pe}}^\sigma(Q)(s, a) = \max \left\{ r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P}V - b(s, a), 0 \right\}, \quad (273)$$

one has as long as $0 \leq Q \leq \frac{1}{1-\gamma} \cdot 1$, it is easily verified $0 \leq \hat{T}_{\text{pe}}^\sigma(Q) \leq \frac{1}{1-\gamma} \cdot 1$. Then, we construct the following sequence of Q-function recursively

$$Q^{(0)} = 0, \quad \text{and} \quad Q^{(t+1)} = \hat{T}_{\text{pe}}^\sigma(Q^{(t)}) \quad \text{for all } t \geq 0, \quad (274)$$

which mimic the iterations of our algorithm DRVI-LCB. As a result, the proof for the Banach fixed-point theorem (Agarwal et al., 2001, Theorem 1) gives that as $t \rightarrow \infty$, $Q^{(t)}$ converges to some point $Q^{(\infty)}$. It can also be verified that $0 \leq Q^{(\infty)} \leq \frac{1}{1-\gamma} \cdot 1$, which indicates the existence of the fixed points. Then, to prove the uniqueness of the fixed points of $\hat{T}_{\text{pe}}^\sigma(\cdot)$, we suppose that there exists another point Q' obeying $Q' = \hat{T}_{\text{pe}}^\sigma(Q')$. The definition of $\hat{T}_{\text{pe}}^\sigma(\cdot)$ directly gives $Q' \geq 0$ and if $\|Q'\|_\infty > \frac{1}{1-\gamma}$, then

$$\begin{aligned} \|Q'\|_\infty &= \left\| \hat{T}_{\text{pe}}^\sigma(Q') \right\|_\infty \leq \|r\|_\infty + \gamma \max_{(s,a)} \left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P}V' \right| \\ &\leq 1 + \gamma \|V'\|_\infty \leq 1 + \gamma \|Q'\|_\infty < (1 - \gamma) \|Q'\|_\infty + \gamma \|Q'\|_\infty = \|Q'\|_\infty \end{aligned} \quad (275)$$

gives contraction. Therefore, we have $0 \leq Q' \leq \frac{1}{1-\gamma} \cdot 1$, which yields

$$\|Q' - Q^{(\infty)}\|_\infty = \left\| \hat{T}_{\text{pe}}^\sigma(Q') - \hat{T}_{\text{pe}}^\sigma(Q^{(\infty)}) \right\|_\infty \leq \gamma \|Q' - Q^{(\infty)}\|_\infty. \quad (276)$$

However, (276) can't happen given $\gamma \in [\frac{1}{2}, 1)$, indicating the uniqueness of the fixed points of $\hat{T}_{\text{pe}}^\sigma(\cdot)$.

F.3.2 PROOF OF LEMMA 12

To begin with, considering any Q, Q' obeying $Q \leq Q', 0 \leq Q \leq \frac{1}{1-\gamma} \cdot 1$, and $0 \leq Q' \leq \frac{1}{1-\gamma} \cdot 1$, we observe that the operator $\hat{T}_{\text{pe}}^\sigma(\cdot)$ (cf. (174)) has the monotone non-decreasing property, namely,

$$\begin{aligned} \hat{T}_{\text{pe}}^\sigma(Q) &= \max \left\{ r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P}V - b(s, a), 0 \right\} \\ &= \max \left\{ r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P} \max_{a'} Q(\cdot, a') - b(s, a), 0 \right\} \\ &\leq \max \left\{ r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P} \max_{a'} Q'(\cdot, a') - b(s, a), 0 \right\} = \hat{T}_{\text{pe}}^\sigma(Q'). \end{aligned} \quad (277)$$

In addition, armed with (277) and the initial rule $\hat{Q}_0 = 0$, we also observe that the fixed-point $\hat{Q}_{\text{pe}}^{*,\sigma}$ of $\hat{T}_{\text{pe}}^\sigma(\cdot)$ obeys $0 \leq \hat{Q}_{\text{pe}}^{*,\sigma} \leq \frac{1}{1-\gamma} \cdot 1$. Consequently, we arrive at

$$\hat{Q}_1 = \hat{T}_{\text{pe}}^\sigma(\hat{Q}_0) \leq \hat{T}_{\text{pe}}^\sigma(\hat{Q}_{\text{pe}}^{*,\sigma}) = \hat{Q}_{\text{pe}}^{*,\sigma}. \quad (278)$$

