

A. Proof of the preliminaries

Before moving forward, let us introduce some additional definitions and facts that will be useful throughout the appendix.

Kullback-Leibler (KL) divergence. First, for any two distributions P and Q , we denote by $\text{KL}(P \parallel Q)$ the Kullback-Leibler (KL) divergence of P and Q . Letting $\text{Ber}(p)$ be the Bernoulli distribution with mean p , we also introduce

$$\text{KL}(p \parallel q) := p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q} \quad \text{and} \quad \chi^2(p \parallel q) := \frac{(p-q)^2}{q} + \frac{(p-q)^2}{1-q} = \frac{(p-q)^2}{q(1-q)}, \quad (75)$$

which represent respectively the KL divergence and the χ^2 divergence of $\text{Ber}(p)$ from $\text{Ber}(q)$ (Tsybakov 2009).

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

LEMMA 5. Consider any $0 \leq V_1, V_2 \leq \frac{1}{1-\gamma}$ obeying $\|V_1 - V_2\|_\infty \leq x$ and any probability vector $P \in \Delta(\mathcal{S})$ (here $\Delta(\mathcal{S})$ represents the simplex over the state space \mathcal{S}), one has

$$|\text{Var}_P(V_1) - \text{Var}_P(V_2)| \leq \frac{2x}{(1-\gamma)}. \quad (76)$$

Proof of Lemma 5: It is immediate to check that

$$\begin{aligned} |\text{Var}_P(V_1) - \text{Var}_P(V_2)| &= |P(V_1 \circ V_1) - (PV_1) \circ (PV_1) - P(V_2 \circ V_2) + (PV_2) \circ (PV_2)| \\ &\leq |P(V_1 \circ V_1 - V_2 \circ V_2)| + |(PV_1 + PV_2)P(V_1 - V_2)| \\ &\leq 2\|V_1 + V_2\|_\infty \|V_1 - V_2\|_\infty \leq \frac{2x}{(1-\gamma)}. \end{aligned} \quad (77)$$

where the penultimate inequality holds by the triangle inequality.

A.1. Proof of Lemma 1 and Lemma 2

Proof of Lemma 1: To begin with, applying (Iyengar 2005, Lemma 4.3), the term of interest obeys

$$\inf_{P \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\mu \in \mathbb{R}^S, \mu \geq 0} \left\{ P(V - \mu) - \sigma \left(\max_{s'} \{V(s') - \mu(s')\} - \min_{s'} \{V(s') - \mu(s')\} \right) \right\}, \quad (78)$$

where $\mu(s')$ represents the s' -th entry of $\mu \in \mathbb{R}^S$. Denoting μ^* as the optimal dual solution, taking $\alpha = \max_{s'} \{V(s') - \mu^*(s')\}$, it is easily verified that μ^* obeys

$$\mu^*(s) = \begin{cases} V(s) - \alpha, & \text{if } V(s) > \alpha \\ 0, & \text{otherwise.} \end{cases} \quad (79)$$

Therefore, (78) can be solved by optimizing α as below (Iyengar 2005, Lemma 4.3):

$$\inf_{P \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P[V]_\alpha - \sigma \left(\alpha - \min_{s'} [V]_\alpha(s') \right) \right\}. \quad (80)$$

Proof of Lemma 2 Due to strong duality (Iyengar 2005, Lemma 4.2), it holds that

$$\inf_{\mathcal{P} \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\mu \in \mathbb{R}^S, \mu \geq 0} \left\{ P(V - \mu) - \sqrt{\sigma \text{Var}_P(V - \mu)} \right\}, \quad (81)$$

and the optimal μ^* obeys

$$\mu^*(s) = \begin{cases} V(s) - \alpha, & \text{if } V(s) > \alpha \\ 0, & \text{otherwise.} \end{cases} \quad (82)$$

for some $\alpha \in [\min_s V(s), \max_s V(s)]$. As a result, solving (81) is equivalent to optimizing the scalar α as below:

$$\inf_{\mathcal{P} \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P[V]_\alpha - \sqrt{\sigma \text{Var}_P([V]_\alpha)} \right\}. \quad (83)$$

A.2. Proof of Lemma 4

Applying the γ -contraction property in Lemma 3 directly yields that for any $t \geq 0$,

$$\begin{aligned} \|\widehat{Q}_t - \widehat{Q}^{*,\sigma}\|_\infty &= \|\widehat{\mathcal{T}}^\sigma(\widehat{Q}_{t-1}) - \widehat{\mathcal{T}}^\sigma(\widehat{Q}^{*,\sigma})\|_\infty \leq \gamma \|\widehat{Q}_{t-1} - \widehat{Q}^{*,\sigma}\|_\infty \\ &\leq \dots \leq \gamma^t \|\widehat{Q}_0 - \widehat{Q}^{*,\sigma}\|_\infty = \gamma^t \|\widehat{Q}^{*,\sigma}\|_\infty \leq \frac{\gamma^t}{1-\gamma}, \end{aligned}$$

where the last inequality holds by the fact $\|\widehat{Q}^{*,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$ (see Lemma 3). In addition,

$$\|\widehat{V}_t - \widehat{V}^{*,\sigma}\|_\infty = \max_{s \in \mathcal{S}} \left\| \max_{a \in \mathcal{A}} \widehat{Q}_t(s, a) - \max_{a \in \mathcal{A}} \widehat{Q}^{*,\sigma}(s, a) \right\|_\infty \leq \|\widehat{Q}_t - \widehat{Q}^{*,\sigma}\|_\infty \leq \frac{\gamma^t}{1-\gamma},$$

where the penultimate inequality holds by the maximum operator is 1-Lipschitz. This completes the proof of (46).

We now move to establish (47). Note that there exists at least one state $s_0 \in \mathcal{S}$ that is associated with the maximum of the value gap, i.e.,

$$\|\widehat{V}^{*,\sigma} - \widehat{V}^{\widehat{\pi}^*,\sigma}\|_\infty = \widehat{V}^{*,\sigma}(s_0) - \widehat{V}^{\widehat{\pi}^*,\sigma}(s_0) \geq \widehat{V}^{*,\sigma}(s) - \widehat{V}^{\widehat{\pi}^*,\sigma}(s), \quad \forall s \in \mathcal{S}.$$

Recall $\widehat{\pi}^*$ is the optimal robust policy for the empirical RMDP $\widehat{\mathcal{M}}_{\text{rob}}$. For convenience, we denote $a_1 = \widehat{\pi}^*(s_0)$ and $a_2 = \widehat{\pi}(s_0)$. Then, since $\widehat{\pi}$ is the greedy policy w.r.t. \widehat{Q}_T , one has

$$r(s_0, a_1) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s_0, a_1}^0)} \mathcal{P}\widehat{V}_{T-1} = \widehat{Q}_T(s_0, a_1) \leq \widehat{Q}_T(s_0, a_2) = r(s_0, a_2) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s_0, a_2}^0)} \mathcal{P}\widehat{V}_{T-1}. \quad (84)$$

Recalling the notation in (39), the above fact and (47) altogether yield

$$\begin{aligned} r(s_0, a_1) + \gamma \widehat{P}_{s_0, a_1}^{\widehat{V}_{T-1}} \left(\widehat{V}^{*,\sigma} - \varepsilon_{\text{opt}} \mathbf{1} \right) &\leq r(s_0, a_1) + \gamma \widehat{P}_{s_0, a_1}^{\widehat{V}_{T-1}} \widehat{V}_{T-1} \\ &\leq r(s_0, a_2) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s_0, a_2}^0)} \mathcal{P}\widehat{V}_{T-1} \\ &\stackrel{(i)}{\leq} r(s_0, a_2) + \gamma \widehat{P}_{s_0, a_2}^{\widehat{V}^{\widehat{\pi}^*,\sigma}} \widehat{V}_{T-1} \\ &\leq r(s_0, a_2) + \gamma \widehat{P}_{s_0, a_2}^{\widehat{V}^{\widehat{\pi}^*,\sigma}} \left(\widehat{V}^{*,\sigma} + \varepsilon_{\text{opt}} \mathbf{1} \right), \end{aligned} \quad (85)$$

where (i) follows from the optimality criteria. The term of interest can be controlled as

$$\begin{aligned}
& \|\widehat{V}^{*,\sigma} - \widehat{V}^{\widehat{\pi},\sigma}\|_{\infty} \\
&= \widehat{V}^{*,\sigma}(s_0) - \widehat{V}^{\widehat{\pi},\sigma}(s_0) \\
&= r(s_0, a_1) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\widehat{P}_{s_0, a_1}^0)} \mathcal{P}\widehat{V}^{*,\sigma} - \left(r(s_0, a_2) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\widehat{P}_{s_0, a_2}^0)} \mathcal{P}\widehat{V}^{\widehat{\pi},\sigma} \right) \\
&= r(s_0, a_1) - r(s_0, a_2) + \gamma \left(\inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\widehat{P}_{s_0, a_1}^0)} \mathcal{P}\widehat{V}^{*,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\widehat{P}_{s_0, a_2}^0)} \mathcal{P}\widehat{V}^{\widehat{\pi},\sigma} \right) \\
&\stackrel{(i)}{\leq} 2\gamma\varepsilon_{\text{opt}} + \gamma \left(\widehat{P}_{s_0, a_2}^{\widehat{V}^{\widehat{\pi},\sigma}} \widehat{V}^{*,\sigma} - \widehat{P}_{s_0, a_1}^{\widehat{V}^{*,\sigma}} \widehat{V}^{*,\sigma} + \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\widehat{P}_{s_0, a_1}^0)} \mathcal{P}\widehat{V}^{*,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\widehat{P}_{s_0, a_2}^0)} \mathcal{P}\widehat{V}^{\widehat{\pi},\sigma} \right) \\
&= 2\gamma\varepsilon_{\text{opt}} + \gamma \left(\widehat{P}_{s_0, a_2}^{\widehat{V}^{\widehat{\pi},\sigma}} \widehat{V}^{*,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\widehat{P}_{s_0, a_2}^0)} \mathcal{P}\widehat{V}^{\widehat{\pi},\sigma} \right) + \gamma \left(\inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\widehat{P}_{s_0, a_1}^0)} \mathcal{P}\widehat{V}^{*,\sigma} - \widehat{P}_{s_0, a_1}^{\widehat{V}^{*,\sigma}} \widehat{V}^{*,\sigma} \right) \\
&\stackrel{(ii)}{\leq} 2\gamma\varepsilon_{\text{opt}} + \gamma \widehat{P}_{s_0, a_2}^{\widehat{V}^{\widehat{\pi},\sigma}} (\widehat{V}^{*,\sigma} - \widehat{V}^{\widehat{\pi},\sigma}) + \gamma \left(\widehat{P}_{s_0, a_1}^{\widehat{V}^{*,\sigma}} \widehat{V}^{*,\sigma} - \widehat{P}_{s_0, a_1}^{\widehat{V}^{\widehat{\pi},\sigma}} \widehat{V}^{\widehat{\pi},\sigma} \right) \\
&\leq 2\gamma\varepsilon_{\text{opt}} + \gamma \|\widehat{V}^{*,\sigma} - \widehat{V}^{\widehat{\pi},\sigma}\|_{\infty}, \tag{86}
\end{aligned}$$

where (i) holds by plugging in (85), and (ii) follows from $\inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(\widehat{P}_{s_0, a_1}^0)} \mathcal{P}\widehat{V}^{*,\sigma} \leq \widehat{P}_{s_0, a_1}^{\widehat{V}^{*,\sigma}}$ for any $\mathcal{P} \in \mathcal{U}^{\sigma}(\widehat{P}_{s_0, a_1}^0)$. Rearranging (86) leads to

$$\|\widehat{V}^{*,\sigma} - \widehat{V}^{\widehat{\pi},\sigma}\|_{\infty} \leq \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}.$$

B. Proof of the upper bound with TV distance: Theorem 1

B.1. Technical lemmas

We begin with a key lemma that is new and distinguishes robust MDPs with TV distance from standard MDPs, which plays a critical role in obtaining the sample complexity upper bound in Theorem 1. This lemma concerns the dynamic range of the robust value function $V^{\pi, \sigma}$ (cf. (5)) for any fixed policy π , which produces tighter control than that in standard MDP (cf. $\frac{1}{1-\gamma}$) when σ is large. The proof is deferred to Appendix B.3.1

LEMMA 6. For any nominal transition kernel $P \in \mathbb{R}^{S \times A \times S}$, any fixed uncertainty level σ , and any policy π , its corresponding robust value function $V^{\pi, \sigma}$ (cf. (5)) satisfies

$$\max_{s \in \mathcal{S}} V^{\pi, \sigma}(s) - \min_{s \in \mathcal{S}} V^{\pi, \sigma}(s) \leq \frac{1}{\gamma \max\{1 - \gamma, \sigma\}}.$$

With the above lemma in hand, we introduce the following lemma that is useful throughout this section, whose proof is postponed to Appendix B.3.2

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

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

REMARK 4. Compared to the results in (Li et al. 2024b, Lemma 11), Lemma 7 provides an improved upper bound, expressed in terms of $\max_s V^{\pi, P}(s) - \min_s V^{\pi, P}(s)$ rather than $\|V^{\pi, P}\|_\infty$.

B.2. Proof of Theorem 1

Throughout this section, for any transition kernel P , the uncertainty set is taken as (see (9))

$$\mathcal{U}^\sigma(P) := \mathcal{U}_{\text{TV}}^\sigma(P) = \otimes \mathcal{U}_{\text{TV}}^\sigma(P_{s,a}), \quad \mathcal{U}_{\text{TV}}^\sigma(P_{s,a}) := \left\{ P'_{s,a} \in \Delta(\mathcal{S}) : \frac{1}{2} \|P'_{s,a} - P_{s,a}\|_1 \leq \sigma \right\}. \quad (87)$$

To control the two main terms in (49), respectively, we first recall (54) which holds for any uncertainty set:

$$\|\widehat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma}\|_\infty \leq \gamma \max \left\{ \left\| \left(I - \gamma \widehat{P}^{\pi^*, \widehat{V}} \right)^{-1} \left(\widehat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \right\|_\infty, \left\| \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \left(\widehat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \right\|_\infty \right\}. \quad (88)$$

B.2.1. Controlling $\|\widehat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma}\|_\infty$. We shall focus on controlling the two terms on the right hand side of the above results separately.

Step 1: controlling $\|\widehat{V}^{\pi^, \sigma} - V^{\pi^*, \sigma}\|_\infty$: bounding the first term in (54).* To control the two terms in (54), we first introduce the following lemma whose proof is postponed to Appendix B.3.3

LEMMA 8. Consider any $\delta \in (0, 1)$. Setting $N \geq \log\left(\frac{18SAN}{\delta}\right)$, with probability at least $1 - \delta$, one has

$$\begin{aligned} \left| \widehat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right| &\leq 2 \sqrt{\frac{\log\left(\frac{18SAN}{\delta}\right)}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{*, \sigma})} + \frac{\log\left(\frac{18SAN}{\delta}\right)}{N(1-\gamma)} 1 \\ &\leq 3 \sqrt{\frac{\log\left(\frac{18SAN}{\delta}\right)}{(1-\gamma)^2 N}} 1, \end{aligned} \quad (89)$$

where $\text{Var}_{P^{\pi^*}}(V^{*, \sigma})$ is defined in (37).

Armed with the above lemma, now we control the first term on the right hand side of (54) as follows:

$$\begin{aligned} &\left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \left(\widehat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \\ &\stackrel{(i)}{\leq} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \left\| \widehat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right\|_\infty \\ &\stackrel{(ii)}{\leq} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \left(2 \sqrt{\frac{\log\left(\frac{18SAN}{\delta}\right)}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{*, \sigma})} + \frac{\log\left(\frac{18SAN}{\delta}\right)}{N(1-\gamma)} 1 \right) \\ &\leq \frac{\log\left(\frac{18SAN}{\delta}\right)}{N(1-\gamma)} \underbrace{\left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} 1 + 2 \sqrt{\frac{\log\left(\frac{18SAN}{\delta}\right)}{N}} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, V}}(V^{*, \sigma})}}_{=: \mathcal{C}_1} \\ &\quad + 2 \underbrace{\sqrt{\frac{\log\left(\frac{18SAN}{\delta}\right)}{N}} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \sqrt{\left| \text{Var}_{\widehat{P}^{\pi^*}}(V^{*, \sigma}) - \text{Var}_{\widehat{P}^{\pi^*, V}}(V^{*, \sigma}) \right|}}_{=: \mathcal{C}_2} \end{aligned}$$

$$+ 2 \underbrace{\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} (I - \gamma \hat{P}^{\pi^*, V})^{-1} \left(\sqrt{\text{Var}_{P^{\pi^*}}(V^{*, \sigma})} - \sqrt{\text{Var}_{\hat{P}^{\pi^*}}(V^{*, \sigma})} \right)}_{=: \mathcal{C}_3}, \quad (90)$$

where (i) holds by $(I - \gamma \hat{P}^{\pi^*, V})^{-1} \geq 0$, (ii) follows from Lemma 8, and the last inequality arise from

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

by applying the triangle inequality.

To continue, observing that each row of $\hat{P}^{\pi^*, V}$ is a probability distribution obeying that the sum is 1, we arrive at

$$(I - \gamma \hat{P}^{\pi^*, V})^{-1} \mathbf{1} = \left(I + \sum_{t=1}^{\infty} \gamma^t (\hat{P}^{\pi^*, V})^t \right) \mathbf{1} = \frac{1}{1 - \gamma} \mathbf{1}. \quad (91)$$

Armed with this fact, we shall control the other three terms $\mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3$ in (90) separately.

- Consider \mathcal{C}_1 . We first introduce the following lemma, whose proof is postponed to Appendix B.3.4.

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

$$(I - \gamma \hat{P}^{\pi^*, V})^{-1} \sqrt{\text{Var}_{\hat{P}^{\pi^*, V}}(V^{*, \sigma})} \leq 4 \sqrt{\frac{\left(1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, \sigma\}}} \mathbf{1} \leq 4 \sqrt{\frac{\left(1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right)}{\gamma^3 (1-\gamma)^3}} \mathbf{1}.$$

Applying Lemma 9 and inserting back to (90) leads to

$$\begin{aligned} \mathcal{C}_1 &= 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} (I - \gamma \hat{P}^{\pi^*, V})^{-1} \sqrt{\text{Var}_{\hat{P}^{\pi^*, V}}(V^{*, \sigma})} \\ &\leq 8 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, \sigma\} N} \left(1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right)} \mathbf{1}. \end{aligned} \quad (92)$$

- Consider \mathcal{C}_2 . First, denote $V' := V^{*, \sigma} - \min_{s' \in \mathcal{S}} V^{*, \sigma}(s') \mathbf{1}$, by Lemma 6, it follows that

$$0 \leq V' \leq \frac{1}{\gamma \max\{1-\gamma, \sigma\}} \mathbf{1}. \quad (93)$$

Then, we have for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, and $P_{s,a} \in \Delta(\mathcal{S})$, and $\tilde{P}_{s,a} \in \mathcal{U}^\sigma(P_{s,a})$:

$$\begin{aligned} \left| \text{Var}_{\tilde{P}_{s,a}}(V^{*, \sigma}) - \text{Var}_{P_{s,a}}(V^{*, \sigma}) \right| &= \left| \text{Var}_{\tilde{P}_{s,a}}(V') - \text{Var}_{P_{s,a}}(V') \right| \\ &\leq \|\tilde{P}_{s,a} - P_{s,a}\|_1 \|V'\|_\infty^2 \\ &\leq \frac{2\sigma}{\gamma^2 (\max\{1-\gamma, \sigma\})^2} \leq \frac{2}{\gamma^2 \max\{1-\gamma, \sigma\}}. \end{aligned} \quad (94)$$

Applying the above relation we obtain

$$\begin{aligned}
\mathcal{C}_2 &= 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \sqrt{\left|\text{Var}_{\hat{P}^{\pi^*}}(V^{*, \sigma}) - \text{Var}_{\hat{P}^{\pi^*, V}}(V^{*, \sigma})\right|} \\
&= 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \sqrt{\left|\Pi^{\pi^*} (\text{Var}_{\hat{P}^0}(V^{*, \sigma}) - \text{Var}_{\hat{P}^{\pi^*, V}}(V^{*, \sigma}))\right|} \\
&\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \sqrt{\|\text{Var}_{\hat{P}^0}(V^{*, \sigma}) - \text{Var}_{\hat{P}^{\pi^*, V}}(V^{*, \sigma})\|_{\infty}} \\
&\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \sqrt{\frac{2}{\gamma^2 \max\{1 - \gamma, \sigma\}}} = 2\sqrt{\frac{2 \log(\frac{18SAN}{\delta})}{\gamma^2 (1 - \gamma)^2 \max\{1 - \gamma, \sigma\} N}} 1,
\end{aligned} \tag{95}$$

where the last equality uses $\left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} 1 = \frac{1}{1 - \gamma}$ (cf. (91)).