Implementing the above result recursively gives

$$\text{for all } m \geq 0 : \quad \hat{Q}_m \leq \hat{Q}_{\text{pe}}^{*,\sigma}. \quad (279)$$

To continue, applying Lemma 11 yields that for any $m \geq 0$,

$$\|\hat{Q}_m - \hat{Q}_{\text{pe}}^{*,\sigma}\|_\infty = \left\| \hat{T}_{\text{pe}}^\sigma(\hat{Q}_{m-1}) - \hat{T}_{\text{pe}}^\sigma(\hat{Q}_{\text{pe}}^{*,\sigma}) \right\|_\infty \leq \gamma \|\hat{Q}_{m-1} - \hat{Q}_{\text{pe}}^{*,\sigma}\|_\infty \quad (280)$$

$$\leq \dots \leq \gamma^m \|\hat{Q}_0 - \hat{Q}_{\text{pe}}^{*,\sigma}\|_\infty = \gamma^m \|\hat{Q}_{\text{pe}}^{*,\sigma}\|_\infty \leq \frac{\gamma^m}{1-\gamma}, \quad (281)$$

where the last inequality holds by the fact $\|\hat{Q}_{\text{pe}}^{*,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$ (see Lemma 11). The final assertion can be directly achieved with the above result by observing

$$\|\hat{Q}_M - \hat{Q}_{\text{pe}}^{*,\sigma}\|_\infty \leq \frac{\gamma^M}{1-\gamma} \leq \frac{1}{\sigma N} \quad (282)$$

when $M \geq \frac{\log \frac{\sigma N}{1-\gamma}}{\log \frac{1}{\gamma}}$.

F.3.3 PROOF OF LEMMA 14

We first note that the second assertion in (199) is a counterpart of (48) which can be verified following the same argument in Appendix C.2.1 except the set of notations are adapted to the infinite-horizon case. Therefore, the rest of the proof will focus on verifying (198).

To begin with, we consider the situation when $N(s, a) = 0$. In this case, (198) can be easily verified that

$$\left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P}V - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^0)} \mathcal{P}V \right| \stackrel{(i)}{=} \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^0)} \mathcal{P}V \leq \|V\|_\infty \stackrel{(ii)}{\leq} \frac{1}{1-\gamma}, \quad (283)$$

where (i) follows from the fact $\hat{P}_{s,a}^0 = 0$ when $N(s, a) = 0$ (see (172)), and (ii) arises from the assumption $\|V\|_\infty \leq \frac{1}{1-\gamma}$. Consequently, in the remainder of the proof, we focus on verifying (198) when $N(s, a) > 0$.

Before continuing, we introduce a counterpart of the fact (47) in Lemma 8 as follows:

Lemma 18. *For all $(s, a) \in \mathcal{S} \times \mathcal{A}$ with $N(s, a) > 0$, consider any vector $V \in \mathbb{R}^{\mathcal{S}}$ independent of $\hat{P}_{s,a}^0$ obeying $\|V\|_\infty \leq \frac{1}{1-\gamma}$. With probability at least $1 - \delta$, one has*

$$\left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} \mathcal{P}V - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^0)} \mathcal{P}V \right| \leq \frac{c_b}{\sigma(1-\gamma)} \sqrt{\frac{\log(\frac{NS}{\delta})}{\hat{P}_{\min}(s, a)N(s, a)}}. \quad (284)$$

Proof. The proof follows from the same pipeline of the proof in Appendix C.2.2. The only difference is the upper bound on $\|V\|_\infty$ is $\frac{1}{1-\gamma}$ (as opposed to H), the union bound is taken over N (as opposed to KH), and some notations are exchanged to that of the infinite-horizon case. We omit the proof details for conciseness. \square

Armed with above point-wise results, we are now ready to derive the union bound over all \tilde{V} desired in Lemma 14 counting on a leave-one-out argument separated into the following steps.