- Consider \mathcal{C}_3 . The following lemma plays an important role.

LEMMA 10. (Panaganti and Kalathil 2022 Lemma 6) Consider any $\delta \in (0, 1)$. For any fixed policy π and fixed value vector $V \in \mathbb{R}^S$, one has with probability at least $1 - \delta$,

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

Applying Lemma 10 with $\pi = \pi^*$ and $V = V^{*, \sigma}$ leads to

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

which can be plugged in (90) to verify

$$\begin{aligned}
\mathcal{C}_3 &= 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \left(\sqrt{\text{Var}_{P^{\pi^*}}(V^{*, \sigma})} - \sqrt{\text{Var}_{\hat{P}^{\pi^*}}(V^{*, \sigma})} \right) \\
&\leq \frac{4}{(1 - \gamma)} \frac{\log(\frac{SAN}{\delta}) \|V^{*, \sigma}\|_{\infty}}{N} 1 \leq \frac{4 \log(\frac{18SAN}{\delta})}{(1 - \gamma)^2 N} 1,
\end{aligned} \tag{96}$$

where the last line uses $\left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} 1 = \frac{1}{1 - \gamma}$ (cf. (91)).

Finally, inserting the results of \mathcal{C}_1 in (92), \mathcal{C}_2 in (95), \mathcal{C}_3 in (96), and (91) back into (90) gives

$$\begin{aligned}
&\left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \left(\hat{P}^{\pi^*, V} V^{\pi^*, \sigma} - P^{\pi^*, V} V^{\pi^*, \sigma}\right) \\
&\leq 8\sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^3 (1 - \gamma)^2 \max\{1 - \gamma, \sigma\} N}} \left(1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1 - \gamma)^2 N}}\right) 1 \\
&\quad + 2\sqrt{\frac{2 \log(\frac{18SAN}{\delta})}{\gamma^2 (1 - \gamma)^2 \max\{1 - \gamma, \sigma\} N}} 1 + \frac{4 \log(\frac{18SAN}{\delta})}{(1 - \gamma)^2 N} 1 + \frac{\log(\frac{18SAN}{\delta})}{N(1 - \gamma)^2} 1
\end{aligned}$$

$$\begin{aligned}
&\leq 10 \sqrt{\frac{2 \log(\frac{18SAN}{\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N} \left(1 + \sqrt{\frac{\log(\frac{SAN}{\delta})}{(1-\gamma)^2N}}\right) 1 + \frac{5 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2N} 1} \\
&\leq 160 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N} 1 + \frac{5 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2N} 1}, \tag{97}
\end{aligned}$$

where the last inequality holds by the fact $\gamma \geq \frac{1}{4}$ and letting $N \geq \frac{\log(\frac{SAN}{\delta})}{(1-\gamma)^2}$.

Step 2: controlling $\|\widehat{V}^{\pi^, \sigma} - V^{\pi^*, \sigma}\|_\infty$:* bounding the second term in (54). To proceed, applying Lemma 8 on the second term of the right hand side of (54) leads to

$$\begin{aligned}
&\left(I - \gamma \widehat{P}^{\pi^*, \widehat{V}}\right)^{-1} \left(\widehat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma}\right) \\
&\leq 2 \left(I - \gamma \widehat{P}^{\pi^*, \widehat{V}}\right)^{-1} \left(\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{\pi^*, \sigma})} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} 1\right) \\
&\leq \underbrace{\frac{2 \log(\frac{18SAN}{\delta})}{N(1-\gamma)} \left(I - \gamma \widehat{P}^{\pi^*, \widehat{V}}\right)^{-1} 1 + 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{P}^{\pi^*, \widehat{V}}\right)^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, \widehat{V}}}(\widehat{V}^{\pi^*, \sigma})}}_{=: \mathcal{C}_4} \\
&\quad + \underbrace{2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{P}^{\pi^*, \widehat{V}}\right)^{-1} \left(\sqrt{\text{Var}_{\widehat{P}^{\pi^*, \widehat{V}}}(V^{\pi^*, \sigma} - \widehat{V}^{\pi^*, \sigma})}\right)}_{=: \mathcal{C}_5} \\
&\quad + \underbrace{2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{P}^{\pi^*, \widehat{V}}\right)^{-1} \left(\sqrt{\left|\text{Var}_{\widehat{P}^{\pi^*}}(V^{\pi^*, \sigma}) - \text{Var}_{\widehat{P}^{\pi^*, \widehat{V}}}(V^{\pi^*, \sigma})\right|}\right)}_{=: \mathcal{C}_6} \\
&\quad + \underbrace{2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{P}^{\pi^*, \widehat{V}}\right)^{-1} \left(\sqrt{\text{Var}_{P^{\pi^*}}(V^{\pi^*, \sigma})} - \sqrt{\text{Var}_{\widehat{P}^{\pi^*}}(V^{\pi^*, \sigma})}\right)}_{=: \mathcal{C}_7}. \tag{98}
\end{aligned}$$

We now bound the above terms separately.

• Applying Lemma 7 with $P = \widehat{P}^{\pi^*, \widehat{V}}$, $\pi = \pi^*$ and taking $V = \widehat{V}^{\pi^*, \sigma}$ which obeys $\widehat{V}^{\pi^*, \sigma} = r_{\pi^*} + \gamma \widehat{P}^{\pi^*, \widehat{V}} \widehat{V}^{\pi^*, \sigma}$, and in view of (91), the term \mathcal{C}_4 in (98) can be controlled as follows:

$$\begin{aligned}
\mathcal{C}_4 &= 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{P}^{\pi^*, \widehat{V}}\right)^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, \widehat{V}}}(\widehat{V}^{\pi^*, \sigma})} \\
&\leq 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\frac{8(\max_s \widehat{V}^{\pi^*, \sigma}(s) - \min_s \widehat{V}^{\pi^*, \sigma}(s))}{\gamma^2(1-\gamma)^2}} 1 \\
&\leq 8 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1, \tag{99}
\end{aligned}$$

where the last inequality holds by applying Lemma 6

- To continue, considering \mathcal{C}_5 , we directly observe that (in view of (91))

$$\begin{aligned} \mathcal{C}_5 &= 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, \hat{V}} \right)^{-1} \sqrt{\text{Var}_{\hat{P}^{\pi^*, \hat{V}}}(V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma})} \\ &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma} \right\|_{\infty} 1. \end{aligned} \quad (100)$$

- Then, it is easily verified that \mathcal{C}_6 can be controlled similarly to (95) as follows:

$$\mathcal{C}_6 \leq 2\sqrt{\frac{2\log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1. \quad (101)$$

- Similarly, \mathcal{C}_7 can be controlled the same as (96) shown below:

$$\mathcal{C}_7 \leq \frac{4\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1. \quad (102)$$

Combining the results in (99), (100), (101), and (102) and inserting back to (98) leads to

$$\begin{aligned} \left(I - \gamma \hat{P}^{\pi^*, \hat{V}} \right)^{-1} \left(\hat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) &\leq 8\sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 \\ &\quad + 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma} \right\|_{\infty} 1 + 2\sqrt{\frac{2\log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 + \frac{4\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1 \\ &\leq 80\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 + 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma} \right\|_{\infty} 1 + \frac{4\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1, \end{aligned} \quad (103)$$

where the last inequality follows from the assumption $\gamma \geq \frac{1}{4}$.

Finally, inserting (97) and (103) back to (54) yields

$$\begin{aligned} &\left\| \hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma} \right\|_{\infty} \\ &\leq \max \left\{ 160\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{5\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}, \right. \\ &\quad \left. 80\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma} \right\|_{\infty} + \frac{4\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} \right\} \\ &\leq 160\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{8\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}, \end{aligned} \quad (104)$$

where the last inequality holds by taking $N \geq \frac{16\log(\frac{SAN}{\delta})}{(1-\gamma)^2}$.

B.2.2. Controlling $\left\| \hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma} \right\|_{\infty}$ Recall the bound in (55) which holds for any uncertainty set:

$$\begin{aligned} \left\| \hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma} \right\|_{\infty} &\leq \gamma \max \left\{ \left\| \left(I - \gamma \underline{P}^{\hat{\pi}, V} \right)^{-1} \left(\hat{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right) \right\|_{\infty}, \right. \\ &\quad \left. \left\| \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \left(\hat{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right) \right\|_{\infty} \right\}. \end{aligned} \quad (105)$$

Step 3: controlling $\|\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma}\|_\infty$: bounding the first term in (55). To begin with, we introduce the following lemma which controls the main term on the right hand side of (55), which is proved in Appendix B.3.5

LEMMA 11. Consider any $\delta \in (0, 1)$. Taking $N \geq \log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)$, with probability at least $1 - \delta$, one has

$$\begin{aligned} \left| \widehat{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right| &\leq 2\sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma})} 1 + \frac{8\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N(1-\gamma)} 1 + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} 1 \\ &\leq 10\sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{(1-\gamma)^2 N}} 1 + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} 1. \end{aligned} \quad (106)$$

With Lemma 11 in hand, we have

$$\begin{aligned} &\left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \left(\widehat{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right) \\ &\stackrel{(i)}{\leq} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \left| \widehat{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right| \\ &\leq 2\sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\text{Var}_{P^{\widehat{\pi}}}(\widehat{V}^{*,\sigma})} + \left(I - \gamma P_Q^{\widehat{\pi},V^{\widehat{\pi}}} \right)^{-1} \left(\frac{8\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N(1-\gamma)} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \right) 1 \\ &\stackrel{(ii)}{\leq} \left(\frac{8\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N(1-\gamma)^2} + \frac{2\gamma\varepsilon_{\text{opt}}}{(1-\gamma)^2} \right) 1 + \underbrace{2\sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})}}_{=: \mathcal{D}_1} \\ &\quad + \underbrace{2\sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\left| \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{*,\sigma}) - \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma}) \right|}}_{=: \mathcal{D}_2} \\ &\quad + \underbrace{2\sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\left| \text{Var}_{P^{\widehat{\pi}}}(\widehat{V}^{*,\sigma}) - \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{*,\sigma}) \right|}}_{=: \mathcal{D}_3}, \end{aligned} \quad (107)$$

where (i) and (ii) hold by the fact that each row of $(1-\gamma) \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1}$ is a probability vector that falls into $\Delta(\mathcal{S})$.

The remainder of the proof will focus on controlling the three terms in (107) separately.

- For \mathcal{D}_1 , we introduce the following lemma, whose proof is postponed to Appendix B.3.6

LEMMA 12. Consider any $\delta \in (0, 1)$. Taking $N \geq \frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{(1-\gamma)^2}$ and $\varepsilon_{\text{opt}} \leq \frac{1-\gamma}{\gamma}$, one has with probability at least $1 - \delta$,

$$\left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})} \leq 6\sqrt{\frac{1}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}}} 1 \leq 6\sqrt{\frac{1}{(1-\gamma)^3 \gamma^2}} 1.$$

Applying Lemma [12](#) and [91](#) to [107](#) leads to

$$\begin{aligned} \mathcal{D}_1 &= 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}\right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}}(\hat{V}^{\hat{\pi}, \sigma})} \\ &\leq 12\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1. \end{aligned} \quad (108)$$

- Applying Lemma [5](#) with $\|\hat{V}^{*,\sigma} - \hat{V}^{\hat{\pi},\sigma}\|_\infty \leq \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}$ and [91](#), \mathcal{D}_2 can be controlled as

$$\begin{aligned} \mathcal{D}_2 &= 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}\right)^{-1} \sqrt{\left| \text{Var}_{\underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}}(\hat{V}^{*,\sigma}) - \text{Var}_{\underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}}(\hat{V}^{\hat{\pi}, \sigma}) \right|} \\ &\leq 4\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}\right)^{-1} \frac{\sqrt{\gamma\varepsilon_{\text{opt}}}}{1-\gamma} \leq 4\sqrt{\frac{\gamma\varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1. \end{aligned} \quad (109)$$

- \mathcal{D}_3 can be controlled similarly to \mathcal{C}_2 in [95](#) as follows:

$$\begin{aligned} \mathcal{D}_3 &= 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}\right)^{-1} \sqrt{\left| \text{Var}_{P^{\hat{\pi}}}(\hat{V}^{*,\sigma}) - \text{Var}_{\underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}}(\hat{V}^{*,\sigma}) \right|} \\ &\leq 4\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}\right)^{-1} \sqrt{\frac{1}{\gamma^2 \max\{1-\gamma, \sigma\}}} 1 \leq 4\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 \end{aligned} \quad (110)$$

Finally, summing up the results in [108](#), [109](#), and [110](#) and inserting them back to [107](#) yields: taking $N \geq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2}$ and $\varepsilon_{\text{opt}} \leq \frac{1-\gamma}{\gamma}$, with probability at least $1 - \delta$,

$$\begin{aligned} &\left(I - \gamma \underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}\right)^{-1} \left(\hat{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}\hat{V}^{\hat{\pi}, \sigma}\right) \\ &\leq \left(\frac{8\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} + \frac{2\gamma\varepsilon_{\text{opt}}}{(1-\gamma)^2}\right) 1 + 12\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 \\ &\quad + 4\sqrt{\frac{\gamma\varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1 + 4\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 \\ &\leq 16\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 + \frac{14\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} 1, \end{aligned} \quad (111)$$

where the last inequality holds by taking $\varepsilon_{\text{opt}} \leq \min\left\{\frac{1-\gamma}{\gamma}, \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}\right\} = \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}$.

Step 4: controlling $\|\hat{V}^{\hat{\pi},\sigma} - V^{\hat{\pi},\sigma}\|_\infty$: bounding the second term in [55](#). Towards this, applying Lemma [11](#) leads to

$$\begin{aligned} &\left(I - \gamma \underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}\right)^{-1} \left(\hat{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}\hat{V}^{\hat{\pi}, \sigma}\right) \leq \left(I - \gamma \underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}\right)^{-1} \left|\hat{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}\hat{V}^{\hat{\pi}, \sigma}\right| \\ &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}\right)^{-1} \sqrt{\text{Var}_{P^{\hat{\pi}}}(\hat{V}^{*,\sigma})} + \left(I - \gamma \underline{P}^{\hat{P}^{\hat{\pi}, \hat{V}}}\right)^{-1} \left(\frac{8\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}\right) 1 \end{aligned}$$

$$\begin{aligned}
&\leq \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} + \frac{2\gamma\varepsilon_{\text{opt}}}{(1-\gamma)^2} \right) 1 + 2 \underbrace{\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} (I - \gamma \underline{P}^{\hat{\pi},V})^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi},V}}(V^{\hat{\pi},\sigma})}}_{=: \mathcal{D}_4} \\
&\quad + 2 \underbrace{\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} (I - \gamma \underline{P}^{\hat{\pi},V})^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi},V}}(\widehat{V}^{\hat{\pi},\sigma} - V^{\hat{\pi},\sigma})}}_{=: \mathcal{D}_5} \\
&\quad + 2 \underbrace{\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} (I - \gamma \underline{P}^{\hat{\pi},\widehat{V}})^{-1} \sqrt{|\text{Var}_{\underline{P}^{\hat{\pi},V}}(\widehat{V}^{*,\sigma}) - \text{Var}_{\underline{P}^{\hat{\pi},V}}(\widehat{V}^{\hat{\pi},\sigma})|}}_{=: \mathcal{D}_6} \\
&\quad + 2 \underbrace{\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} (I - \gamma \underline{P}^{\hat{\pi},\widehat{V}})^{-1} \sqrt{|\text{Var}_{P^{\hat{\pi}}}(\widehat{V}^{*,\sigma}) - \text{Var}_{\underline{P}^{\hat{\pi},V}}(\widehat{V}^{*,\sigma})|}}_{=: \mathcal{D}_7}. \tag{112}
\end{aligned}$$

We shall bound each of the terms separately.

- Applying Lemma 7 with $P = \underline{P}^{\hat{\pi},V}$, $\pi = \hat{\pi}$, and taking $V = V^{\hat{\pi},\sigma}$ which obeys $V^{\hat{\pi},\sigma} = r_{\hat{\pi}} + \gamma \underline{P}^{\hat{\pi},V} V^{\hat{\pi},\sigma}$, the term \mathcal{D}_4 can be controlled similarly to (99) as follows:

$$\mathcal{D}_4 \leq 8 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1. \tag{113}$$

- For \mathcal{D}_5 , it is observed that

$$\begin{aligned}
\mathcal{D}_5 &= 2 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} (I - \gamma \underline{P}^{\hat{\pi},V})^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi},V}}(\widehat{V}^{\hat{\pi},\sigma} - V^{\hat{\pi},\sigma})} \\
&\leq 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\hat{\pi},\sigma} - \widehat{V}^{\hat{\pi},\sigma} \right\|_{\infty} 1. \tag{114}
\end{aligned}$$

- Next, observing that \mathcal{D}_6 and \mathcal{D}_7 are almost the same as the terms \mathcal{D}_2 (controlled in (109)) and \mathcal{D}_3 (controlled in (110)) in (107), it is easily verified that they can be controlled as follows

$$\mathcal{D}_6 \leq 4 \sqrt{\frac{\gamma\varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1, \quad \mathcal{D}_7 \leq 4 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1. \tag{115}$$

Then inserting the results in (113), (114), and (115) back to (112) leads to

$$\begin{aligned}
&(I - \gamma \underline{P}^{\hat{\pi},V})^{-1} (\widehat{\underline{P}}^{\hat{\pi},\widehat{V}} \widehat{V}^{\hat{\pi},\sigma} - \underline{P}^{\hat{\pi},\widehat{V}} \widehat{V}^{\hat{\pi},\sigma}) \\
&\leq \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} + \frac{2\gamma\varepsilon_{\text{opt}}}{(1-\gamma)^2} \right) 1 + 8 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1 \\
&\quad + 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\hat{\pi},\sigma} - \widehat{V}^{\hat{\pi},\sigma} \right\|_{\infty} 1 + 4 \sqrt{\frac{\gamma\varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1 + 4 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1 \\
&\leq 12 \sqrt{\frac{2 \log(\frac{8SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1 + \frac{14 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} 1 + 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\hat{\pi},\sigma} - \widehat{V}^{\hat{\pi},\sigma} \right\|_{\infty} 1, \tag{116}
\end{aligned}$$

where the last inequality holds by letting $\varepsilon_{\text{opt}} \leq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}$, which directly satisfies $\varepsilon_{\text{opt}} \leq \frac{1-\gamma}{\gamma}$ by letting $N \geq \frac{\log(\frac{54SAN^2}{\delta})}{1-\gamma}$.

Finally, inserting (111) and (116) back to (55) yields: taking $\varepsilon_{\text{opt}} \leq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}$ and $N \geq \frac{16 \log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2}$, with probability at least $1 - \delta$, one has

$$\begin{aligned} & \left\| \widehat{V}^{\widehat{\pi}, \sigma} - V^{\widehat{\pi}, \sigma} \right\|_{\infty} \\ & \leq \max \left\{ 16 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} + \frac{14 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2}, \right. \\ & \quad \left. 12 \sqrt{\frac{2 \log(\frac{8SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} + \frac{14 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} + 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\widehat{\pi}, \sigma} - \widehat{V}^{\widehat{\pi}, \sigma} \right\|_{\infty} \right\} \\ & \leq 24 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} + \frac{28 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2}. \end{aligned} \quad (117)$$

Step 5: Summing up the results. Summing up the results in (104) and (117) and inserting back to (49) complete the proof as follows: taking $\varepsilon_{\text{opt}} \leq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}$ and $N \geq \frac{16 \log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2}$, with probability at least $1 - \delta$,

$$\begin{aligned} \left\| V^{*, \sigma} - V^{\widehat{\pi}, \sigma} \right\|_{\infty} & \leq \left\| V^{*, \sigma} - \widehat{V}^{\pi^*, \sigma} \right\|_{\infty} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} + \left\| \widehat{V}^{\widehat{\pi}, \sigma} - V^{\widehat{\pi}, \sigma} \right\|_{\infty} \\ & \leq \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} + 160 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} + \frac{8 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} \\ & \quad + 24 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} + \frac{28 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} \\ & \leq 184 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} + \frac{36 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} \\ & \leq 1508 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}}, \end{aligned} \quad (118)$$

where the last inequality holds by $\gamma \geq \frac{1}{4}$ and $N \geq \frac{16 \log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2}$.

B.3. Proof of the auxiliary lemmas for Theorem 1

B.3.1. Proof of Lemma 6 To begin, note that there at least exists one state s_0 for any $V^{\pi, \sigma}$ such that $V^{\pi, \sigma}(s_0) = \min_{s \in \mathcal{S}} V^{\pi, \sigma}(s)$. With this in mind, for any policy π , one has by the definition in (5) and the Bellman's equation (7a),

$$\begin{aligned} \max_{s \in \mathcal{S}} V^{\pi, \sigma}(s) & = \max_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi(\cdot | s)} \left[r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s, a})} \mathcal{P} V^{\pi, \sigma} \right] \\ & \leq \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} \left(1 + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s, a})} \mathcal{P} V^{\pi, \sigma} \right), \end{aligned}$$

where the second line holds since the reward function $r(s, a) \in [0, 1]$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. To continue, note that for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, there exists some $\tilde{P}_{s,a} \in \mathbb{R}^{\mathcal{S}}$ constructed by reducing the values of some elements of $P_{s,a}$ to obey $P_{s,a} \geq \tilde{P}_{s,a} \geq 0$ and $\sum_{s'} (P_{s,a}(s') - \tilde{P}_{s,a}(s')) = \sigma$. This implies $\tilde{P}_{s,a} + \sigma e_{s_0}^\top \in \mathcal{U}^\sigma(P_{s,a})$, where e_{s_0} is the standard basis vector supported on s_0 , since $\frac{1}{2} \|\tilde{P}_{s,a} + \sigma e_{s_0}^\top - P_{s,a}\|_1 \leq \frac{1}{2} \|\tilde{P}_{s,a} - P_{s,a}\|_1 + \frac{\sigma}{2} = \sigma$. Consequently,

$$\begin{aligned} \inf_{P \in \mathcal{U}^\sigma(P_{s,a})} \mathcal{P}V^{\pi,\sigma} &\leq \left(\tilde{P}_{s,a} + \sigma e_{s_0}^\top \right) V^{\pi,\sigma} \leq \|\tilde{P}_{s,a}\|_1 \|V^{\pi,\sigma}\|_\infty + \sigma V^{\pi,\sigma}(s_0) \\ &\leq (1 - \sigma) \max_{s \in \mathcal{S}} V^{\pi,\sigma}(s) + \sigma \min_{s \in \mathcal{S}} V^{\pi,\sigma}(s), \end{aligned} \quad (119)$$

where the second inequality holds by $\|\tilde{P}_{s,a}\|_1 = \sum_{s'} \tilde{P}_{s,a}(s') = -\sum_{s'} (P_{s,a}(s') - \tilde{P}_{s,a}(s')) + \sum_{s'} P_{s,a}(s') = 1 - \sigma$. Plugging this back to the previous relation gives

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

which, by rearranging terms, immediately yields

$$\begin{aligned} \max_{s \in \mathcal{S}} V^{\pi,\sigma}(s) &\leq \frac{1 + \gamma\sigma \min_{s \in \mathcal{S}} V^{\pi,\sigma}(s)}{1 - \gamma(1 - \sigma)} \\ &\leq \frac{1}{(1 - \gamma) + \gamma\sigma} + \min_{s \in \mathcal{S}} V^{\pi,\sigma}(s) \leq \frac{1}{\gamma \max\{1 - \gamma, \sigma\}} + \min_{s \in \mathcal{S}} V^{\pi,\sigma}(s). \end{aligned}$$

B.3.2. Proof of Lemma 7 Observing that each row of P_π belongs to $\Delta(\mathcal{S})$, it can be directly verified that each row of $(1 - \gamma)(I - \gamma P_\pi)^{-1}$ falls into $\Delta(\mathcal{S})$. As a result,

$$\begin{aligned} (I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} &= \frac{1}{1 - \gamma} (1 - \gamma)(I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} \\ &\stackrel{(i)}{\leq} \frac{1}{1 - \gamma} \sqrt{(1 - \gamma)(I - \gamma P_\pi)^{-1} \text{Var}_{P_\pi}(V^{\pi,P})} \\ &= \sqrt{\frac{1}{1 - \gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \text{Var}_{P_\pi}(V^{\pi,P})}, \end{aligned} \quad (120)$$

where (i) holds by Jensen's inequality.

To continue, denoting the minimum value of V as $V_{\min} = \min_{s \in \mathcal{S}} V^{\pi,P}(s)$ and $V' := V^{\pi,P} - V_{\min} \mathbf{1}$. We control $\text{Var}_{P_\pi}(V^{\pi,P})$ as follows:

$$\begin{aligned} &\text{Var}_{P_\pi}(V^{\pi,P}) \\ &\stackrel{(i)}{=} \text{Var}_{P_\pi}(V') = P_\pi(V' \circ V') - (P_\pi V') \circ (P_\pi V') \\ &\stackrel{(ii)}{=} P_\pi(V' \circ V') - \frac{1}{\gamma^2} (V' - r_\pi + (1 - \gamma)V_{\min} \mathbf{1}) \circ (V' - r_\pi + (1 - \gamma)V_{\min} \mathbf{1}) \\ &= P_\pi(V' \circ V') - \frac{1}{\gamma^2} V' \circ V' + \frac{2}{\gamma^2} V' \circ (r_\pi - (1 - \gamma)V_{\min} \mathbf{1}) \\ &\quad - \frac{1}{\gamma^2} (r_\pi - (1 - \gamma)V_{\min} \mathbf{1}) \circ (r_\pi - (1 - \gamma)V_{\min} \mathbf{1}) \\ &\leq P_\pi(V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_\infty \mathbf{1}, \end{aligned} \quad (121)$$

where (i) holds by the fact that $\text{Var}_{P_\pi}(V^{\pi,P} - b1) = \text{Var}_{P_\pi}(V^{\pi,P})$ for any scalar b and $V^{\pi,P} \in \mathbb{R}^S$, (ii) follows from $V' = r_\pi + \gamma P_\pi V^{\pi,P} - V_{\min}1 = r_\pi - (1 - \gamma)V_{\min}1 + \gamma P_\pi V'$, and the last line arises from $\frac{1}{\gamma^2}V' \circ V' \geq \frac{1}{\gamma}V' \circ V'$ and $\|r_\pi - (1 - \gamma)V_{\min}1\|_\infty \leq 1$. Plugging (121) back to (120) leads to

$$\begin{aligned}
(I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} &\leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \left(P_\pi(V' \circ V') - \frac{1}{\gamma}V' \circ V' + \frac{2}{\gamma^2}\|V'\|_\infty 1 \right)} \\
&\stackrel{(i)}{\leq} \sqrt{\frac{1}{1-\gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \left(P_\pi(V' \circ V') - \frac{1}{\gamma}V' \circ V' \right) \right|} + \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \frac{2}{\gamma^2}\|V'\|_\infty 1} \\
&\leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\left| \left(\sum_{t=0}^{\infty} \gamma^t (P_\pi)^{t+1} - \sum_{t=0}^{\infty} \gamma^{t-1} (P_\pi)^t \right) (V' \circ V') \right|} + \sqrt{\frac{2\|V'\|_\infty 1}{\gamma^2(1-\gamma)^2}} \\
&\stackrel{(ii)}{\leq} \sqrt{\frac{\|V'\|_\infty^2 1}{\gamma(1-\gamma)}} + \sqrt{\frac{2\|V'\|_\infty 1}{\gamma^2(1-\gamma)^2}} \\
&\leq \sqrt{\frac{8\|V'\|_\infty 1}{\gamma^2(1-\gamma)^2}}, \tag{122}
\end{aligned}$$

where (i) holds by the triangle inequality, (ii) holds by following recursion, and the last inequality holds by $\|V'\|_\infty \leq \frac{1}{1-\gamma}$.

B.3.3. Proof of Lemma 8

Step 1: controlling the point-wise concentration. We first consider a more general term w.r.t. any fixed (independent from \hat{P}^0) value vector V obeying $0 \leq V \leq \frac{1}{1-\gamma}1$ and any policy π . Invoking Lemma 1 leads to that for any $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned}
\left| \hat{P}_{s,a}^{\pi,V} V - P_{s,a}^{\pi,V} V \right| &\leq \left| \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ \hat{P}_{s,a}^0 [V]_\alpha - \sigma \left(\alpha - \min_{s'} [V]_\alpha(s') \right) \right\} \right. \\
&\quad \left. - \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P_{s,a}^0 [V]_\alpha - w\sigma \left(\alpha - \min_{s'} [V]_\alpha(s') \right) \right\} \right| \\
&\leq \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \underbrace{\left| \left(P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [V]_\alpha \right|}_{=: g_{s,a}(\alpha, V)}, \tag{123}
\end{aligned}$$

where the last inequality holds by that the maximum operator is 1-Lipschitz.

Then for a fixed α and any vector V that is independent with \hat{P}^0 , using the Bernstein's inequality, one has with probability at least $1 - \delta$,

$$\begin{aligned}
g_{s,a}(\alpha, V) &= \left| \left(P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [V]_\alpha \right| \leq \sqrt{\frac{2 \log(\frac{2}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} + \frac{2 \log(\frac{2}{\delta})}{3N(1-\gamma)} \\
&\leq \sqrt{\frac{2 \log(\frac{2}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{2 \log(\frac{2}{\delta})}{3N(1-\gamma)}. \tag{124}
\end{aligned}$$

Step 2: deriving the uniform concentration. To obtain the union bound, we first notice that $g_{s,a}(\alpha, V)$ is 1-Lipschitz w.r.t. α for any V obeying $\|V\|_\infty \leq \frac{1}{1-\gamma}$. In addition, we can construct an ε_1 -net N_{ε_1} over $[0, \frac{1}{1-\gamma}]$ whose size satisfies $|N_{\varepsilon_1}| \leq \frac{3}{\varepsilon_1(1-\gamma)}$ (Vershynin 2018). By the union bound and (124), it holds with probability at least $1 - \frac{\delta}{SA}$ that for all $\alpha \in N_{\varepsilon_1}$,

$$g_{s,a}(\alpha, V) \leq \sqrt{\frac{2\log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{2\log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{3N(1-\gamma)}. \quad (125)$$

Combined with (123), it yields that,

$$\begin{aligned} \left| \widehat{P}_{s,a}^{\pi, V} V - P_{s,a}^{\pi, V} V \right| &\leq \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [V]_\alpha \right| \\ &\stackrel{(i)}{\leq} \varepsilon_1 + \sup_{\alpha \in N_{\varepsilon_1}} \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [V]_\alpha \right| \\ &\stackrel{(ii)}{\leq} \varepsilon_1 + \sqrt{\frac{2\log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{2\log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{3N(1-\gamma)} \end{aligned} \quad (126)$$

$$\begin{aligned} &\stackrel{(iii)}{\leq} \sqrt{\frac{2\log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{\log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N(1-\gamma)} \\ &\stackrel{(iv)}{\leq} 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} \end{aligned} \quad (127)$$

$$\begin{aligned} &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \|V\|_\infty + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} \\ &\leq 3\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \end{aligned} \quad (128)$$

where (i) follows from that the optimal α^* falls into the ε_1 -ball centered around some point inside N_{ε_1} and $g_{s,a}(\alpha, V)$ is 1-Lipschitz, (ii) holds by (125), (iii) arises from taking $\varepsilon_1 = \frac{\log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{3N(1-\gamma)}$, (iv) is verified by $|N_{\varepsilon_1}| \leq \frac{3}{\varepsilon_1(1-\gamma)} \leq 9N$, and the last inequality is due to the fact $\|V^{*,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$ and letting $N \geq \log(\frac{18SAN}{\delta})$.

To continue, applying (127) and (128) with $\pi = \pi^*$ and $V = V^{*,\sigma}$ (independent with \widehat{P}^0) and taking the union bound over $(s, a) \in \mathcal{S} \times \mathcal{A}$ gives that with probability at least $1 - \delta$, it holds simultaneously for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ that

$$\begin{aligned} \left| \widehat{P}_{s,a}^{\pi^*, V} V^{*,\sigma} - P_{s,a}^{\pi^*, V} V^{*,\sigma} \right| &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V^{*,\sigma})} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} \\ &\leq 3\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}}. \end{aligned} \quad (129)$$

By converting (129) to the matrix form, one has with probability at least $1 - \delta$,

$$\begin{aligned} \left| \widehat{P}^{\pi^*, V} V^{\pi^*, \sigma} - P^{\pi^*, V} V^{\pi^*, \sigma} \right| &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{*,\sigma})} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} \mathbf{1} \\ &\leq 3\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \mathbf{1}. \end{aligned} \quad (130)$$

B.3.4. Proof of Lemma 9 Following the same argument as (120), it follows

$$\left(I - \gamma \widehat{\underline{P}}^{\pi^*, V}\right)^{-1} \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*, V}}(V^{*, \sigma})} = \sqrt{\frac{1}{1 - \gamma} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left(\widehat{\underline{P}}^{\pi^*, V}\right)^t \text{Var}_{\widehat{\underline{P}}^{\pi^*, V}}(V^{*, \sigma})}}. \quad (131)$$

To continue, we first focus on controlling $\text{Var}_{\widehat{\underline{P}}^{\pi^*, V}}(V^{*, \sigma})$. Towards this, denoting the minimum value of $V^{*, \sigma}$ as $V_{\min} := \min_{s \in \mathcal{S}} V^{*, \sigma}(s)$ and $V' := V^{*, \sigma} - V_{\min} \mathbf{1}$, we arrive at (see the robust Bellman's consistency equation in (45))

$$\begin{aligned} V' &= V^{*, \sigma} - V_{\min} \mathbf{1} = r_{\pi^*} + \gamma \underline{P}^{\pi^*, V} V^{*, \sigma} - V_{\min} \mathbf{1} \\ &= r_{\pi^*} + \gamma \widehat{\underline{P}}^{\pi^*, V} V^{*, \sigma} + \gamma \left(\underline{P}^{\pi^*, V} - \widehat{\underline{P}}^{\pi^*, V}\right) V^{*, \sigma} - V_{\min} \mathbf{1} \\ &= r_{\pi^*} - (1 - \gamma) V_{\min} \mathbf{1} + \gamma \widehat{\underline{P}}^{\pi^*, V} V' + \gamma \left(\underline{P}^{\pi^*, V} - \widehat{\underline{P}}^{\pi^*, V}\right) V^{*, \sigma} \\ &= r'_{\pi^*} + \gamma \widehat{\underline{P}}^{\pi^*, V} V' + \gamma \left(\underline{P}^{\pi^*, V} - \widehat{\underline{P}}^{\pi^*, V}\right) V^{*, \sigma}, \end{aligned} \quad (132)$$

where the last line holds by letting $r'_{\pi^*} := r_{\pi^*} - (1 - \gamma) V_{\min} \mathbf{1} \leq r_{\pi^*}$. With the above fact in hand, we control $\text{Var}_{\widehat{\underline{P}}^{\pi^*, V}}(V^{*, \sigma})$ as follows:

$$\begin{aligned} \text{Var}_{\widehat{\underline{P}}^{\pi^*, V}}(V^{*, \sigma}) &\stackrel{(i)}{=} \text{Var}_{\widehat{\underline{P}}^{\pi^*, V}}(V') = \widehat{\underline{P}}^{\pi^*, V} (V' \circ V') - \left(\widehat{\underline{P}}^{\pi^*, V} V'\right) \circ \left(\widehat{\underline{P}}^{\pi^*, V} V'\right) \\ &\stackrel{(ii)}{=} \widehat{\underline{P}}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma^2} \left(V' - r'_{\pi^*} - \gamma \left(\underline{P}^{\pi^*, V} - \widehat{\underline{P}}^{\pi^*, V}\right) V^{*, \sigma}\right)^{\circ 2} \\ &= \widehat{\underline{P}}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma^2} V' \circ V' + \frac{2}{\gamma^2} V' \circ \left(r'_{\pi^*} + \gamma \left(\underline{P}^{\pi^*, V} - \widehat{\underline{P}}^{\pi^*, V}\right) V^{*, \sigma}\right) \\ &\quad - \frac{1}{\gamma^2} \left(r'_{\pi^*} + \gamma \left(\underline{P}^{\pi^*, V} - \widehat{\underline{P}}^{\pi^*, V}\right) V^{*, \sigma}\right)^{\circ 2} \\ &\stackrel{(iii)}{\leq} \widehat{\underline{P}}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} \mathbf{1} + \frac{2}{\gamma} \|V'\|_{\infty} \left| \left(\underline{P}^{\pi^*, V} - \widehat{\underline{P}}^{\pi^*, V}\right) V^{*, \sigma} \right| \end{aligned} \quad (133)$$

$$\leq \widehat{\underline{P}}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} \mathbf{1} + \frac{6}{\gamma} \|V'\|_{\infty} \sqrt{\frac{\log\left(\frac{18SAN}{\delta}\right)}{(1 - \gamma)^2 N}} \mathbf{1}, \quad (134)$$

where (i) holds by the fact that $\text{Var}_{P_{\pi}}(V - b \mathbf{1}) = \text{Var}_{P_{\pi}}(V)$ for any scalar b and $V \in \mathbb{R}^{\mathcal{S}}$, (ii) follows from (132), (iii) arises from $\frac{1}{\gamma^2} V' \circ V' \geq \frac{1}{\gamma} V' \circ V'$ and $-1 \leq r_{\pi^*} - (1 - \gamma) V_{\min} \mathbf{1} = r'_{\pi^*} \leq r_{\pi^*} \leq 1$, and the last inequality holds by Lemma 8.