Step 1: construction of auxiliary robust MDPs with state-absorbing nominal transitions. To begin with, we denote the empirical infinite-horizon robust MDP with the nominal transition kernel \hat{P}^0 as $\hat{\mathcal{M}}_{\text{rob}}$. Then, for each state s and each scalar $u \geq 0$, we can construct an auxiliary robust MDP $\hat{\mathcal{M}}_{\text{rob}}^{s,u}$ so that it is the same as $\hat{\mathcal{M}}_{\text{rob}}$ except the properties in state s . Specifically, the reward function of the auxiliary robust MDP $\hat{\mathcal{M}}_{\text{rob}}^{s,u}$ is denoted as $r^{s,u}$ which obeys

$$\begin{cases} r^{s,u}(s, a) = u & \text{for all } a \in \mathcal{A}, \\ r^{s,u}(\tilde{s}, a) = r(\tilde{s}, a) & \text{for all } (\tilde{s}, a) \in \mathcal{S} \times \mathcal{A} \text{ and } \tilde{s} \neq s. \end{cases} \quad (285)$$

In addition, the nominal transition kernel of $\hat{\mathcal{M}}_{\text{rob}}^{s,u}$ is denoted as $P^{s,u}$ such that

$$\begin{cases} P^{s,u}(s' | s, a) = \mathbb{1}(s' = s) & \text{for all } (s', a) \in \mathcal{S} \times \mathcal{A}, \\ P^{s,u}(\cdot | \tilde{s}, a) = \hat{P}^0(\cdot | \tilde{s}, a) & \text{for all } (\tilde{s}, a) \in \mathcal{S} \times \mathcal{A} \text{ and } \tilde{s} \neq s. \end{cases} \quad (286)$$

It can be observed that the nominal transition kernel $P^{s,u}$ of the auxiliary $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ drops all the randomness of $\widehat{P}_{s,a}^0$ for all $a \in \mathcal{A}$ in state s and makes s an absorbing state, while keeps other parts the same as \widehat{P}^0 .

With the robust MDP $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ in hand, we can define the corresponding penalty term for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ as follows

$$b^{s,u}(s, a) := \begin{cases} \min \left\{ \frac{c_b}{\sigma(1-\gamma)} \sqrt{\frac{\log \left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta} \right)}{P_{\min}^{s,u}(s, a)N(s, a)}} + \frac{4}{N\sigma(1-\gamma)}, \frac{1}{1-\gamma} \right\} + \frac{2}{\sigma N} & \text{if } N(s, a) > 0, \\ \frac{1}{1-\gamma} + \frac{2}{\sigma N} & \text{otherwise,} \end{cases} \quad (287)$$

where $P_{\min}^{s,u}$ is defined as the smallest positive state transition probability over the nominal kernel $P^{s,u}$ as follows:

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad P_{\min}^{s,u}(s, a) := \min_{s'} \left\{ P^{s,u}(s' | s, a) : P^{s,u}(s' | s, a) > 0 \right\}. \quad (288)$$

Armed with the penalty term, the pessimistic robust Bellman operator $\widehat{T}_{s,u}^\sigma(Q)(\cdot)$ w.r.t. $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ is defined as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \widehat{T}_{s,u}^\sigma(Q)(s, a) = \max \left\{ r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^{s,u})} \mathcal{P}V - b^{s,u}(s, a), 0 \right\}. \quad (289)$$

Step 2: verifying the relation between $\widehat{\mathcal{M}}_{\text{rob}}$ and the auxiliary robust MDP $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$. Recall that $\widehat{Q}_{\text{pe}}^{*,\sigma}$ is the unique fixed-point of operator $\widehat{T}_{\text{pe}}^\sigma(\cdot)$ with the corresponding value $\widehat{V}_{\text{pe}}^{*,\sigma}$. In particular, given a state s , we introduce a special reward

$$u^* := (1 - \gamma)\widehat{V}_{\text{pe}}^{*,\sigma}(s) + \min \left\{ \frac{c_b}{\sigma(1-\gamma)} \sqrt{\frac{\log \left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta} \right)}{P_{\min}^{s,u}(s, a)N(s, a)}} + \frac{4}{N\sigma(1-\gamma)}, \frac{1}{1-\gamma} \right\} + \frac{2}{\sigma N}. \quad (290)$$