Plugging (134) into (131) leads to

$$\begin{aligned} &\left(I - \gamma \widehat{\underline{P}}^{\pi^*, V}\right)^{-1} \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*, V}}(V^{*, \sigma})} \\ &\leq \sqrt{\frac{1}{1 - \gamma} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left(\widehat{\underline{P}}^{\pi^*, V}\right)^t \left(\widehat{\underline{P}}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} \mathbf{1} + \frac{6}{\gamma} \|V'\|_{\infty} \sqrt{\frac{\log\left(\frac{18SAN}{\delta}\right)}{(1 - \gamma)^2 N}} \mathbf{1}\right)}} \\ &\stackrel{(i)}{\leq} \sqrt{\frac{1}{1 - \gamma} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left(\widehat{\underline{P}}^{\pi^*, V}\right)^t \left(\widehat{\underline{P}}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V'\right)}} \end{aligned}$$

$$\begin{aligned}
& + \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left(\widehat{P}^{\pi^*, V} \right)^t \left(\frac{2}{\gamma^2} \|V'\|_{\infty} 1 + \frac{6}{\gamma} \|V'\|_{\infty} \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} 1 \right)} \\
& \leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t \left(\widehat{P}^{\pi^*, V} \right)^t \left[\widehat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' \right] \right|} + \sqrt{\frac{\left(2 + 6 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right) \|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1,
\end{aligned} \tag{135}$$

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

$$\begin{aligned}
& \left| \sum_{t=0}^{\infty} \gamma^t \left(\widehat{P}^{\pi^*, V} \right)^t \left(\widehat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' \right) \right| \\
& = \left| \left(\sum_{t=0}^{\infty} \gamma^t \left(\widehat{P}^{\pi^*, V} \right)^{t+1} - \sum_{t=0}^{\infty} \gamma^{t-1} \left(\widehat{P}^{\pi^*, V} \right)^t \right) (V' \circ V') \right| \leq \frac{1}{\gamma} \|V'\|_{\infty}^2 1
\end{aligned} \tag{136}$$

by recursion. Inserting (136) back to (135) leads to

$$\begin{aligned}
& \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, V}}(V^{*, \sigma})} \\
& \leq \sqrt{\frac{\|V'\|_{\infty}^2}{\gamma(1-\gamma)}} 1 + 3 \sqrt{\frac{\left(1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right) \|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1 \\
& \leq 4 \sqrt{\frac{\left(1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right) \|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1 \leq 4 \sqrt{\frac{\left(1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, \sigma\}}} 1 \leq 4 \sqrt{\frac{\left(1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right)}{\gamma^3 (1-\gamma)^3}} 1,
\end{aligned} \tag{137}$$

where the penultimate inequality follows from applying Lemma 6 with $P = P^0$ and $\pi = \pi^*$:

$$\|V'\|_{\infty} = \max_{s \in \mathcal{S}} V^{*, \sigma}(s) - \min_{s \in \mathcal{S}} V^{*, \sigma}(s) \leq \frac{1}{\gamma \max\{1-\gamma, \sigma\}}.$$

B.3.5. Proof of Lemma 11 To begin with, for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, invoking the results in (123), we have

$$\begin{aligned}
& \left| \widehat{P}_{s,a}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - P_{s,a}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right| \leq \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha} \right| \\
& \stackrel{(i)}{\leq} \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left(\left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{*, \sigma}]_{\alpha} \right| + \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) \left([\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha} - [\widehat{V}^{*, \sigma}]_{\alpha} \right) \right| \right) \\
& \leq \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left(\left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{*, \sigma}]_{\alpha} \right| + \left\| P_{s,a}^0 - \widehat{P}_{s,a}^0 \right\|_1 \left\| [\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha} - [\widehat{V}^{*, \sigma}]_{\alpha} \right\|_{\infty} \right) \\
& \stackrel{(ii)}{\leq} \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{*, \sigma}]_{\alpha} \right| + 2 \left\| \widehat{V}^{\widehat{\pi}, \sigma} - \widehat{V}^{*, \sigma} \right\|_{\infty} \\
& \stackrel{(iii)}{\leq} \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{*, \sigma}]_{\alpha} \right| + \frac{2\gamma \varepsilon_{\text{opt}}}{1-\gamma},
\end{aligned} \tag{138}$$

where (i) holds by the triangle inequality, and (ii) follows from $\|P_{s,a}^0 - \widehat{P}_{s,a}^0\|_1 \leq 2$ and $\|[\widehat{V}^{\widehat{\pi},\sigma}]_\alpha - [\widehat{V}^{*,\sigma}]_\alpha\|_\infty \leq \|\widehat{V}^{\widehat{\pi},\sigma} - \widehat{V}^{*,\sigma}\|_\infty$, and (iii) follows from (48).

To control $\left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{*,\sigma}]_\alpha \right|$ in (138) for any given $\alpha \in [0, \frac{1}{1-\gamma}]$, and tame the dependency between $\widehat{V}^{*,\sigma}$ and \widehat{P}^0 , we resort to the following leave-one-out argument motivated by (Agarwal et al. 2020, Li et al. 2024a, Shi and Chi 2024). Specifically, we first construct a set of auxiliary RMDPs which simultaneously have the desired statistical independence between robust value functions and the estimated nominal transition kernel, and are minimally different from the original RMDPs under consideration. Then we control the term of interest associated with these auxiliary RMDPs and show the value is close to the target quantity for the desired RMDP. The process is divided into several steps as below.

Step 1: construction of auxiliary RMDPs with deterministic empirical nominal transitions. Recall that we target the empirical infinite-horizon robust MDP $\widehat{\mathcal{M}}_{\text{rob}}$ with the nominal transition kernel \widehat{P}^0 . Towards this, we can construct an auxiliary robust MDP $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ for each state s and any non-negative scalar $u \geq 0$, so that it is the same as $\widehat{\mathcal{M}}_{\text{rob}}$ except for the transition properties in state s . In particular, we define the nominal transition kernel and reward function of $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ as $P^{s,u}$ and $r^{s,u}$, which are expressed as follows

$$\begin{cases} P^{s,u}(s' | s, a) = \mathbb{1}(s' = s) & \text{for all } (s', a) \in \mathcal{S} \times \mathcal{A}, \\ P^{s,u}(\cdot | \tilde{s}, a) = \widehat{P}^0(\cdot | \tilde{s}, a) & \text{for all } (\tilde{s}, a) \in \mathcal{S} \times \mathcal{A} \text{ and } \tilde{s} \neq s, \end{cases} \quad (139)$$

and

$$\begin{cases} r^{s,u}(s, a) = u & \text{for all } a \in \mathcal{A}, \\ r^{s,u}(\tilde{s}, a) = r(\tilde{s}, a) & \text{for all } (\tilde{s}, a) \in \mathcal{S} \times \mathcal{A} \text{ and } \tilde{s} \neq s. \end{cases} \quad (140)$$

It is evident that the nominal transition probability at state s of the auxiliary $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$, i.e. it never leaves state s once entered. This useful property removes the randomness of $\widehat{P}_{s,a}^0$ for all $a \in \mathcal{A}$ in state s , which will be leveraged later.

Correspondingly, the robust Bellman operator $\widehat{\mathcal{T}}_{s,u}^\sigma(\cdot)$ associated with the RMDP $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ is defined as

$$\forall (\tilde{s}, a) \in \mathcal{S} \times \mathcal{A}: \quad \widehat{\mathcal{T}}_{s,u}^\sigma(Q)(\tilde{s}, a) = r^{s,u}(\tilde{s}, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{\tilde{s},a}^{s,u})} \mathcal{P}V, \quad \text{with } V(\tilde{s}) = \max_a Q(\tilde{s}, a). \quad (141)$$

Step 2: fixed-point equivalence between $\widehat{\mathcal{M}}_{\text{rob}}$ and the auxiliary RMDP $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$. Recall that $\widehat{Q}^{*,\sigma}$ is the unique fixed point of $\widehat{\mathcal{T}}^\sigma(\cdot)$ with the corresponding robust value $\widehat{V}^{*,\sigma}$. We assert that the corresponding robust value function $\widehat{V}_{s,u^*}^{*,\sigma}$ obtained from the fixed point of $\widehat{\mathcal{T}}_{s,u^*}^\sigma(\cdot)$ aligns with the robust value function $\widehat{V}^{*,\sigma}$ derived from $\widehat{\mathcal{T}}^\sigma(\cdot)$, as long as we choose u in the following manner:

$$u^* := u^*(s) = \widehat{V}^{*,\sigma}(s) - \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(e_s)} \mathcal{P}\widehat{V}^{*,\sigma}. \quad (142)$$

where e_s is the s -th standard basis vector in \mathbb{R}^S . Towards verifying this, we shall break our arguments in two different cases.

- **For state s :** One has for any $a \in \mathcal{A}$:

$$\begin{aligned} r^{s,u^*}(s,a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^{s,u^*})} \mathcal{P}\widehat{V}^{*,\sigma} &= u^* + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(e_s)} \mathcal{P}\widehat{V}^{*,\sigma} \\ &= \widehat{V}^{*,\sigma}(s) - \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(e_s)} \mathcal{P}\widehat{V}^{*,\sigma} + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(e_s)} \mathcal{P}\widehat{V}^{*,\sigma} = \widehat{V}^{*,\sigma}(s), \end{aligned} \quad (143)$$

where the first equality follows from the definition of $P_{s,a}^{s,u^*}$ in (139), and the second equality follows from plugging in the definition of u^* in (142).

- **For state $s' \neq s$:** It is easily verified that for all $a \in \mathcal{A}$,

$$\begin{aligned} r^{s,u^*}(s',a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s',a}^{s,u^*})} \mathcal{P}\widehat{V}^{*,\sigma} &= r(s',a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s',a}^0)} \mathcal{P}\widehat{V}^{*,\sigma} \\ &= \widehat{\mathcal{T}}^\sigma(\widehat{Q}^{*,\sigma})(s',a) = \widehat{Q}^{*,\sigma}(s',a), \end{aligned} \quad (144)$$

where the first equality follows from the definitions in (140) and (139), and the last line arises from the definition of the robust Bellman operator in (15), and that $\widehat{Q}^{*,\sigma}$ is the fixed point of $\widehat{\mathcal{T}}^\sigma(\cdot)$ (see Lemma 3).

Combining the facts in the above two cases, we establish that there exists a fixed point $\widehat{Q}_{s,u^*}^{*,\sigma}$ of the operator $\widehat{\mathcal{T}}_{s,u^*}^\sigma(\cdot)$ by taking

$$\begin{cases} \widehat{Q}_{s,u^*}^{*,\sigma}(s,a) = \widehat{V}^{*,\sigma}(s) & \text{for all } a \in \mathcal{A}, \\ \widehat{Q}_{s,u^*}^{*,\sigma}(s',a) = \widehat{Q}^{*,\sigma}(s',a) & \text{for all } s' \neq s \text{ and } a \in \mathcal{A}. \end{cases} \quad (145)$$

Consequently, we 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 $\widehat{V}^{*,\sigma}$. Note that the corresponding facts between $\widehat{\mathcal{M}}_{\text{rob}}$ and $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ in Step 1 and Step 2 hold in fact for any uncertainty set.

Step 3: building an ε -net for all reward values u . It is easily verified that

$$0 \leq u^* \leq \widehat{V}^{*,\sigma}(s) \leq \frac{1}{1-\gamma}. \quad (146)$$

We can construct a N_{ε_2} -net over the interval $[0, \frac{1}{1-\gamma}]$, where the size is bounded by $|N_{\varepsilon_2}| \leq \frac{3}{\varepsilon_2(1-\gamma)}$ (Ver-shynin 2018). Following the same arguments in the proof of Lemma 3, we can demonstrate that for each $u \in N_{\varepsilon_2}$, there exists a unique fixed point $\widehat{Q}_{s,u}^{*,\sigma}$ of the operator $\widehat{\mathcal{T}}_{s,u}^\sigma(\cdot)$, which satisfies $0 \leq \widehat{Q}_{s,u}^{*,\sigma} \leq \frac{1}{1-\gamma} \cdot 1$. Consequently, the corresponding robust value function also satisfies $\|\widehat{V}_{s,u}^{*,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$.

By the definitions in (139) and (140), we observe that for all $u \in N_{\varepsilon_2}$, $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ is statistically independent from $\widehat{P}_{s,a}^0$. This independence indicates that $[\widehat{V}_{s,u}^{*,\sigma}]_\alpha$ and $\widehat{P}_{s,a}^0$ are independent for a fixed α . With this in mind, invoking the fact in (127) and (128) and taking the union bound over all $(s,a,\alpha) \in \mathcal{S} \times \mathcal{A} \times N_{\varepsilon_1}$, $u \in N_{\varepsilon_2}$ yields that, with probability at least $1 - \delta$, it holds for all $(s,a,u) \in \mathcal{S} \times \mathcal{A} \times N_{\varepsilon_2}$ that

$$\begin{aligned} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}_{s,u}^{*,\sigma}]_\alpha \right| &\leq \varepsilon_2 + 2 \sqrt{\frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}_{s,u}^{*,\sigma})} + \frac{2 \log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{3N(1-\gamma)} \\ &\leq \varepsilon_2 + 3 \sqrt{\frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{(1-\gamma)^2 N}}, \end{aligned} \quad (147)$$

where the last inequality holds by the fact $\text{Var}_{P_{s,a}^0}(\widehat{V}_{s,u}^{*,\sigma}) \leq \|\widehat{V}_{s,u}^{*,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$ and letting $N \geq \log\left(\frac{18SAN|N_{\varepsilon_2}|}{\delta}\right)$.

Step 4: uniform concentration. Recalling that $u^* \in [0, \frac{1}{1-\gamma}]$ (see (146)), we can always find some $\bar{u} \in N_{\varepsilon_2}$ such that $|\bar{u} - u^*| \leq \varepsilon_2$. Consequently, plugging in the operator $\widehat{\mathcal{T}}_{s,u}^\sigma(\cdot)$ in (141) yields

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

With this in mind, we observe that the fixed points of $\widehat{\mathcal{T}}_{s,\bar{u}}^\sigma(\cdot)$ and $\widehat{\mathcal{T}}_{s,u^*}^\sigma(\cdot)$ obey

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

where the last inequality holds by the fact that $\widehat{\mathcal{T}}_{s,u}^\sigma(\cdot)$ is a γ -contraction. It directly indicates that

$$\left\| \widehat{Q}_{s,\bar{u}}^{*,\sigma} - \widehat{Q}_{s,u^*}^{*,\sigma} \right\|_\infty \leq \frac{\varepsilon_2}{(1-\gamma)} \quad \text{and} \quad \left\| \widehat{V}_{s,\bar{u}}^{*,\sigma} - \widehat{V}_{s,u^*}^{*,\sigma} \right\|_\infty \leq \left\| \widehat{Q}_{s,\bar{u}}^{*,\sigma} - \widehat{Q}_{s,u^*}^{*,\sigma} \right\|_\infty \leq \frac{\varepsilon_2}{(1-\gamma)}. \quad (148)$$

Armed with the above facts, to control the first term in (138), invoking the identity $\widehat{V}^{*,\sigma} = \widehat{V}_{s,u^*}^{*,\sigma}$ established in Step 2 gives that: for all $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned} &\max_{\alpha \in [\min_s \widehat{V}^{\pi,\sigma}(s), \max_s \widehat{V}^{\pi,\sigma}(s)]} \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}^{*,\sigma}]_\alpha \right| \\ &\leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}^{*,\sigma}]_\alpha \right| = \max_{\alpha \in [0, 1/(1-\gamma)]} \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}_{s,u^*}^{*,\sigma}]_\alpha \right| \\ &\stackrel{(i)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left\{ \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}_{s,\bar{u}}^{*,\sigma}]_\alpha \right| + \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) \left([\widehat{V}_{s,\bar{u}}^{*,\sigma}]_\alpha - [\widehat{V}_{s,u^*}^{*,\sigma}]_\alpha \right) \right| \right\} \\ &\stackrel{(ii)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}_{s,\bar{u}}^{*,\sigma}]_\alpha \right| + \frac{2\varepsilon_2}{(1-\gamma)} \\ &\stackrel{(iii)}{\leq} \frac{3\varepsilon_2}{(1-\gamma)} + \varepsilon_2 + 2\sqrt{\frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}_{s,\bar{u}}^{*,\sigma})} + \frac{2\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{3N(1-\gamma)} \\ &\leq \frac{3\varepsilon_2}{(1-\gamma)} + 2\sqrt{\frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma})} + \frac{2\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{3N(1-\gamma)} \\ &\quad + 2\sqrt{\frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\left| \text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma}) - \text{Var}_{P_{s,a}^0}(\widehat{V}_{s,\bar{u}}^{*,\sigma}) \right|} \\ &\stackrel{(iv)}{\leq} \frac{3\varepsilon_2}{(1-\gamma)} + 2\sqrt{\frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma})} + \frac{2\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{3N(1-\gamma)} + 2\sqrt{\frac{2\varepsilon_2 \log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{N(1-\gamma)^2}} \\ &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma})} + \frac{8\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} \end{aligned} \quad (149)$$

$$\leq 10\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2 N}}, \quad (150)$$

where (i) holds by the triangle inequality, (ii) arises from the bound

$$\begin{aligned} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) \left([\widehat{V}_{s,\bar{u}}^{*,\sigma}]_\alpha - [\widehat{V}_{s,u^*}^{*,\sigma}]_\alpha \right) \right| &\leq \left\| P_{s,a}^0 - \widehat{P}_{s,a}^0 \right\|_1 \left\| [\widehat{V}_{s,\bar{u}}^{*,\sigma}]_\alpha - [\widehat{V}_{s,u^*}^{*,\sigma}]_\alpha \right\|_\infty \\ &\leq 2 \left\| \widehat{V}_{s,\bar{u}}^{*,\sigma} - \widehat{V}_{s,u^*}^{*,\sigma} \right\|_\infty \leq \frac{2\varepsilon_2}{(1-\gamma)}, \end{aligned} \quad (151)$$

(iii) follows from (147), (iv) can be verified by applying Lemma 5 with (148). Here, the penultimate inequality holds by letting $\varepsilon_2 = \frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{N}$, which leads to $|N_{\varepsilon_2}| \leq \frac{3}{\varepsilon_2(1-\gamma)} \leq \frac{3N}{1-\gamma}$, and the last inequality holds by the fact $\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma}) \leq \|\widehat{V}^{*,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$ and letting $N \geq \log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)$.

Step 5: finishing up. Inserting (149) and (150) back into (138) and combining with (150) give that with probability at least $1 - \delta$,

$$\begin{aligned} \left| \widehat{P}_{s,a}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}_{s,a}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right| &\leq \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi},\sigma}(s), \max_s \widehat{V}^{\widehat{\pi},\sigma}(s)]} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{*,\sigma}]_\alpha \right| + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \\ &\leq \max_{\alpha \in [0,1/(1-\gamma)]} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{*,\sigma}]_\alpha \right| + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \\ &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma})} + \frac{8\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \\ &\leq 10\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \end{aligned} \quad (152)$$

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

Finally, we complete the proof by compiling everything into the matrix form as follows:

$$\begin{aligned} \left| \widehat{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right| &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma})} \mathbf{1} + \frac{8\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} \mathbf{1} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \mathbf{1} \\ &\leq 10\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2N}} \mathbf{1} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \mathbf{1}. \end{aligned} \quad (153)$$

B.3.6. Proof of Lemma 12 The proof can be achieved by directly applying the same routine as Appendix B.3.4. Towards this, similar to (131), we arrive at

$$\left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})} \leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left(\underline{P}^{\widehat{\pi},\widehat{V}} \right)^t \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})}. \quad (154)$$

To control $\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})$, we denote the minimum value of $\widehat{V}^{\widehat{\pi},\sigma}$ as $V_{\min} = \min_{s \in \mathcal{S}} \widehat{V}^{\widehat{\pi},\sigma}(s)$ and $V' := \widehat{V}^{\widehat{\pi},\sigma} - V_{\min} \mathbf{1}$. By the same argument as (133), we arrive at

$$\begin{aligned} &\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma}) \\ &\leq \underline{P}^{\widehat{\pi},\widehat{V}} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_\infty \mathbf{1} + \frac{2}{\gamma} \|V'\|_\infty \left| \left(\widehat{P}^{\widehat{\pi},\widehat{V}} - \underline{P}^{\widehat{\pi},\widehat{V}} \right) \widehat{V}^{\widehat{\pi},\sigma} \right| \\ &\leq \underline{P}^{\widehat{\pi},\widehat{V}} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_\infty \mathbf{1} + \frac{2}{\gamma} \|V'\|_\infty \left(10\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \right) \mathbf{1}, \end{aligned} \quad (155)$$

where the last inequality makes use of Lemma [11](#). Plugging [\(155\)](#) back into [\(154\)](#) leads to

$$\begin{aligned}
\left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}}\right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} &\stackrel{(i)}{\leq} \sqrt{\frac{1}{1-\gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t \left(\underline{P}^{\hat{\pi}, \hat{V}}\right)^t \left(\underline{P}^{\hat{\pi}, \hat{V}}(V' \circ V') - \frac{1}{\gamma} V' \circ V'\right) \right|} \\
&+ \sqrt{\frac{1}{(1-\gamma)^2 \gamma^2} \left(2 + 20 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2 N}} + \frac{2\gamma \varepsilon_{\text{opt}}}{1-\gamma}\right) \|V'\|_{\infty}} \\
&\stackrel{(ii)}{\leq} \sqrt{\frac{\|V'\|_{\infty}^2}{\gamma(1-\gamma)}} 1 + \sqrt{\frac{\left(2 + 20 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2 N}} + \frac{2\gamma \varepsilon_{\text{opt}}}{1-\gamma}\right) \|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1 \\
&\stackrel{(iii)}{\leq} \sqrt{\frac{\|V'\|_{\infty}^2}{\gamma(1-\gamma)}} 1 + \sqrt{\frac{24\|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1 \leq 6 \sqrt{\frac{\|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1, \tag{156}
\end{aligned}$$

where (i) arises from following the routine of [\(135\)](#), (ii) holds by repeating the argument of [\(136\)](#), (iii) follows by taking $N \geq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2}$ and $\varepsilon_{\text{opt}} \leq \frac{1-\gamma}{\gamma}$, and the last inequality holds by $\|V'\|_{\infty} \leq \|V^{*,\sigma}\|_{\infty} \leq \frac{1}{1-\gamma}$.

Finally, applying Lemma [6](#) with $P = \hat{P}^0$ and $\pi = \hat{\pi}$ yields

$$\|V'\|_{\infty} \leq \max_{s \in \mathcal{S}} \hat{V}^{\hat{\pi}, \sigma}(s) - \min_{s \in \mathcal{S}} \hat{V}^{\hat{\pi}, \sigma}(s) \leq \frac{1}{\gamma \max\{1-\gamma, \sigma\}},$$

which can be inserted into [\(156\)](#) and gives

$$\left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}}\right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} \leq 6 \sqrt{\frac{1}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, \sigma\}}} 1 \leq 6 \sqrt{\frac{1}{(1-\gamma)^3 \gamma^2}} 1.$$

C. Proof of the lower bound with TV distance: Theorem [2](#)

C.1. Construction of the hard problem instances

First, note that we shall use the same MDPs defined in Section [5.3](#) as follows

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

In particular, we shall keep the structure of the transition kernel in [\(56\)](#), reward function in [\(61\)](#) and initial state distribution in [\(62\)](#), while p and Δ shall be specified according to TV distance case.

Uncertainty set of the transition kernels. Recalling the uncertainty set assumed throughout this section is defined as $\mathcal{U}^{\sigma}(P^{\phi})$ with TV distance:

$$\mathcal{U}^{\sigma}(P^{\phi}) := \mathcal{U}_{\text{TV}}^{\sigma}(P^{\phi}) = \otimes \mathcal{U}_{\text{TV}}^{\sigma}(P_{s,a}^{\phi}), \quad \mathcal{U}_{\text{TV}}^{\sigma}(P_{s,a}^{\phi}) := \left\{ P'_{s,a} \in \Delta(\mathcal{S}) : \frac{1}{2} \|P'_{s,a} - P_{s,a}^{\phi}\|_1 \leq \sigma \right\}, \tag{157}$$

where $P_{s,a}^{\phi} := P^{\phi}(\cdot \mid s, a)$ is defined similar to [\(4\)](#). In addition, without loss of generality, we recall the radius $\sigma \in (0, 1 - c_0]$ with $0 < c_0 < 1$. With the uncertainty level in hand, taking $c_1 := \frac{c_0}{2}$, p and Δ which determines the instances obey (recall [\(58\)](#))

$$p = (1 + c_1) \max\{1 - \gamma, \sigma\} \quad \text{and} \quad \Delta \leq c_1 \max\{1 - \gamma, \sigma\}, \tag{158}$$

which ensure $0 \leq p \leq 1$ as follows:

$$(1 + c_1)\sigma \leq 1 - c_0 + c_1\sigma \leq 1 - \frac{c_0}{2} < 1, \quad (1 + c_1)(1 - \gamma) \leq \frac{3}{2}(1 - \gamma) \leq \frac{3}{4} < 1. \quad (159)$$

Consequently, applying (57) directly leads to

$$p \geq q \geq \max\{1 - \gamma, \sigma\}. \quad (160)$$

To continue, for any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, we denote the infimum probability of moving to the next state s' associated with any perturbed transition kernel $P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^\phi)$ as

$$\underline{P}^\phi(s' | s, a) := \inf_{P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^\phi)} P(s' | s, a) = \max\{P(s' | s, a) - \sigma, 0\}, \quad (161)$$

where the last equation can be easily verified by the definition of $\mathcal{U}^\sigma(P^\phi)$ in (157). As shall be seen, the transition from state 0 to state 1 plays an important role in the analysis, for convenience, we denote

$$\underline{p} := \underline{P}^\phi(1 | 0, \phi) = p - \sigma, \quad \underline{q} := \underline{P}^\phi(1 | 0, 1 - \phi) = q - \sigma, \quad (162)$$

which follows from the fact that $p \geq q \geq \sigma$ in (160).

Robust value functions and robust optimal policies. To proceed, we are ready to derive the corresponding robust value functions, identify the optimal policies, and characterize the optimal values. For any MDP \mathcal{M}_ϕ with the above uncertainty set, we denote π_ϕ^* as the optimal policy, and the robust value function of any policy π (resp. the optimal policy π_ϕ^*) as $V_\phi^{\pi, \sigma}$ (resp. $V_\phi^{*, \sigma}$). Then, we introduce the following lemma which describes some important properties of the robust (optimal) value functions and optimal policies. The proof is postponed to Appendix C.3.1.

LEMMA 13. *For any $\phi \in \{0, 1\}$ and any policy π , the robust value function obeys*

$$V_\phi^{\pi, \sigma}(0) = \frac{\gamma(z_\phi^\pi - \sigma)}{(1 - \gamma) \left(1 + \frac{\gamma(z_\phi^\pi - \sigma)}{1 - \gamma(1 - \sigma)}\right) (1 - \gamma(1 - \sigma))}, \quad (163)$$

where z_ϕ^π is defined as

$$z_\phi^\pi := p\pi(\phi | 0) + q\pi(1 - \phi | 0). \quad (164)$$

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

$$V_\phi^{*, \sigma}(0) = \frac{\gamma(p - \sigma)}{(1 - \gamma) \left(1 + \frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)}\right) (1 - \gamma(1 - \sigma))}, \quad (165a)$$

$$\pi_\phi^*(\phi | s) = 1, \quad \text{for } s \in \mathcal{S}. \quad (165b)$$

C.2. Establishing the minimax lower bound

Note that our goal is to control the quantity w.r.t. any policy estimator $\hat{\pi}$ based on the chosen initial distribution φ in (62) and the dataset consisting of N samples over each state-action pair generated from the nominal transition kernel P^ϕ , which gives

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

Step 1: converting the goal to estimate ϕ . We make the following useful claim which shall be verified in Appendix C.3.2: With $\varepsilon \leq \frac{c_1}{32(1-\gamma)}$, letting

$$\Delta = 32(1-\gamma) \max\{1-\gamma, \sigma\} \varepsilon \leq c_1 \max\{1-\gamma, \sigma\} \quad (166)$$

which satisfies (158), it leads to that for any policy $\hat{\pi}$,

$$\langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\hat{\pi},\sigma} \rangle \geq 2\varepsilon(1 - \hat{\pi}(\phi|0)). \quad (167)$$

With this connection established between the policy $\hat{\pi}$ and its sub-optimality gap as depicted in (167), we can now proceed to build an estimate for ϕ . Here, we denote \mathbb{P}_ϕ as the probability distribution when the MDP is \mathcal{M}_ϕ , where ϕ can take on values in the set $\{0, 1\}$.

Let's assume momentarily that an estimated policy $\hat{\pi}$ achieves

$$\mathbb{P}_\phi \{ \langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\hat{\pi},\sigma} \rangle \leq \varepsilon \} \geq \frac{7}{8}, \quad (168)$$

then in view of (167), we necessarily have $\hat{\pi}(\phi|0) \geq \frac{1}{2}$ with probability at least $\frac{7}{8}$. With this in mind, we are motivated to construct the following estimate $\hat{\phi}$ for $\phi \in \{0, 1\}$:

$$\hat{\phi} = \arg \max_{a \in \{0,1\}} \hat{\pi}(a|0), \quad (169)$$

which obeys

$$\mathbb{P}_\phi \{ \hat{\phi} = \phi \} \geq \mathbb{P}_\phi \{ \hat{\pi}(\phi|0) > 1/2 \} \geq \frac{7}{8}. \quad (170)$$

Subsequently, our aim is to demonstrate that (170) cannot occur without an adequate number of samples, which would in turn contradict (167).

Step 2: probability of error in testing two hypotheses. Equipped with the aforementioned groundwork, we can now delve into differentiating between the two hypotheses $\phi \in \{0, 1\}$. To achieve this, we consider the concept of minimax probability of error, defined as follows:

$$p_e := \inf_{\psi} \max \{ \mathbb{P}_0(\psi \neq 0), \mathbb{P}_1(\psi \neq 1) \}. \quad (171)$$

Here, the infimum is taken over all possible tests ψ constructed from the samples generated from the nominal transition kernel P^ϕ .

Moving forward, let us denote μ_ϕ (resp. $\mu_\phi(s)$) as the distribution of a sample tuple (s_i, a_i, s'_i) under the nominal transition kernel P^ϕ associated with \mathcal{M}_ϕ and the samples are generated independently. Applying standard results from [Tsybakov \(2009\)](#), Theorem 2.2) and the additivity of the KL divergence (cf. [Tsybakov \(2009\)](#), Page 85)), we obtain

$$\begin{aligned} p_e &\geq \frac{1}{4} \exp\left(-NSA \cdot \text{KL}(\mu_0 \parallel \mu_1)\right) \\ &= \frac{1}{4} \exp\left\{-N\left(\text{KL}(P^0(\cdot|0,0) \parallel P^1(\cdot|0,0)) + \text{KL}(P^0(\cdot|0,1) \parallel P^1(\cdot|0,1))\right)\right\}, \end{aligned} \quad (172)$$

where the last inequality holds by observing that

$$\begin{aligned} \text{KL}(\mu_0 \parallel \mu_1) &= \frac{1}{SA} \sum_{s,a,s'} \text{KL}(P^0(s'|s,a) \parallel P^1(s'|s,a)) \\ &= \frac{1}{SA} \sum_{a \in \{0,1\}} \text{KL}(P^0(\cdot|0,a) \parallel P^1(\cdot|0,a)), \end{aligned}$$

Here, the last equality holds by the fact that $P^0(\cdot|s,a)$ and $P^1(\cdot|s,a)$ only differ when $s=0$.

Now, our focus shifts towards bounding the terms involving the KL divergence in [\(172\)](#). Given $p \geq q \geq \max\{1-\gamma, \sigma\}$ (cf. [\(160\)](#)), applying [Tsybakov \(2009\)](#), Lemma 2.7) gives

$$\begin{aligned} \text{KL}(P^0(\cdot|0,1) \parallel P^1(\cdot|0,1)) &= \text{KL}(p \parallel q) \leq \frac{(p-q)^2}{(1-p)p} \stackrel{(i)}{=} \frac{\Delta^2}{p(1-p)} \\ &\stackrel{(ii)}{=} \frac{1024(1-\gamma)^2 \max\{1-\gamma, \sigma\}^2 \varepsilon^2}{p(1-p)} \\ &\leq \frac{1024(1-\gamma)^2 \max\{1-\gamma, \sigma\}^2 \varepsilon^2}{1-p} \leq \frac{4096}{c_1} (1-\gamma)^2 \max\{1-\gamma, \sigma\}^2 \varepsilon^2, \end{aligned} \quad (173)$$

where (i) stems from the definition in [\(57\)](#), (ii) follows by the expression of Δ in [\(166\)](#), and the last inequality arises from $1-q \geq 1-p \geq \frac{c_0}{4}$ (see [\(159\)](#)).

Note that it can be shown that $\text{KL}(P^0(\cdot|0,0) \parallel P^1(\cdot|0,0))$ can be upper bounded in the same manner. Substituting [\(173\)](#) back into [\(172\)](#) demonstrates that: if the sample size is selected as

$$N \leq \frac{c_1 \log 2}{8192(1-\gamma)^2 \max\{1-\gamma, \sigma\}^2 \varepsilon^2}, \quad (174)$$

then one necessarily has

$$p_e \geq \frac{1}{4} \exp\left\{-N \frac{8192}{c_1} (1-\gamma)^2 \max\{1-\gamma, \sigma\}^2 \varepsilon^2\right\} \geq \frac{1}{8}, \quad (175)$$

Step 3: putting the results together. Lastly, suppose that there exists an estimator $\widehat{\pi}$ such that

$$\mathbb{P}_0\{\langle \varphi, V_0^{*,\sigma} - V_0^{\widehat{\pi},\sigma} \rangle > \varepsilon\} < \frac{1}{8} \quad \text{and} \quad \mathbb{P}_1\{\langle \varphi, V_1^{*,\sigma} - V_1^{\widehat{\pi},\sigma} \rangle > \varepsilon\} < \frac{1}{8}.$$

According to Step 1, the estimator $\widehat{\phi}$ defined in (169) must satisfy

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

However, this cannot occur under the sample size condition (174) to avoid contradiction with (175). Thus, we have completed the proof.

C.3. Proof of the auxiliary facts

C.3.1. Proof of Lemma 13

Deriving the robust value function over different states. For any \mathcal{M}_ϕ with $\phi \in \{0, 1\}$, we first characterize the robust value function of any policy π over different states. Before proceeding, we denote the minimum of the robust value function over states as below:

$$V_{\phi,\min}^{\pi,\sigma} := \min_{s \in \mathcal{S}} V_\phi^{\pi,\sigma}(s). \quad (176)$$

Clearly, there exists at least one state $s_{\phi,\min}^\pi$ that satisfies $V_\phi^{\pi,\sigma}(s_{\phi,\min}^\pi) = V_{\phi,\min}^{\pi,\sigma}$.

With this in mind, it is easily observed that for any policy π , the robust value function at state $s = 1$ obeys

$$\begin{aligned} V_\phi^{\pi,\sigma}(1) &= \mathbb{E}_{a \sim \pi(\cdot|1)} \left[r(1, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{1,a}^\phi)} \mathcal{P} V_\phi^{\pi,\sigma} \right] \\ &\stackrel{(i)}{=} 1 + \gamma \mathbb{E}_{a \sim \pi(\cdot|1)} \left[\underline{P}^\phi(1|1, a) V_\phi^{\pi,\sigma}(1) \right] + \gamma \sigma V_{\phi,\min}^{\pi,\sigma} \stackrel{(ii)}{=} 1 + \gamma(1 - \sigma) V_\phi^{\pi,\sigma}(1) + \gamma \sigma V_{\phi,\min}^{\pi,\sigma}, \end{aligned} \quad (177)$$

where (i) holds by $r(1, a) = 1$ for all $a \in \mathcal{A}'$ and (161), and (ii) follows from $P^\phi(1|1, a) = 1$ for all $a \in \mathcal{A}'$.

Similarly, for any $s \in \{2, 3, \dots, S-1\}$, we have

$$\begin{aligned} V_\phi^{\pi,\sigma}(s) &= 0 + \gamma \mathbb{E}_{a \sim \pi(\cdot|s)} \left[\underline{P}^\phi(1|s, a) V_\phi^{\pi,\sigma}(1) \right] + \gamma \sigma V_{\phi,\min}^{\pi,\sigma} \\ &= \gamma(1 - \sigma) V_\phi^{\pi,\sigma}(1) + \gamma \sigma V_{\phi,\min}^{\pi,\sigma}, \end{aligned} \quad (178)$$

since $r(s, a) = 0$ for all $s \in \{2, 3, \dots, S-1\}$ and the definition in (161).

Finally, we move onto compute $V_\phi^{\pi,\sigma}(0)$, the robust value function at state 0 associated with any policy π . First, it obeys

$$\begin{aligned} V_\phi^{\pi,\sigma}(0) &= \mathbb{E}_{a \sim \pi(\cdot|0)} \left[r(0, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,a}^\phi)} \mathcal{P} V_\phi^{\pi,\sigma} \right] \\ &= 0 + \gamma \pi(\phi|0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,\phi}^\phi)} \mathcal{P} V_\phi^{\pi,\sigma} + \gamma \pi(1-\phi|0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,1-\phi}^\phi)} \mathcal{P} V_\phi^{\pi,\sigma}. \end{aligned} \quad (179)$$

Recall the transition kernel defined in (56) and the fact about the uncertainty set over state 0 in (162), it is easily verified that the following probability vector $P_1 \in \Delta(\mathcal{S})$ obeys $P_1 \in \mathcal{U}^\sigma(P_{0,\phi}^\phi)$, which is defined as

$$\begin{aligned} P_1(0) &= 1 - p + \sigma \mathbb{1}(0 = s_{\phi,\min}^\pi), & P_1(1) &= \underline{p} = p - \sigma, \\ P_1(s) &= \sigma \mathbb{1}(s = s_{\phi,\min}^\pi), & \forall s \in \{2, 3, \dots, S-1\}, \end{aligned} \quad (180)$$

where $\underline{p} = p - \sigma$ due to (162). Similarly, the following probability vector $P_2 \in \Delta(\mathcal{S})$ also falls into the uncertainty set $\mathcal{U}^\sigma(P_{0,1-\phi}^\phi)$:

$$\begin{aligned} P_2(0) &= 1 - q + \sigma \mathbb{1}(0 = s_{\phi,\min}^\pi), & P_2(1) &= \underline{q} = q - \sigma, \\ P_2(s) &= \sigma \mathbb{1}(s = s_{\phi,\min}^\pi) & \forall s \in \{2, 3, \dots, S-1\}. \end{aligned} \quad (181)$$

It is noticed that P_1 and P_2 defined above are the worst-case perturbations, since the probability mass at state 1 will be moved to the state with the least value. Plugging the above facts about $P_1 \in \mathcal{U}^\sigma(P_{0,\phi}^\phi)$ and $P_2 \in \mathcal{U}^\sigma(P_{0,1-\phi}^\phi)$ into (179), we arrive at

$$\begin{aligned} V_\phi^{\pi,\sigma}(0) &\leq \gamma\pi(\phi|0)P_1V_\phi^{\pi,\sigma} + \gamma\pi(1-\phi|0)P_2V_\phi^{\pi,\sigma} \\ &= \gamma\pi(\phi|0) \left[(p-\sigma)V_\phi^{\pi,\sigma}(1) + (1-p)V_\phi^{\pi,\sigma}(0) + \sigma V_{\phi,\min}^{\pi,\sigma} \right] \\ &\quad + \gamma\pi(1-\phi|0) \left[(q-\sigma)V_\phi^{\pi,\sigma}(1) + (1-q)V_\phi^{\pi,\sigma}(0) + \sigma V_{\phi,\min}^{\pi,\sigma} \right] \\ &\stackrel{(i)}{=} \gamma(z_\phi^\pi - \sigma)V_\phi^{\pi,\sigma}(1) + \gamma\sigma V_{\phi,\min}^{\pi,\sigma} + \gamma(1-z_\phi^\pi)V_\phi^{\pi,\sigma}(0), \end{aligned} \quad (182)$$

where the last equality holds by the definition of z_ϕ^π in (164). To continue, recursively applying (182) yields

$$\begin{aligned} &V_\phi^{\pi,\sigma}(0) \\ &\leq \gamma(z_\phi^\pi - \sigma)V_\phi^{\pi,\sigma}(1) + \gamma\sigma V_{\phi,\min}^{\pi,\sigma} + \gamma(1-z_\phi^\pi) \left[\gamma(z_\phi^\pi - \sigma)V_\phi^{\pi,\sigma}(1) + \gamma\sigma V_{\phi,\min}^{\pi,\sigma} + \gamma(1-z_\phi^\pi)V_\phi^{\pi,\sigma}(0) \right] \\ &\stackrel{(i)}{\leq} \gamma(z_\phi^\pi - \sigma)V_\phi^{\pi,\sigma}(1) + \gamma\sigma V_{\phi,\min}^{\pi,\sigma} + \gamma(1-z_\phi^\pi) \left[\gamma z_\phi^\pi V_\phi^{\pi,\sigma}(1) + \gamma(1-z_\phi^\pi)V_\phi^{\pi,\sigma}(0) \right] \\ &\leq \dots \\ &\leq \gamma(z_\phi^\pi - \sigma)V_\phi^{\pi,\sigma}(1) + \gamma\sigma V_{\phi,\min}^{\pi,\sigma} + \gamma z_\phi^\pi \sum_{t=1}^{\infty} \gamma^t (1-z_\phi^\pi)^t V_\phi^{\pi,\sigma}(1) + \lim_{t \rightarrow \infty} \gamma^t (1-z_\phi^\pi)^t V_\phi^{\pi,\sigma}(0) \\ &\stackrel{(ii)}{\leq} \gamma(z_\phi^\pi - \sigma)V_\phi^{\pi,\sigma}(1) + \gamma\sigma V_{\phi,\min}^{\pi,\sigma} + \gamma(1-z_\phi^\pi) \frac{\gamma z_\phi^\pi}{1-\gamma(1-z_\phi^\pi)} V_\phi^{\pi,\sigma}(1) + 0 \\ &< \gamma(z_\phi^\pi - \sigma)V_\phi^{\pi,\sigma}(1) + \gamma\sigma V_{\phi,\min}^{\pi,\sigma} + \gamma(1-z_\phi^\pi)V_\phi^{\pi,\sigma}(1) \\ &= \gamma(1-\sigma)V_\phi^{\pi,\sigma}(1) + \gamma\sigma V_{\phi,\min}^{\pi,\sigma}, \end{aligned} \quad (183)$$

where (i) uses $V_{\phi,\min}^{\pi,\sigma} \leq V_\phi^{\pi,\sigma}(1)$, (ii) follows from $\gamma(1-z_\phi^\pi) < 1$, and the penultimate line follows from the trivial fact that $\frac{\gamma z_\phi^\pi}{1-\gamma(1-z_\phi^\pi)} < 1$.