With it in mind, we shall justify that there exists a fixed-point $\widehat{Q}_{s,u^*}^{*,\sigma}$ of the operator $\widehat{T}_{s,u^*}^\sigma(\cdot)$ whose corresponding value $\widehat{V}_{s,u^*}^{*,\sigma}$ is identical to $\widehat{V}_{\text{pe}}^{*,\sigma}$. Towards this, we shall show the facts in two different cases:

- **For state $s' \neq s$.** In this case, for any $s' \neq s$ and $a \in \mathcal{A}$, it can be verified that

$$\begin{aligned} & \max \left\{ r^{s,u^*}(s', a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s',a}^{s,u^*})} \mathcal{P}\widehat{V}_{\text{pe}}^{*,\sigma} - b^{s,u^*}(s', a), 0 \right\} \\ &= \max \left\{ r(s', a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s',a}^0)} \mathcal{P}\widehat{V}_{\text{pe}}^{*,\sigma} - b(s', a), 0 \right\} \\ &= \widehat{T}_{\text{pe}}^\sigma(\widehat{Q}_{\text{pe}}^{*,\sigma})(s', a) = \widehat{Q}_{\text{pe}}^{*,\sigma}(s', a), \end{aligned} \quad (291)$$

where the first identity follows from the definitions in (285) and (286), the penultimate equality arises from (174), and the final equality holds by that $\widehat{Q}_{\text{pe}}^{*,\sigma}$ is the fixed-point.

- **For state s .** In this case, for any u and $a \in \mathcal{A}$, observing that $P^{s,u}(s' | s, a)$ has only one positive entry equal to 1 (cf. (286)), applying (288) yields

$$P_{\min}^{s,u}(s, a) = 1. \quad (292)$$

Plugging above fact into (287) leads to

$$b^{s,u}(s, a) = \begin{cases} \min \left\{ \frac{c_b}{\sigma(1-\gamma)} \sqrt{\frac{\log \left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta} \right)}{N(s, a)}} + \frac{4}{N\sigma(1-\gamma)}, \frac{1}{1-\gamma} \right\} + \frac{2}{\sigma N} & \text{if } N(s, a) > 0, \\ \frac{1}{1-\gamma} & \text{otherwise} \end{cases} \quad (293)$$

for all $a \in \mathcal{A}$. As a result, we have for any $a \in \mathcal{A}$:

$$\begin{aligned} & \max \left\{ r^{s,u^*}(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\mathcal{P}_{s,a}^{s,u^*})} \mathcal{P} \widehat{V}_{\text{pe}}^{*,\sigma} - b^{s,u^*}(s, a), 0 \right\} \\ &= \max \left\{ u^* + \gamma \widehat{V}_{\text{pe}}^{*,\sigma}(s) - b^{s,u^*}(s, a), 0 \right\} \\ &\stackrel{(i)}{=} \max \left\{ (1 - \gamma) \widehat{V}_{\text{pe}}^{*,\sigma}(s) + \gamma \widehat{V}_{\text{pe}}^{*,\sigma}(s), 0 \right\} = \widehat{V}_{\text{pe}}^{*,\sigma}(s), \end{aligned} \quad (294)$$

where (i) follows from plugging in the definition of u^* in (290) and $b^{s,u^*}(s, a)$ in (293).

Summing up the above results, we observe that there exists a fixed point $\widehat{Q}_{s,u^*}^{*,\sigma}$ of the operator $\widehat{\mathcal{T}}_{s,u^*}^\sigma(\cdot)$ if we let

$$\begin{cases} \widehat{Q}_{s,u^*}^{*,\sigma}(s, a) = \widehat{V}_{\text{pe}}^{*,\sigma}(s) & \text{for all } a \in \mathcal{A} \\ \widehat{Q}_{s,u^*}^{*,\sigma}(s', a) = \widehat{Q}_{\text{pe}}^{*,\sigma}(s', a) & \text{for all } s' \neq s \text{ and } a \in \mathcal{A}. \end{cases} \quad (295)$$

Consequently, we already confirm the existence of a fixed point of the operator $\widehat{\mathcal{T}}_{s,u^*}^\sigma(\cdot)$. In addition, its corresponding value function $\widehat{V}_{s,u^*}^{*,\sigma}$ also coincides with $\widehat{V}_{\text{pe}}^{*,\sigma}$.