Combining (177), (178), and (183), we have that for any policy π ,

$$V_{\phi}^{\pi, \sigma}(0) = V_{\phi, \min}^{\pi, \sigma}, \quad (184)$$

which directly leads to

$$V_{\phi}^{\pi, \sigma}(1) = 1 + \gamma(1 - \sigma)V_{\phi}^{\pi, \sigma}(1) + \gamma\sigma V_{\phi, \min}^{\pi, \sigma} = \frac{1 + \gamma\sigma V_{\phi}^{\pi, \sigma}(0)}{1 - \gamma(1 - \sigma)}. \quad (185)$$

Let's now return to the characterization of $V_{\phi}^{\pi, \sigma}(0)$. In view of (184), the equality in (182) holds, and we have

$$\begin{aligned} V_{\phi}^{\pi, \sigma}(0) &= \gamma(z_{\phi}^{\pi} - \sigma)V_{\phi}^{\pi, \sigma}(1) + \gamma(1 - z_{\phi}^{\pi} + \sigma)V_{\phi}^{\pi, \sigma}(0) \\ &\stackrel{(i)}{=} \gamma(z_{\phi}^{\pi} - \sigma) \frac{1 + \gamma\sigma V_{\phi}^{\pi, \sigma}(0)}{1 - \gamma(1 - \sigma)} + \gamma(1 - z_{\phi}^{\pi} + \sigma)V_{\phi}^{\pi, \sigma}(0) \\ &= \frac{\gamma(z_{\phi}^{\pi} - \sigma)}{1 - \gamma(1 - \sigma)} + \gamma \left(1 + (z_{\phi}^{\pi} - \sigma) \frac{\gamma\sigma - (1 - \gamma(1 - \sigma))}{1 - \gamma(1 - \sigma)} \right) V_{\phi}^{\pi, \sigma}(0) \\ &= \frac{\gamma(z_{\phi}^{\pi} - \sigma)}{1 - \gamma(1 - \sigma)} + \gamma \left(1 - \frac{(1 - \gamma)(z_{\phi}^{\pi} - \sigma)}{1 - \gamma(1 - \sigma)} \right) V_{\phi}^{\pi, \sigma}(0), \end{aligned}$$

where (i) arises from (185). Solving this relation gives

$$V_{\phi}^{\pi, \sigma}(0) = \frac{\frac{\gamma(z_{\phi}^{\pi} - \sigma)}{1 - \gamma(1 - \sigma)}}{(1 - \gamma) \left(1 + \frac{\gamma(z_{\phi}^{\pi} - \sigma)}{1 - \gamma(1 - \sigma)} \right)}. \quad (186)$$

The optimal robust policy and optimal robust value function. We move on to characterize the robust optimal policy and its corresponding robust value function. To begin with, denoting

$$z := \frac{\gamma(z_{\phi}^{\pi} - \sigma)}{1 - \gamma(1 - \sigma)}, \quad (187)$$

we rewrite (186) as

$$V_{\phi}^{\pi, \sigma}(0) = \frac{z}{(1 - \gamma)(1 + z)} =: f(z).$$

Plugging in the fact that $z_{\phi}^{\pi} \geq q \geq \sigma > 0$ in (160), it follows that $z > 0$. So for any $z > 0$, the derivative of $f(z)$ w.r.t. z obeys

$$\frac{(1 - \gamma)(1 + z) - (1 - \gamma)z}{(1 - \gamma)^2(1 + z)^2} = \frac{1}{(1 - \gamma)(1 + z)^2} > 0. \quad (188)$$

Observing that $f(z)$ is increasing in z , z is increasing in z_{ϕ}^{π} , and z_{ϕ}^{π} is also increasing in $\pi(\phi | 0)$ (see the fact $p \geq q$ in (160)), the optimal policy in state 0 thus obeys

$$\pi_{\phi}^*(\phi | 0) = 1. \quad (189)$$

Considering that the action does not influence the state transition for all states $s > 0$, without loss of generality, we choose the robust optimal policy to obey

$$\forall s > 0: \quad \pi_\phi^*(\phi | s) = 1. \quad (190)$$

Taking $\pi = \pi_\phi^*$, we complete the proof by showing that the corresponding optimal robust value function at state 0 as follows:

$$V_\phi^{*,\sigma}(0) = \frac{\frac{\gamma(z_\phi^{\pi^*} - \sigma)}{1-\gamma(1-\sigma)}}{(1-\gamma)\left(1 + \frac{\gamma(z_\phi^{\pi^*} - \sigma)}{1-\gamma(1-\sigma)}\right)} = \frac{\frac{\gamma(p-\sigma)}{1-\gamma(1-\sigma)}}{(1-\gamma)\left(1 + \frac{\gamma(p-\sigma)}{1-\gamma(1-\sigma)}\right)}. \quad (191)$$

C.3.2. Proof of the claim (167) Plugging in the definition of φ , we arrive at that for any policy π ,

$$\langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\pi,\sigma} \rangle = V_\phi^{*,\sigma}(0) - V_\phi^{\pi,\sigma}(0) = \frac{\frac{\gamma(p-z_\phi^\pi)}{1-\gamma(1-\sigma)}}{(1-\gamma)\left(1 + \frac{\gamma(p-\sigma)}{1-\gamma(1-\sigma)}\right)\left(1 + \frac{\gamma(z_\phi^\pi - \sigma)}{1-\gamma(1-\sigma)}\right)}, \quad (192)$$

which follows from applying (163) and basic calculus. Then, we proceed to control the above term in two cases separately in terms of the uncertainty level σ .

- When $\sigma \in (0, 1 - \gamma]$. Then regarding the important terms in (192), we observe that

$$1 - \gamma < 1 - \gamma(1 - \sigma) \leq 1 - \gamma(1 - (1 - \gamma)) = (1 - \gamma)(1 + \gamma) \leq 2(1 - \gamma), \quad (193)$$

which directly leads to

$$\frac{\gamma(z_\phi^\pi - \sigma)}{1 - \gamma(1 - \sigma)} \stackrel{(i)}{\leq} \frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)} \leq \frac{\gamma c_1(1 - \gamma)}{1 - \gamma(1 - \sigma)} \stackrel{(ii)}{<} c_1 \gamma, \quad (194)$$

where (i) holds by $z_\phi^\pi < p$, and (ii) is due to (193). Inserting (193) and (194) back into (192), we arrive at

$$\begin{aligned} \langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\pi,\sigma} \rangle &\geq \frac{\frac{\gamma(p-z_\phi^\pi)}{2(1-\gamma)}}{(1-\gamma)(1+c_1\gamma)^2} \geq \frac{\gamma(p-z_\phi^\pi)}{8(1-\gamma)^2} \\ &= \frac{\gamma(p-q)(1-\pi(\phi|0))}{8(1-\gamma)^2} = \frac{\gamma\Delta(1-\pi(\phi|0))}{8(1-\gamma)^2} \geq 2\varepsilon(1-\pi(\phi|0)), \end{aligned} \quad (195)$$

where the last inequality holds by setting $(\gamma \geq 1/2)$

$$\Delta = 32(1-\gamma)^2\varepsilon. \quad (196)$$

Finally, it is easily verified that

$$\varepsilon \leq \frac{c_1}{32(1-\gamma)} \implies \Delta \leq c_1(1-\gamma).$$

- When $\sigma \in (1 - \gamma, 1 - c_1]$. Regarding (192), we observe that

$$\gamma\sigma < 1 - \gamma(1 - \sigma) = 1 - \gamma + \gamma\sigma \leq (1 + \gamma)\sigma \leq 2\sigma, \quad (197)$$

which directly leads to

$$\frac{\gamma(z_\phi^\pi - \sigma)}{1 - \gamma(1 - \sigma)} \leq \frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)} \leq \frac{\gamma c_1 \sigma}{1 - \gamma(1 - \sigma)} \stackrel{(i)}{<} c_1, \quad (198)$$

where (i) holds by (197). Inserting (197) and (198) back into (192), we arrive at

$$\begin{aligned} \langle \varphi, V_\phi^{*\sigma} - V_\phi^{\pi,\sigma} \rangle &\geq \frac{\frac{\gamma(p - z_\phi^\pi)}{2\sigma}}{(1 - \gamma)(1 + c_1)^2} \geq \frac{\gamma(p - z_\phi^\pi)}{8(1 - \gamma)\sigma} = \frac{\gamma(p - q)(1 - \pi(\phi|0))}{8(1 - \gamma)\sigma} \\ &= \frac{\gamma\Delta(1 - \pi(\phi|0))}{8(1 - \gamma)\sigma} \geq 2\varepsilon(1 - \pi(\phi|0)), \end{aligned} \quad (199)$$

where the last inequality holds by letting $(\gamma \geq 1/2)$

$$\Delta = 32(1 - \gamma)\sigma\varepsilon. \quad (200)$$

Finally, it is easily verified that

$$\varepsilon \leq \frac{c_1}{32(1 - \gamma)} \implies \Delta \leq c_1\sigma. \quad (201)$$

C.3.3. Proof of Lemma 15 The proof follows the same routine as that of Lemma 9. Taking the same pipeline as that in Appendix B.3.4, similar to (133), we have

$$\begin{aligned} \text{Var}_{\underline{\hat{P}}^{\pi^*,V}}(V^{*\sigma}) &\leq \underline{\hat{P}}^{\pi^*,V}(V' \circ V') - \frac{1}{\gamma}V' \circ V' + \frac{2}{\gamma^2}\|V'\|_\infty 1 + \frac{2}{\gamma}\|V'\|_\infty \left| \left(\underline{\hat{P}}^{\pi^*,V} - \widehat{\underline{P}}^{\pi^*,V} \right) V^{*\sigma} \right| \\ &\leq \underline{\hat{P}}^{\pi^*,V}(V' \circ V') - \frac{1}{\gamma}V' \circ V' + \frac{2}{\gamma^2}\|V'\|_\infty 1 + \frac{4}{\gamma}\|V'\|_\infty (1 + 2\sqrt{\sigma}) \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1 - \gamma)^2 N}} 1, \end{aligned} \quad (202)$$

where the last inequality holds by Lemma 14.

Plugging the above results leads to

$$\begin{aligned} &\left(I - \gamma \widehat{\underline{P}}^{\pi^*,V} \right)^{-1} \sqrt{\text{Var}_{\underline{\hat{P}}^{\pi^*,V}}(V^{*\sigma})} \\ &\leq \sqrt{\frac{1}{1 - \gamma}} \left(\sum_{t=0}^{\infty} \gamma^t \left(\widehat{\underline{P}}^{\pi^*,V} \right)^t \left(\underline{\hat{P}}^{\pi^*,V}(V' \circ V') - \frac{1}{\gamma}V' \circ V' + \frac{2}{\gamma^2}\|V'\|_\infty 1 \right. \right. \\ &\quad \left. \left. + \frac{4}{\gamma}\|V'\|_\infty (1 + 2\sqrt{\sigma}) \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1 - \gamma)^2 N}} 1 \right) \right)^{1/2} \\ &\stackrel{(i)}{\leq} \sqrt{\frac{1}{1 - \gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t \left(\widehat{\underline{P}}^{\pi^*,V} \right)^t \left(\underline{\hat{P}}^{\pi^*,V}(V' \circ V') - \frac{1}{\gamma}V' \circ V' \right) \right|} \end{aligned}$$

$$\begin{aligned}
 & + \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left(\widehat{P}^{\pi^*, V} \right)^t \left(\frac{2}{\gamma^2} \|V'\|_{\infty} 1 + \frac{4}{\gamma} \|V'\|_{\infty} (1 + 2\sqrt{\sigma}) \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} 1 \right)} \\
 & \leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left(\widehat{P}^{\pi^*, V} \right)^t \left[\widehat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' \right]} \\
 & \quad + \sqrt{\frac{\left(2 + 4(1 + 2\sqrt{\sigma}) \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right) \|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1 \\
 & \stackrel{(ii)}{\leq} \sqrt{\frac{\|V'\|_{\infty}^2}{\gamma(1-\gamma)}} 1 + \sqrt{\frac{\left(2 + 4(1 + 2\sqrt{\sigma}) \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right) \|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1 \\
 & \leq 2 \sqrt{\frac{\left(2 + 4(1 + 2\sqrt{\sigma}) \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right)}{(1-\gamma)^3 \gamma^2}} 1 \tag{203}
 \end{aligned}$$

where (i) holds by the triangle inequality, (ii) follows from (136), and the last inequality is obtained by the fact $\|V'\|_{\infty} \leq \frac{1}{1-\gamma}$.

D. Proof of the upper bound with χ^2 divergence: Theorem 3

The proof of Theorem 3 mainly follows the structure of the proof of Theorem 1 in Appendix B.2. Throughout this section, for any nominal transition kernel P , the uncertainty set is taken as (see (10))

$$\mathcal{U}^{\sigma}(P) = \mathcal{U}_{\chi^2}^{\sigma}(P) := \otimes \mathcal{U}_{\chi^2}^{\sigma}(P_{s,a}), \quad \mathcal{U}_{\chi^2}^{\sigma}(P_{s,a}) := \left\{ P'_{s,a} \in \Delta(\mathcal{S}) : \sum_{s' \in \mathcal{S}} \frac{(P'(s' | s, a) - P(s' | s, a))^2}{P(s' | s, a)} \leq \sigma \right\}. \tag{204}$$

D.1. Proof of Theorem 3

In order to control the performance gap $\|V^{*,\sigma} - V^{\widehat{\pi},\sigma}\|_{\infty}$, recall the error decomposition in (49): as long as the iteration number $T \geq \log\left(\frac{1}{(1-\gamma)\varepsilon_{\text{opt}}}\right)$,

$$V^{*,\sigma} - V^{\widehat{\pi},\sigma} \leq \left(V^{\pi^*,\sigma} - \widehat{V}^{\pi^*,\sigma} \right) + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} 1 + \left(\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma} \right), \tag{205}$$

where ε_{opt} (cf. (48)) shall be specified later (which justifies Remark 2). To further control (205), we bound the remaining two terms separately.

D.1.1. Controlling $\|\widehat{V}^{\pi^*,\sigma} - V^{\pi^*,\sigma}\|_{\infty}$ Towards this, recall the bound in (54) which holds for any uncertainty set:

$$\begin{aligned}
 \|\widehat{V}^{\pi^*,\sigma} - V^{\pi^*,\sigma}\|_{\infty} & \leq \gamma \max \left\{ \left\| \left(I - \gamma \widehat{P}^{\pi^*, \widehat{V}} \right)^{-1} \left(\widehat{P}^{\pi^*, V} V^{\pi^*,\sigma} - \underline{P}^{\pi^*, V} V^{\pi^*,\sigma} \right) \right\|_{\infty}, \right. \\
 & \quad \left. \left\| \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \left(\widehat{P}^{\pi^*, V} V^{\pi^*,\sigma} - \underline{P}^{\pi^*, V} V^{\pi^*,\sigma} \right) \right\|_{\infty} \right\}. \tag{206}
 \end{aligned}$$

To control the main term $\widehat{P}^{\pi^*, V} V^{\pi^*,\sigma} - \underline{P}^{\pi^*, V} V^{\pi^*,\sigma}$ in (206), we first introduce an important lemma whose proof is postponed to Appendix D.2.1

LEMMA 14. Consider any $\sigma > 0$ and the uncertainty set $\mathcal{U}^\sigma(\cdot) := \mathcal{U}_{\chi^2}^\sigma(\cdot)$. For any $\delta \in (0, 1)$, one has with probability at least $1 - \delta$,

$$\begin{aligned} \left| \widehat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right|_\infty &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{*, \sigma})} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} 1 + 4\sqrt{\frac{\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}} 1 \\ &\leq \sqrt{\frac{2 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} 1 + 4\sqrt{\frac{\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}} 1. \end{aligned}$$

Step 1: controlling the first term in (206). Armed with the above lemma, now we control the first term on the right hand side of (206) as follows:

$$\begin{aligned} &(I - \gamma \widehat{\underline{P}}^{\pi^*, V})^{-1} (\widehat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma}) \\ &\stackrel{(i)}{\leq} (I - \gamma \widehat{\underline{P}}^{\pi^*, V})^{-1} \left| \widehat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right|_\infty \\ &\stackrel{(ii)}{\leq} (I - \gamma \widehat{\underline{P}}^{\pi^*, V})^{-1} \left(2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{*, \sigma})} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} 1 + 4\sqrt{\frac{\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}} 1 \right) \\ &= 2(I - \gamma \widehat{\underline{P}}^{\pi^*, V})^{-1} \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{*, \sigma})} \\ &\quad + \left(\frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} + 4\sqrt{\frac{\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}} \right) (I - \gamma \widehat{\underline{P}}^{\pi^*, V})^{-1} 1 \\ &\leq \left(\frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} + 4\sqrt{\frac{\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}} \right) (I - \gamma \widehat{\underline{P}}^{\pi^*, V})^{-1} 1 \\ &\quad + \underbrace{2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} (I - \gamma \widehat{\underline{P}}^{\pi^*, V})^{-1} \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*, V}}(V^{*, \sigma})}}_{=: \mathcal{F}_1} \\ &\quad + \underbrace{2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} (I - \gamma \widehat{\underline{P}}^{\pi^*, V})^{-1} \sqrt{|\text{Var}_{\widehat{P}^{\pi^*}}(V^{*, \sigma}) - \text{Var}_{\widehat{\underline{P}}^{\pi^*, V}}(V^{*, \sigma})|}}_{=: \mathcal{F}_2} \\ &\quad + \underbrace{2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} (I - \gamma \widehat{\underline{P}}^{\pi^*, V})^{-1} (\sqrt{\text{Var}_{P^{\pi^*}}(V^{*, \sigma})} - \sqrt{\text{Var}_{\widehat{P}^{\pi^*}}(V^{*, \sigma})})}_{=: \mathcal{F}_3}, \end{aligned} \tag{207}$$

where (i) holds by $(I - \gamma \widehat{\underline{P}}^{\pi^*, V})^{-1} \geq 0$, (ii) follows from Lemma 14, and the last inequality can be obtained similarly as (90).

We shall control the three terms $\mathcal{F}_1, \mathcal{F}_2, \mathcal{F}_3$ in (207) separately.