Step 3: building an ε -net for all rewards u . Before continuing, it is easily verified that the reward obeys

$$u^* \leq 1 + \min \left\{ \frac{c_b}{\sigma(1-\gamma)} \sqrt{\frac{\log \left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta} \right)}{\widehat{P}_{\min}^{s,u}(s, a)N(s, a)} + \frac{4}{\sigma N(1-\gamma)}}, \frac{1}{1-\gamma} \right\} + \frac{2}{\sigma N} \leq \frac{2}{\sigma} + \frac{2}{(1-\gamma)}. \quad (296)$$

As a result, we construct an ε -net (Vershynin, 2018) of the range $[0, \frac{2}{\sigma} + \frac{2}{(1-\gamma)}]$ with $\varepsilon = \frac{1}{\sigma N}$ as follows:

$$\mathcal{U}_\varepsilon := \left\{ \frac{i}{\sigma N} \mid 1 \leq i \leq \sigma N \left(\frac{2}{\sigma} + \frac{2}{(1-\gamma)} \right) \right\}. \quad (297)$$

Armed with this covering net \mathcal{U}_ε , we can construct an auxiliary robust MDP $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ and its corresponding pessimistic robust Bellman operator for each $u \in \mathcal{U}_\varepsilon$ (see Step 1). Following the same pipeline in the proof of Lemma 11 (cf. Appendix F.3.1), for each $u \in \mathcal{U}_\varepsilon$, it can be verified that there exists a unique fixed-point $\widehat{Q}_{s,u}^{*,\sigma}$ of the operator $\widehat{\mathcal{T}}_{s,u}^\sigma(\cdot)$ which satisfies $0 \leq \widehat{Q}_{s,u}^{*,\sigma} \leq \frac{1}{1-\gamma}$. Consequently, the corresponding value function also satisfies $\|\widehat{V}_{s,u}^{*,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$.

To continue, in view of the definitions in (285) and (286), we notice that for all $u \in \mathcal{U}_\varepsilon$, $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ is statistically independent from $\widehat{P}_{s,a}^0$, which indicates the independence between $\widehat{V}_{s,u}^{*,\sigma}$ and $\widehat{P}_{s,a}^0$. So invoking Lemma 18 and taking the union bound over all samples N and $u \in \mathcal{U}_\varepsilon$ give that, with probability at least $1 - \delta$,

$$\left| \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s,a}^0)} \mathcal{P} \widehat{V}_{s,u}^{*,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s,a}^0)} \mathcal{P} \widehat{V}_{s,u}^{*,\sigma} \right| \leq \frac{c_b}{\sigma(1-\gamma)} \sqrt{\frac{\log \left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta} \right)}{\widehat{P}_{\min}(s, a)N(s, a)}} \quad (298)$$

hold simultaneously for all $(s, a, u) \in \mathcal{S} \times \mathcal{A} \times \mathcal{U}_\varepsilon$ with $N(s, a) > 0$.

Step 4: implementing a covering argument. To continue, we shall control the gap between the value functions of the fixed-points of $\widehat{\mathcal{T}}_{\text{pe}}^\sigma(\cdot)$ and the auxiliary operator $\widehat{\mathcal{T}}_{s,u}^\sigma(\cdot)$, i.e., $\|\widehat{V}_{s,u}^{*,\sigma} - \widehat{V}_{\text{pe}}^{*,\sigma}\|_\infty$. First, recalling that $u^* \in [0, \frac{2}{\sigma} + \frac{2}{(1-\gamma)}]$ (see (290)), we can always find some $\tilde{u} \in \mathcal{U}_\varepsilon$ such that $|\tilde{u} - u^*| \leq \frac{1}{\sigma N}$. Consequently, plugging in the operator in (289) yields

$$\forall Q \in \mathbb{R}^{\mathcal{S}\mathcal{A}} : \quad \left\| \widehat{\mathcal{T}}_{s,\tilde{u}}^\sigma(Q) - \widehat{\mathcal{T}}_{s,u^*}^\sigma(Q) \right\|_\infty \stackrel{(i)}{\leq} |\tilde{u} - u^*| \leq \frac{1}{\sigma N}, \quad (299)$$

where (i) holds by $b^{s,\tilde{u}}(s, a) = b^{s,u^*}(s, a)$ for s (see (293)) and $b^{s,\tilde{u}}(s', a) = b^{s,u^*}(s', a) = b(s', a)$ for all $s' \neq s$.