- Consider \mathcal{F}_1 . We first introduce the following lemma, whose proof can be found in Appendix C.3.3

LEMMA 15. Consider any $\delta \in (0, 1)$. With probability at least $1 - \delta$, one has

$$(I - \gamma \widehat{\underline{P}}^{\pi^*, V})^{-1} \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*, V}}(V^{*, \sigma})} \leq 2\sqrt{\frac{\left(2 + 4(1 + 2\sqrt{\sigma}) \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right)}{(1-\gamma)^3 \gamma^2}} 1.$$

Applying Lemma [15](#) leads to

$$\begin{aligned} \mathcal{F}_1 &= 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, V}}(V^{*, \sigma})} \\ &\leq 4\sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^3 N}} \left(2 + 4(1 + 2\sqrt{\sigma}) \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right) 1. \end{aligned} \quad (208)$$

- Consider \mathcal{F}_2 . For all $(s, a) \in \mathcal{S} \times \mathcal{A}$, $P_{s,a} \in \Delta(\mathcal{S})$, and $\tilde{P}_{s,a} \in \mathcal{U}^\sigma(P_{s,a})$:

$$|\text{Var}_{\tilde{P}_{s,a}}(V^{*, \sigma}) - \text{Var}_{P_{s,a}}(V^{*, \sigma})| \leq \|\tilde{P}_{s,a} - P_{s,a}\|_1 \|V^{*, \sigma}\|_\infty^2 \leq \frac{\sqrt{\sigma}}{(1-\gamma)^2}, \quad (209)$$

where the last inequality holds by the fact that $\|\tilde{P}_{s,a} - P_{s,a}\|_1 \leq \sqrt{\rho_{\chi^2}(\tilde{P}_{s,a}, P_{s,a})}$, and $\|V'\|_\infty \leq \frac{1}{1-\gamma}$.

Applying the above relation and following the same routine in [95](#) give

$$\begin{aligned} \mathcal{F}_2 &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \sqrt{\|\text{Var}_{\tilde{P}_0}(V^{*, \sigma}) - \text{Var}_{\widehat{P}^{\pi^*, V}}(V^{*, \sigma})\|_\infty} 1 \\ &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \sqrt{\frac{\sqrt{\sigma}}{(1-\gamma)^2}} 1 = 2\sqrt{\frac{\sqrt{\sigma} \log(\frac{18SAN}{\delta})}{(1-\gamma)^4 N}} 1, \end{aligned} \quad (210)$$

where the last equality uses $\left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} 1 = \frac{1}{1-\gamma}$ (cf. [91](#)).

- Consider \mathcal{F}_3 . Applying Lemma [10](#) with $\pi = \pi^*$ and $V = V^{*, \sigma}$ leads to

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

which can be plugged in to verify similar to [96](#) as

$$\mathcal{F}_3 \leq \frac{4 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1. \quad (211)$$

Finally, inserting the results in [208](#), [210](#), and [211](#) back into [207](#) gives

$$\begin{aligned} &\left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \left(\widehat{P}^{\pi^*, V} V^{\pi^*, \sigma} - P^{\pi^*, V} V^{\pi^*, \sigma} \right) \leq \left(\frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)^2} + 4\sqrt{\frac{\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^4 N}} \right) 1 \\ &\quad + 4\sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^3 N}} \left(2 + 4(1 + 2\sqrt{\sigma}) \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right) 1 + \frac{4 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1 + 2\sqrt{\frac{\sqrt{\sigma} \log(\frac{18SAN}{\delta})}{(1-\gamma)^4 N}} 1 \\ &\leq \frac{5 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1 + 4\sqrt{\frac{\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^4 N}} 1 + 2\sqrt{\frac{\sqrt{\sigma} \log(\frac{18SAN}{\delta})}{(1-\gamma)^4 N}} 1 \\ &\quad + 4\sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^3 N}} \left(2 + 4(1 + 2\sqrt{\sigma}) \right) 1 \\ &\leq 48\sqrt{\frac{\log(\frac{24SAN}{\delta})}{(1-\gamma)^3 N}} \left(1 + \frac{\sigma^{1/2} + \sigma^{1/4}}{\sqrt{1-\gamma}} \right) 1 + \frac{5 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1, \end{aligned} \quad (212)$$

where the last inequality holds by the fact $\gamma \geq \frac{1}{4}$ and letting $N \geq \frac{\log(\frac{SAN}{\delta})}{(1-\gamma)^2}$.

Step 2: bounding the second term in (206). Applying Lemma 14 to the second term on the right hand side of (206) leads to

$$\begin{aligned}
& \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \left(\widehat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma}\right) \\
& \leq \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \left(2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{\pi^*, \sigma})} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} 1 + 4\sqrt{\frac{\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}} 1\right) \\
& \leq \left(\frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} + 4\sqrt{\frac{\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}}\right) \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} 1 \\
& \quad + \underbrace{2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*, \widehat{V}}}(V^{\pi^*, \sigma})}}_{=: \mathcal{F}_4} \\
& \quad + \underbrace{2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \left(\sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*, \widehat{V}}}(V^{\pi^*, \sigma} - \widehat{V}^{\pi^*, \sigma})}\right)}_{=: \mathcal{F}_5} \\
& \quad + \underbrace{2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \left(\sqrt{\left|\text{Var}_{\widehat{\underline{P}}^{\pi^*}}(V^{\pi^*, \sigma}) - \text{Var}_{\widehat{\underline{P}}^{\pi^*, \widehat{V}}}(V^{\pi^*, \sigma})\right|}\right)}_{=: \mathcal{F}_6} \\
& \quad + \underbrace{2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \left(\sqrt{\text{Var}_{P^{\pi^*}}(V^{\pi^*, \sigma})} - \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*}}(V^{\pi^*, \sigma})}\right)}_{=: \mathcal{F}_7}. \tag{213}
\end{aligned}$$

We now control the above terms separately.

- Applying Lemma 7 and in view of (91), the term \mathcal{F}_4 in (213) can be controlled similarly to (99) as follows:

$$\mathcal{F}_4 = 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*, \widehat{V}}}(V^{\pi^*, \sigma})} \leq 8\sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^3 N}} 1. \tag{214}$$

- In view of (91), we have

$$\begin{aligned}
\mathcal{F}_5 &= 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{\underline{P}}^{\pi^*, \widehat{V}}\right)^{-1} \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*, \widehat{V}}}(V^{\pi^*, \sigma} - \widehat{V}^{\pi^*, \sigma})} \\
&\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \|V^{\pi^*, \sigma} - \widehat{V}^{\pi^*, \sigma}\|_{\infty} 1. \tag{215}
\end{aligned}$$

- Then, it is easily verified that \mathcal{F}_6 can be controlled similarly to (210) as follows:

$$\mathcal{F}_6 \leq 2\sqrt{\frac{\sqrt{\sigma} \log(\frac{18SAN}{\delta})}{(1-\gamma)^4 N}} 1. \tag{216}$$

- Similarly, \mathcal{F}_7 can be controlled the same as (211) shown below:

$$\mathcal{F}_7 \leq \frac{4\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1. \tag{217}$$

Plugging in the results in (214), (215), (216), and (217) to (213) gives

$$\begin{aligned}
& \left(I - \gamma \widehat{P}^{\pi^*, \widehat{V}} \right)^{-1} \left(\widehat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \\
& \leq \left(\frac{\log\left(\frac{18SAN}{\delta}\right)}{N(1-\gamma)^2} + 4\sqrt{\frac{\sigma \log\left(\frac{24SAN}{\delta}\right)}{(1-\gamma)^4 N}} \right) 1 + 8\sqrt{\frac{\log\left(\frac{18SAN}{\delta}\right)}{\gamma^2(1-\gamma)^3 N}} 1 \\
& \quad + 2\sqrt{\frac{\log\left(\frac{18SAN}{\delta}\right)}{(1-\gamma)^2 N}} \left\| V^{*, \sigma} - \widehat{V}^{\pi^*, \sigma} \right\|_{\infty} 1 + 2\sqrt{\frac{\sqrt{\sigma} \log\left(\frac{18SAN}{\delta}\right)}{(1-\gamma)^4 N}} 1 + \frac{4 \log\left(\frac{18SAN}{\delta}\right)}{(1-\gamma)^2 N} 1 \\
& \leq 2\sqrt{\frac{\log\left(\frac{18SAN}{\delta}\right)}{(1-\gamma)^2 N}} \left\| V^{*, \sigma} - \widehat{V}^{\pi^*, \sigma} \right\|_{\infty} 1 + \frac{5 \log\left(\frac{18SAN}{\delta}\right)}{(1-\gamma)^2 N} 1 + 32\sqrt{\frac{\log\left(\frac{24SAN}{\delta}\right)}{(1-\gamma)^3 N}} \left(1 + \frac{\sigma^{1/2} + \sigma^{1/4}}{\sqrt{1-\gamma}} \right) 1,
\end{aligned} \tag{218}$$

where the last inequality follows from the assumption $\gamma \geq \frac{1}{4}$.

Finally, inserting (212) and (218) back to (206) yields

$$\begin{aligned}
& \left\| \widehat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma} \right\|_{\infty} \\
& \leq \max \left\{ 48\sqrt{\frac{\log\left(\frac{24SAN}{\delta}\right)}{(1-\gamma)^3 N}} \left(1 + \frac{\sigma^{1/2} + \sigma^{1/4}}{\sqrt{1-\gamma}} \right) + \frac{5 \log\left(\frac{18SAN}{\delta}\right)}{(1-\gamma)^2 N}, \right. \\
& \quad \left. 2\sqrt{\frac{\log\left(\frac{18SAN}{\delta}\right)}{(1-\gamma)^2 N}} \left\| V^{*, \sigma} - \widehat{V}^{\pi^*, \sigma} \right\|_{\infty} + \frac{5 \log\left(\frac{18SAN}{\delta}\right)}{(1-\gamma)^2 N} + 32\sqrt{\frac{\log\left(\frac{24SAN}{\delta}\right)}{(1-\gamma)^3 N}} \left(1 + \frac{\sigma^{1/2} + \sigma^{1/4}}{\sqrt{1-\gamma}} \right) \right\} \\
& \leq 96\sqrt{\frac{\log\left(\frac{24SAN}{\delta}\right)}{(1-\gamma)^3 N}} \left(1 + \frac{\sigma^{1/2} + \sigma^{1/4}}{\sqrt{1-\gamma}} \right) + \frac{10 \log\left(\frac{18SAN}{\delta}\right)}{(1-\gamma)^2 N},
\end{aligned} \tag{219}$$

where the last inequality holds by taking $N \geq \frac{16 \log\left(\frac{SAN}{\delta}\right)}{(1-\gamma)^2}$ and rearranging terms.

D.1.2. Controlling $\left\| \widehat{V}^{\widehat{\pi}, \sigma} - V^{\widehat{\pi}, \sigma} \right\|_{\infty}$ Recall the bound in (55) which holds for any uncertainty set:

$$\begin{aligned}
\left\| \widehat{V}^{\widehat{\pi}, \sigma} - V^{\widehat{\pi}, \sigma} \right\|_{\infty} & \leq \gamma \max \left\{ \left\| \left(I - \gamma \underline{P}^{\widehat{\pi}, V} \right)^{-1} \left(\widehat{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - \underline{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right) \right\|_{\infty}, \right. \\
& \quad \left. \left\| \left(I - \gamma \widehat{P}^{\widehat{\pi}, \widehat{V}} \right)^{-1} \left(\widehat{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - \underline{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right) \right\|_{\infty} \right\}.
\end{aligned} \tag{220}$$

To begin with, we introduce the following lemma which controls the main term on the right hand side of (220), which is proved in Appendix D.2.2.

LEMMA 16. Consider any $\delta \in (0, 1)$. Taking $N \geq \log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)$, with probability at least $1 - \delta$, one has

$$\begin{aligned}
& \left| \widehat{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - \underline{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right| \\
& \leq 2\sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N}} \sqrt{\text{Var}_{P^{\widehat{\pi}, \sigma}}(\widehat{V}^{*, \sigma})} 1 + \frac{8 \log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N(1-\gamma)} 1 + 6\sqrt{\frac{2\sigma \log\left(\frac{36SAN^2}{\delta}\right)}{(1-\gamma)^2 N}} 1 + \frac{2\gamma \varepsilon_{\text{opt}} + 4\sqrt{\sigma \varepsilon_{\text{opt}}}}{1-\gamma} 1 \\
& \leq 10\sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{(1-\gamma)^2 N}} 1 + 6\sqrt{\frac{2\sigma \log\left(\frac{36SAN^2}{\delta}\right)}{(1-\gamma)^2 N}} 1 + \frac{2\gamma \varepsilon_{\text{opt}} + 4\sqrt{\sigma \varepsilon_{\text{opt}}}}{1-\gamma} 1.
\end{aligned} \tag{221}$$

Step 3: controlling the first term in (220). Applying Lemma 16 leads to

$$\begin{aligned}
& \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \left(\underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right) \\
& \stackrel{(i)}{\leq} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \left| \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right| \\
& \leq 2 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \sqrt{\text{Var}_{P^{\hat{\pi}}}(\hat{V}^{\star, \sigma})} \\
& \quad + \left(I - \gamma P_Q^{\hat{\pi}, V^{\hat{\pi}}} \right)^{-1} \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} + 6 \sqrt{\frac{2\sigma \log(\frac{36SAN^2}{\delta})}{(1-\gamma)^2 N}} + \frac{2\gamma \varepsilon_{\text{opt}} + 4\sqrt{\sigma \varepsilon_{\text{opt}}}}{1-\gamma} \right) 1 \\
& \stackrel{(ii)}{\leq} \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} + 6 \sqrt{\frac{2\sigma \log(\frac{36SAN^2}{\delta})}{(1-\gamma)^4 N}} + \frac{2\gamma \varepsilon_{\text{opt}} + 4\sqrt{\sigma \varepsilon_{\text{opt}}}}{(1-\gamma)^2} \right) 1 \\
& \quad + 2 \underbrace{\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})}}_{=: \mathcal{G}_1} \\
& \quad + 2 \underbrace{\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \sqrt{\left| \text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\star, \sigma}) - \text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma}) \right|}}_{=: \mathcal{G}_2} \\
& \quad + 2 \underbrace{\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \sqrt{\left| \text{Var}_{P^{\hat{\pi}}}(\hat{V}^{\star, \sigma}) - \text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\star, \sigma}) \right|}}_{=: \mathcal{G}_3}, \tag{222}
\end{aligned}$$

where (i) and (ii) hold by the fact that each row of $(1-\gamma) \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1}$ is a probability vector that falls into $\Delta(\mathcal{S})$.

Therefore, the remainder of the proof will focus on controlling $\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3$ separately.

- For \mathcal{G}_1 , we introduce the following lemma, whose proof is postponed to Appendix D.2.3

LEMMA 17. Consider any $\delta \in (0, 1)$. Taking $N \geq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2}$ and $\varepsilon_{\text{opt}} \leq \frac{1-\gamma}{\gamma}$, one has with probability at least $1 - \delta$,

$$\left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} \leq 7 \sqrt{\frac{1}{(1-\gamma)^3 \gamma^2}} + 6 \sqrt{\frac{\sqrt{\sigma}}{(1-\gamma)^3 \gamma^2}} 1.$$

Applying Lemma 17 and (91) to (222) leads to

$$\begin{aligned}
\mathcal{G}_1 &= 2 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} \\
&\leq 14 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2 (1-\gamma)^3 N}} + 12 \sqrt{\frac{\sqrt{\sigma} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2 (1-\gamma)^3 N}} 1. \tag{223}
\end{aligned}$$

- Applying Lemma 5 with $\|\widehat{V}^{*,\sigma} - \widehat{V}^{\widehat{\pi},\sigma}\|_\infty \leq \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}$ and (91), \mathcal{G}_2 can be controlled as

$$\begin{aligned}\mathcal{G}_2 &= 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}}\right)^{-1} \sqrt{\left|\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{*,\sigma}) - \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})\right|} \\ &\leq 4\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}}\right)^{-1} \frac{\sqrt{\gamma\varepsilon_{\text{opt}}}}{1-\gamma} \leq 4\sqrt{\frac{\gamma\varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1.\end{aligned}\quad (224)$$

- \mathcal{G}_3 can be controlled similarly to \mathcal{F}_2 in (95) as follows:

$$\begin{aligned}\mathcal{G}_3 &= 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}}\right)^{-1} \sqrt{\left|\text{Var}_{\underline{P}^{\widehat{\pi}}}(\widehat{V}^{*,\sigma}) - \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{*,\sigma})\right|} \\ &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}}\right)^{-1} \sqrt{\frac{\sqrt{\sigma}}{(1-\gamma)^2}} 1 = 2\sqrt{\frac{\sqrt{\sigma} \log(\frac{54SAN^2}{\delta})}{(1-\gamma)^4 N}} 1.\end{aligned}\quad (225)$$

To proceed, summing up the results in (223), (224), and (225) and inserting them back to (222) yields: taking $N \geq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2}$ and $\varepsilon_{\text{opt}} \leq \frac{(1-\gamma)^2}{\gamma}$, with probability at least $1 - \delta$,

$$\begin{aligned}&\left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}}\right)^{-1} \left(\widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma}\right) \\ &\leq \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)^2\delta})}{N(1-\gamma)} + 6\sqrt{\frac{2\sigma \log(\frac{36SAN^2}{\delta})}{(1-\gamma)^4 N}} + \frac{2\gamma\varepsilon_{\text{opt}} + 4\sqrt{\sigma\varepsilon_{\text{opt}}}}{(1-\gamma)^2}\right) 1 \\ &\quad + 14\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^3 N}} + 12\sqrt{\frac{\sqrt{\sigma} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^3 N}} 1 + 4\sqrt{\frac{\gamma\varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1 + 2\sqrt{\frac{\sqrt{\sigma} \log(\frac{54SAN^2}{\delta})}{(1-\gamma)^4 N}} 1 \\ &\leq 18 \left(1 + \frac{\sigma^{1/2} + \sigma^{1/4}}{\sqrt{1-\gamma}}\right) \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^3 N}} 1 + \frac{10 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} 1,\end{aligned}\quad (226)$$

where the last inequality holds by taking $\varepsilon_{\text{opt}} \leq \min\left\{\frac{1-\gamma}{\gamma^3}, \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}\right\} = \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}$.

Step 4: bounding the second term in (220). Towards this, applying Lemma 16 leads to

$$\begin{aligned}&\left(I - \gamma \underline{P}^{\widehat{\pi},V}\right)^{-1} \left(\widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma}\right) \leq \left(I - \gamma \underline{P}^{\widehat{\pi},V}\right)^{-1} \left|\widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma}\right| \\ &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},V}\right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi}}}(\widehat{V}^{*,\sigma})} \\ &\quad + \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} + 6\sqrt{\frac{2\sigma \log(\frac{36SAN^2}{\delta})}{(1-\gamma)^2 N}} + \frac{2\gamma\varepsilon_{\text{opt}} + 4\sqrt{\sigma\varepsilon_{\text{opt}}}}{1-\gamma}\right) 1 \\ &\leq \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} + 6\sqrt{\frac{2\sigma \log(\frac{36SAN^2}{\delta})}{(1-\gamma)^2 N}} + \frac{2\gamma\varepsilon_{\text{opt}} + 4\sqrt{\sigma\varepsilon_{\text{opt}}}}{1-\gamma}\right) 1 \\ &\quad + 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \underbrace{\left(I - \gamma \underline{P}^{\widehat{\pi},V}\right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi},V}}(V^{\widehat{\pi},\sigma})}}_{=:\mathcal{G}_4}\end{aligned}$$

$$\begin{aligned}
& + 2 \underbrace{\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} (I - \gamma \underline{P}^{\hat{\pi}, V})^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, V}}(\hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma})}}_{=: \mathcal{G}_5} \\
& + 2 \underbrace{\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} (I - \gamma \underline{P}^{\hat{\pi}, \hat{V}})^{-1} \sqrt{|\text{Var}_{\underline{P}^{\hat{\pi}, V}}(\hat{V}^{*, \sigma}) - \text{Var}_{\underline{P}^{\hat{\pi}, V}}(\hat{V}^{\hat{\pi}, \sigma})|}}_{=: \mathcal{G}_6} \\
& + 2 \underbrace{\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} (I - \gamma \underline{P}^{\hat{\pi}, \hat{V}})^{-1} \sqrt{|\text{Var}_{\underline{P}^{\hat{\pi}}}(\hat{V}^{*, \sigma}) - \text{Var}_{\underline{P}^{\hat{\pi}, V}}(\hat{V}^{*, \sigma})|}}_{=: \mathcal{G}_7}. \tag{227}
\end{aligned}$$

We shall bound each of the terms separately.