With this in mind, we observe that the fixed-points obey

$$\begin{aligned} \left\| \hat{Q}_{s,\tilde{u}}^{*,\sigma} - \hat{Q}_{s,u^*}^{*,\sigma} \right\|_{\infty} &= \left\| \hat{T}_{s,\tilde{u}}^{\sigma}(\hat{Q}_{s,\tilde{u}}^{*,\sigma}) - \hat{T}_{s,u^*}^{\sigma}(\hat{Q}_{s,u^*}^{*,\sigma}) \right\|_{\infty} \\ &\leq \left\| \hat{T}_{s,\tilde{u}}^{\sigma}(\hat{Q}_{s,\tilde{u}}^{*,\sigma}) - \hat{T}_{s,\tilde{u}}^{\sigma}(\hat{Q}_{s,u^*}^{*,\sigma}) \right\|_{\infty} + \left\| \hat{T}_{s,\tilde{u}}^{\sigma}(\hat{Q}_{s,u^*}^{*,\sigma}) - \hat{T}_{s,u^*}^{\sigma}(\hat{Q}_{s,u^*}^{*,\sigma}) \right\|_{\infty} \\ &\leq \gamma \left\| \hat{Q}_{s,\tilde{u}}^{*,\sigma} - \hat{Q}_{s,u^*}^{*,\sigma} \right\|_{\infty} + \frac{1}{\sigma N}, \end{aligned} \quad (300)$$

which directly indicates that

$$\left\| \hat{Q}_{s,\tilde{u}}^{*,\sigma} - \hat{Q}_{s,u^*}^{*,\sigma} \right\|_{\infty} \leq \frac{1}{(1-\gamma)\sigma N} \quad (301)$$

and then

$$\left\| \hat{V}_{s,\tilde{u}}^{*,\sigma} - \hat{V}_{s,u^*}^{*,\sigma} \right\|_{\infty} \leq \left\| \hat{Q}_{s,\tilde{u}}^{*,\sigma} - \hat{Q}_{s,u^*}^{*,\sigma} \right\|_{\infty} \leq \frac{1}{(1-\gamma)\sigma N}. \quad (302)$$

Armed with above facts, invoking the identity between $\hat{V}_{pe}^{*,\sigma}$ and $\hat{V}_{s,u^*}^{*,\sigma}$ established in Step 2 gives

$$\begin{aligned} \left| \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\hat{P}_{s,a}^0)} \mathcal{P} \hat{V}_{pe}^{*,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s,a}^0)} \mathcal{P} \hat{V}_{pe}^{*,\sigma} \right| &= \left| \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\hat{P}_{s,a}^0)} \mathcal{P} \hat{V}_{s,u^*}^{*,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s,a}^0)} \mathcal{P} \hat{V}_{s,u^*}^{*,\sigma} \right| \\ &\stackrel{(i)}{\leq} \left| \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\hat{P}_{s,a}^0)} \mathcal{P} \hat{V}_{s,\tilde{u}}^{*,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s,a}^0)} \mathcal{P} \hat{V}_{s,\tilde{u}}^{*,\sigma} \right| \\ &\quad + \left| \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\hat{P}_{s,a}^0)} \mathcal{P} \hat{V}_{s,\tilde{u}}^{*,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\hat{P}_{s,a}^0)} \mathcal{P} \hat{V}_{s,u^*}^{*,\sigma} \right| + \left| \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s,a}^0)} \mathcal{P} \hat{V}_{s,\tilde{u}}^{*,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s,a}^0)} \mathcal{P} \hat{V}_{s,u^*}^{*,\sigma} \right| \\ &\stackrel{(ii)}{\leq} \left| \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\hat{P}_{s,a}^0)} \mathcal{P} \hat{V}_{s,\tilde{u}}^{*,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s,a}^0)} \mathcal{P} \hat{V}_{s,\tilde{u}}^{*,\sigma} \right| + \frac{2}{N\sigma(1-\gamma)} \\ &\leq \frac{c_b}{\sigma(1-\gamma)} \sqrt{\frac{\log\left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta}\right)}{\hat{P}_{\min}(s,a)N(s,a)}} + \frac{2}{N\sigma(1-\gamma)}, \end{aligned} \quad (303)$$

where (i) holds by applying the triangle inequality, (ii) arises from (302) and the basic fact that infimum operator is a 1-contraction w.r.t. $\|\cdot\|_{\infty}$, and the final inequality follows from (298).

Step 5: union bound for all \tilde{V} . Now we are positioned to show the union bound for all vector \tilde{V} obeying $\|\tilde{V} - \hat{V}_{pe}^{*,\sigma}\|_{\infty} \leq \frac{1}{\sigma N}$ and $\|\tilde{V}\|_{\infty} \leq \frac{1}{1-\gamma}$. For any \tilde{V} mentioned above, we invoke (303) and apply the triangle inequality to reach

$$\begin{aligned} \left| \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\hat{P}_{s,a}^0)} \mathcal{P} \tilde{V} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s,a}^0)} \mathcal{P} \tilde{V} \right| &\leq \left| \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\hat{P}_{s,a}^0)} \mathcal{P} \hat{V}_{pe}^{*,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s,a}^0)} \mathcal{P} \hat{V}_{pe}^{*,\sigma} \right| \\ &\quad + \left| \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\hat{P}_{s,a}^0)} \mathcal{P} \tilde{V} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\hat{P}_{s,a}^0)} \mathcal{P} \hat{V}_{pe}^{*,\sigma} \right| + \left| \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s,a}^0)} \mathcal{P} \tilde{V} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s,a}^0)} \mathcal{P} \hat{V}_{pe}^{*,\sigma} \right| \end{aligned} \quad (304)$$

$$\leq \frac{c_b}{\sigma(1-\gamma)} \sqrt{\frac{\log\left(\frac{2(1+\sigma)N^3S}{(1-\gamma)\delta}\right)}{\hat{P}_{\min}(s,a)N(s,a)}} + \frac{4}{N\sigma(1-\gamma)}. \quad (305)$$

Finally, we complete the proof by verifying that

$$\left| \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\hat{P}_{s,a}^0)} \mathcal{P} \tilde{V} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s,a}^0)} \mathcal{P} \tilde{V} \right| \leq \|\tilde{V}\|_{\infty} \leq \frac{1}{1-\gamma} \quad (306)$$

which holds by that the infimum operator is a 1-contraction w.r.t. $\|\cdot\|_{\infty}$ and the assumption $\|\tilde{V}\|_{\infty} \leq \frac{1}{1-\gamma}$.

F.3.4 PROOF OF (195)

For all $(s, a) \in \mathcal{C}^b$, one has

$$Nd^{b,P^0}(s, a) \stackrel{(i)}{\geq} \frac{c_1 d^{b,P^0}(s, a) \log(NS/\delta)}{d_{\min}^b P_{\min}^b} \stackrel{(ii)}{\geq} \frac{c_1 \log(NS/\delta)}{P_{\min}^b} \stackrel{(iii)}{\geq} \frac{c_1 \log(NS/\delta)}{P_{\min}(s, a)}, \quad (307)$$

where (i) follows from the condition (187), (ii) arises from the definition that $d_{\min}^b \leq d^{b,P^0}(s, a)$ for all $(s, a) \in \mathcal{C}^b$, and (iii) follows from the definition in (192).

Armed with above result, when c_1 is large enough, one has $\frac{2}{3} \log \frac{NS}{\delta} < \frac{Nd^{b,P^0}(s, a)}{12}$. Consequently, Lemma 13 tells us that with probability at least $1 - \delta$,

$$N(s, a) \geq \frac{Nd^{b,P^0}(s, a)}{12} \geq \frac{c_1 \log(NS/\delta)}{12P_{\min}(s, a)}. \quad (308)$$

Regarding the basic fact $x \leq -\log(1 - x)$ for all $x \in [0, 1]$, the last inequality of (195) can be verified by

$$\frac{c_1 \log(NS/\delta)}{12P_{\min}(s, a)} \geq -\frac{\log \frac{2NS}{\delta}}{\log(1 - P_{\min}(s, a))}. \quad (309)$$