- The term \mathcal{G}_4 can be controlled similarly to (214) as follows:

$$\mathcal{G}_4 \leq 8 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^3 N}} 1. \tag{228}$$

- For \mathcal{G}_5 , it is observed that

$$\begin{aligned}
\mathcal{G}_5 & = 2 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} (I - \gamma \underline{P}^{\hat{\pi}, V})^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, V}}(\hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma})} \\
& \leq 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \|V^{\hat{\pi}, \sigma} - \hat{V}^{\hat{\pi}, \sigma}\|_{\infty} 1. \tag{229}
\end{aligned}$$

- Observing that \mathcal{G}_6 and \mathcal{G}_7 are almost the same as the terms \mathcal{G}_2 (controlled in (224)) and \mathcal{G}_3 (controlled in (225)), it is easily verified that they can be controlled as follows

$$\mathcal{G}_6 \leq 4 \sqrt{\frac{\gamma \varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1, \quad \mathcal{G}_7 \leq 2 \sqrt{\frac{\sqrt{\sigma} \log(\frac{54SAN^2}{\delta})}{(1-\gamma)^4 N}} 1. \tag{230}$$

To continue, inserting the results in (228), (229), and (230) back to (227) leads to

$$\begin{aligned}
& (I - \gamma \underline{P}^{\hat{\pi}, V})^{-1} (\hat{\underline{P}}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma}) \\
& \leq \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} + 6 \sqrt{\frac{2\sigma \log(\frac{36SAN^2}{\delta})}{(1-\gamma)^2 N}} + \frac{2\gamma \varepsilon_{\text{opt}} + 4\sqrt{\sigma} \varepsilon_{\text{opt}}}{1-\gamma} \right) 1 + 8 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^3 N}} 1 \\
& \quad + 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \|V^{\hat{\pi}, \sigma} - \hat{V}^{\hat{\pi}, \sigma}\|_{\infty} 1 + 4 \sqrt{\frac{\gamma \varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1 + 2 \sqrt{\frac{\sqrt{\sigma} \log(\frac{54SAN^2}{\delta})}{(1-\gamma)^4 N}} 1 \\
& \leq 16 \left(1 + \frac{\sigma^{1/2} + \sigma^{1/4}}{\sqrt{1-\gamma}} \right) \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^3 N}} 1 + \frac{10 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} 1 + 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \|V^{\hat{\pi}, \sigma} - \hat{V}^{\hat{\pi}, \sigma}\|_{\infty} 1, \tag{231}
\end{aligned}$$

where the last inequality holds by letting $\varepsilon_{\text{opt}} \leq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}$, which directly satisfies $\varepsilon_{\text{opt}} \leq \frac{1-\gamma}{\gamma^3}$ by letting $N \geq \frac{\log(\frac{54SAN^2}{\delta})}{1-\gamma}$.

Finally, inserting (226) and (231) back to (220) yields: taking $\varepsilon_{\text{opt}} \leq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}$ and $N \geq \frac{16 \log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2}$, with probability at least $1 - \delta$, one has

$$\begin{aligned}
& \left\| \widehat{V}^{\widehat{\pi}, \sigma} - V^{\widehat{\pi}, \sigma} \right\|_{\infty} \\
& \leq \max \left\{ 18 \left(1 + \frac{\sigma^{1/2} + \sigma^{1/4}}{\sqrt{1-\gamma}} \right) \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^3 N}} + \frac{10 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2}, \right. \\
& \quad \left. 16 \left(1 + \frac{\sigma^{1/2} + \sigma^{1/4}}{\sqrt{1-\gamma}} \right) \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^3 N}} + \frac{10 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} + 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\widehat{\pi}, \sigma} - \widehat{V}^{\widehat{\pi}, \sigma} \right\|_{\infty} \right\} \\
& \leq 36 \left(1 + \frac{\sigma^{1/2} + \sigma^{1/4}}{\sqrt{1-\gamma}} \right) \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^3 N}} + \frac{20 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2}. \tag{232}
\end{aligned}$$

Step 5: summing up the results. Inserting the results in (232) and (219) back to (205) completes the proof as follows: taking $\varepsilon_{\text{opt}} \leq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}$ and $N \geq \frac{16 \log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2}$, with probability at least $1 - \delta$,

$$\begin{aligned}
\left\| V^{*, \sigma} - V^{\widehat{\pi}, \sigma} \right\|_{\infty} & \leq \left\| V^{\pi^*, \sigma} - \widehat{V}^{\pi^*, \sigma} \right\|_{\infty} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} + \left\| \widehat{V}^{\widehat{\pi}, \sigma} - V^{\widehat{\pi}, \sigma} \right\|_{\infty} \\
& \leq \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} + 96 \sqrt{\frac{\log(\frac{24SAN}{\delta})}{(1-\gamma)^3 N}} \left(1 + \frac{\sigma^{1/2} + \sigma^{1/4}}{\sqrt{1-\gamma}} \right) + \frac{10 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} \\
& \quad + 36 \left(1 + \frac{\sigma^{1/2} + \sigma^{1/4}}{\sqrt{1-\gamma}} \right) \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^3 N}} + \frac{20 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} \\
& \leq 132 \left(1 + \frac{\sigma^{1/2} + \sigma^{1/4}}{\sqrt{1-\gamma}} \right) \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^3 N}} + \frac{22 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} \\
& \leq 550 \left(1 + \frac{\sigma^{1/2} + \sigma^{1/4}}{\sqrt{1-\gamma}} \right) \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^3 N}}, \tag{233}
\end{aligned}$$

where the last inequality holds by $\gamma \geq \frac{1}{4}$ and $N \geq \frac{16 \log(\frac{54SAN^2}{\delta})}{(1-\gamma)}$.

D.2. Proof of the auxiliary lemmas

D.2.1. Proof of Lemma 14 Without loss of generality, we focus on a more general form that considers any fixed deterministic policy π .

Step 1: controlling the point-wise concentration. Consider any fixed policy π and the corresponding robust value vector $V := V^{\pi, \sigma}$ (independent from \widehat{P}^0). Invoking Lemma 2 leads to that for any $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned}
& \left| \widehat{P}_{s,a}^{\pi, V} V^{\pi, \sigma} - P_{s,a}^{\pi, V} V^{\pi, \sigma} \right| \\
& = \left| \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P_{s,a}^0 [V]_{\alpha} - \sqrt{\sigma \text{Var}_{P_{s,a}^0}([V]_{\alpha})} \right\} \right. \\
& \quad \left. - \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ \widehat{P}_{s,a}^0 [V]_{\alpha} - \sqrt{\sigma \text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha})} \right\} \right|
\end{aligned}$$

$$\begin{aligned}
&\leq \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_\alpha + \sqrt{\sigma \text{Var}_{\widehat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\sigma \text{Var}_{P_{s,a}^0}([V]_\alpha)} \right| \\
&\leq \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_\alpha \right| + \\
&\quad + \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \sqrt{\sigma} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right|, \tag{234}
\end{aligned}$$

where the first inequality follows by that the maximum operator is 1-Lipschitz, and the second inequality follows from the triangle inequality. Observing that the first term in (234) is exactly the same as (123), recalling the fact in (128) directly leads to: with probability at least $1 - \delta$,

$$\begin{aligned}
\max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_\alpha \right| &\leq 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} \\
&\leq 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \tag{235}
\end{aligned}$$

holds for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. Then the remainder of the proof focuses on controlling the second term in (234).

Step 2: controlling the second term in (234). For any given $(s, a) \in \mathcal{S} \times \mathcal{A}$ and fixed $\alpha \in [0, \frac{1}{1-\gamma}]$, applying the concentration inequality (Panaganti and Kalathil 2022, Lemma 6) with $\|[V]_\alpha\|_\infty \leq \frac{1}{1-\gamma}$, we arrive at

$$\left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right| \leq \sqrt{\frac{2 \log(\frac{2}{\delta})}{(1-\gamma)^2 N}} \tag{236}$$

holds with probability at least $1 - \delta$. To obtain a uniform bound, we first observe the following lemma proven in Appendix D.2.4.

LEMMA 18. For any V obeying $\|V\|_\infty \leq \frac{1}{1-\gamma}$, the function $J_{s,a}(\alpha, V) := \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right|$ w.r.t. α obeys

$$|J_{s,a}(\alpha_1, V) - J_{s,a}(\alpha_2, V)| \leq 4 \sqrt{\frac{|\alpha_1 - \alpha_2|}{1-\gamma}}.$$

In addition, we can construct an ε_3 -net N_{ε_3} over $[0, \frac{1}{1-\gamma}]$ whose size is $|N_{\varepsilon_3}| \leq \frac{3}{\varepsilon_3(1-\gamma)}$ (Vershynin 2018). Armed with the above, we can derive the uniform bound over $\alpha \in [\min_s V(s), \max_s V(s)] \subset [0, 1/(1-\gamma)]$: with probability at least $1 - \frac{\delta}{SA}$, it holds that for any $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned}
&\max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right| \\
&\leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right| \\
&\stackrel{(i)}{\leq} 4 \sqrt{\frac{\varepsilon_3}{1-\gamma}} + \sup_{\alpha \in N_{\varepsilon_3}} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right| \\
&\stackrel{(ii)}{\leq} 4 \sqrt{\frac{\varepsilon_3}{1-\gamma}} + \sqrt{\frac{2 \log(\frac{2SA|N_{\varepsilon_3}|}{\delta})}{(1-\gamma)^2 N}}
\end{aligned}$$

$$\stackrel{\text{(iii)}}{\leq} 2\sqrt{\frac{2\log(\frac{2SA|N\varepsilon_3|}{\delta})}{(1-\gamma)^2N}} \leq 2\sqrt{\frac{2\log(\frac{24SAN}{\delta})}{(1-\gamma)^2N}}, \quad (237)$$

where (i) holds by the property of N_{ε_3} , (ii) follows from (236), (iii) arises from taking $\varepsilon_3 = \frac{\log(\frac{2SA|N\varepsilon_3|}{\delta})}{8N(1-\gamma)}$, and the last inequality is verified by $|N_{\varepsilon_3}| \leq \frac{3}{\varepsilon_3(1-\gamma)} \leq 24N$.

Inserting (235) and (237) back to (234) and taking the union bound over $(s, a) \in \mathcal{S} \times \mathcal{A}$, we arrive at that for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, with probability at least $1 - \delta$,

$$\begin{aligned} \left| \widehat{P}_{s,a}^{\pi,V} V - P_{s,a}^{\pi,V} V \right| &\leq \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_{\alpha} \right| + \\ &\quad + \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \sqrt{\sigma \text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha})} - \sqrt{\sigma \text{Var}_{P_{s,a}^0}([V]_{\alpha})} \right| \\ &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} + 4\sqrt{\frac{\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^2N}} \\ &\leq \sqrt{\frac{2\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N}} + 4\sqrt{\frac{\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^2N}}. \end{aligned}$$

Finally, recalling the matrix form leads to: with probability at least $1 - \delta$,

$$\begin{aligned} \left| \widehat{P}^{\pi,V} V - P^{\pi,V} V \right|_{\infty} &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P^{\pi}}(V)} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} 1 + 4\sqrt{\frac{\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^2N}} 1 \\ &\leq \sqrt{\frac{2\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N}} 1 + 4\sqrt{\frac{\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^2N}} 1. \end{aligned}$$

Applying the above results with $\pi = \pi^*$ and $V = V^{*,\sigma}$ completes the proof.

D.2.2. Proof of Lemma 16

Step 1: decomposing the term of interest. The proof follows the routine of the proof of Lemma 11 in Appendix B.3.5. To begin with, for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, following the same arguments of (234) yields

$$\begin{aligned} \left| \widehat{P}_{s,a}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - P_{s,a}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right| &\leq \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha} \right| + \\ &\quad + \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \sqrt{\sigma} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha})} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha})} \right|. \quad (238) \end{aligned}$$

Invoking the fact in the first line of (138) and (152) (for proving Lemma 11), the first term in (238) obeys

$$\begin{aligned} &\max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha} \right| \\ &\leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha} \right| \\ &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma})} + \frac{8\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \end{aligned}$$

$$\leq 10\sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}, \quad (239)$$

by letting $N \geq \log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)$. The remainder of the proof will focus on controlling the second term of (238).

Step 2: controlling the second term of (238). Towards this, we recall the auxiliary robust MDP $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ defined in Appendix B.3.5. Taking the uncertainty set $\mathcal{U}^\sigma(\cdot) := \mathcal{U}_{\chi^2}^\sigma(\cdot)$ for both $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ and $\widehat{\mathcal{M}}_{\text{rob}}$, we recall the corresponding robust Bellman operator $\widehat{\mathcal{T}}_{s,u}^\sigma(\cdot)$ in (141) and the following definition in (142)

$$u^* := \widehat{V}^{*,\sigma}(s) - \gamma \inf_{P \in \mathcal{U}^\sigma(e_s)} \mathcal{P} \widehat{V}^{*,\sigma}. \quad (240)$$

Following the arguments in Appendix B.3.5, 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 $0 \leq \widehat{Q}_{s,u}^{*,\sigma} \leq \frac{1}{1-\gamma}$. In addition, the corresponding robust value function coincides with that of the operator $\widehat{\mathcal{T}}^\sigma(\cdot)$, i.e., $\widehat{V}_{s,u}^{*,\sigma} = \widehat{V}^{*,\sigma}$.

We recall the N_{ε_2} -net over $\left[0, \frac{1}{1-\gamma}\right]$ whose size obeying $|N_{\varepsilon_2}| \leq \frac{3}{\varepsilon_2(1-\gamma)}$ (Vershynin 2018). Then for all $u \in N_{\varepsilon_2}$ and a fixed α , $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ is statistically independent from $\widehat{P}_{s,a}^0$, which indicates the independence between $[\widehat{V}_{s,u}^{*,\sigma}]_\alpha$ and $\widehat{P}_{s,a}^0$. With this in mind, invoking the fact in (237) and taking the union bound over all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and $u \in N_{\varepsilon_2}$ yields that, with probability at least $1 - \delta$,

$$\max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}_{s,u}^{*,\sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}_{s,u}^{*,\sigma}]_\alpha)} \right| \leq 2\sqrt{\frac{2\log\left(\frac{24SAN|N_{\varepsilon_2}|}{\delta}\right)}{(1-\gamma)^2N}} \quad (241)$$

holds for all $(s, a, u) \in \mathcal{S} \times \mathcal{A} \times N_{\varepsilon_2}$.

To continue, we decompose the main part of the second term in (238) as follows:

$$\begin{aligned} & \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{\widehat{\pi}, \sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{\widehat{\pi}, \sigma}]_\alpha)} \right| \\ & \leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{\widehat{\pi}, \sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{\widehat{\pi}, \sigma}]_\alpha)} \right| \\ & \stackrel{(i)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{*,\sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{*,\sigma}]_\alpha)} \right| \\ & \quad + \max_{\alpha \in [0, 1/(1-\gamma)]} \left[\left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{\widehat{\pi}, \sigma}]_\alpha)} - \text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{*,\sigma}]_\alpha) \right| \right. \\ & \quad \left. + \sqrt{\left| \text{Var}_{P_{s,a}^0}([\widehat{V}^{\widehat{\pi}, \sigma}]_\alpha) - \text{Var}_{P_{s,a}^0}([\widehat{V}^{*,\sigma}]_\alpha) \right|} \right] \\ & \stackrel{(ii)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{*,\sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{*,\sigma}]_\alpha)} \right| \\ & \quad + \max_{\alpha \in [0, 1/(1-\gamma)]} 2\sqrt{\frac{2}{(1-\gamma)}} \left\| [\widehat{V}^{\widehat{\pi}, \sigma}]_\alpha - [\widehat{V}^{*,\sigma}]_\alpha \right\|_\infty \\ & \leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{*,\sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{*,\sigma}]_\alpha)} \right| + 4\sqrt{\frac{\varepsilon_{\text{opt}}}{(1-\gamma)^2}}, \end{aligned} \quad (242)$$

where (i) holds by the triangle inequality, (ii) arises from applying Lemma 5, and the last inequality holds by (48).

Armed with the above facts, invoking the identity $\widehat{V}^{*,\sigma} = \widehat{V}_{s,u^*}^{*,\sigma}$ leads to that for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, with probability at least $1 - \delta$,

$$\begin{aligned}
 & \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{*,\sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{*,\sigma}]_\alpha)} \right| \\
 &= \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}_{s,u^*}^{*,\sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}_{s,u^*}^{*,\sigma}]_\alpha)} \right| \\
 &\stackrel{(i)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}_{s,\bar{u}}^{*,\sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}_{s,\bar{u}}^{*,\sigma}]_\alpha)} \right| \\
 &\quad + \max_{\alpha \in [0, 1/(1-\gamma)]} \left[\sqrt{\left| \text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}_{s,u^*}^{*,\sigma}]_\alpha) - \text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}_{s,\bar{u}}^{*,\sigma}]_\alpha) \right|} \right] \\
 &\quad + \sqrt{\left| \text{Var}_{P_{s,a}^0}([\widehat{V}_{s,u^*}^{*,\sigma}]_\alpha) - \text{Var}_{P_{s,a}^0}([\widehat{V}_{s,\bar{u}}^{*,\sigma}]_\alpha) \right|} \\
 &\stackrel{(ii)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}_{s,\bar{u}}^{*,\sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}_{s,\bar{u}}^{*,\sigma}]_\alpha)} \right| + 4\sqrt{\frac{\varepsilon_2}{(1-\gamma)}} \\
 &\stackrel{(iii)}{\leq} 2\sqrt{\frac{2\log(\frac{24SAN|N\varepsilon_2|}{\delta})}{(1-\gamma)^2N}} + 4\sqrt{\frac{\varepsilon_2}{(1-\gamma)}} \\
 &\leq 6\sqrt{\frac{2\log(\frac{36SAN^2}{\delta})}{(1-\gamma)^2N}}, \tag{243}
 \end{aligned}$$

where (i) holds by the triangle inequality, (ii) arises from applying Lemma 5 and the fact $\left\| \widehat{V}_{s,\bar{u}}^{*,\sigma} - \widehat{V}_{s,u^*}^{*,\sigma} \right\| \leq \frac{\varepsilon_2}{(1-\gamma)}$ (see (148)), (iii) follows from (241), and the last inequality holds by letting $\varepsilon_2 = \frac{2\log(\frac{24SAN|N\varepsilon_2|}{\delta})}{(1-\gamma)N}$, which leads to $|N\varepsilon_2| \leq \frac{3}{\varepsilon_2(1-\gamma)} \leq \frac{3N}{2}$.

In summary, inserting (243) back to (242) leads to with probability at least $1 - \delta$,

$$\begin{aligned}
 & \max_{\alpha \in [\min_s \widehat{V}^{\hat{\pi},\sigma}(s), \max_s \widehat{V}^{\hat{\pi},\sigma}(s)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{\hat{\pi},\sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{\hat{\pi},\sigma}]_\alpha)} \right| \\
 &\leq 6\sqrt{\frac{2\sigma \log(\frac{36SAN^2}{\delta})}{(1-\gamma)^2N}} + 4\sqrt{\frac{\sigma \varepsilon_{\text{opt}}}{(1-\gamma)^2}} \tag{244}
 \end{aligned}$$

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

Step 4: finishing up. Inserting (244) and (239) back to (238), we have for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, with probability at least $1 - \delta$,

$$\begin{aligned}
 & \left| \widehat{P}_{s,a}^{\hat{\pi},\widehat{V}} \widehat{V}^{\hat{\pi},\sigma} - P_{s,a}^{\hat{\pi},\widehat{V}} \widehat{V}^{\hat{\pi},\sigma} \right| \\
 &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma})} + \frac{8\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} + 6\sqrt{\frac{2\sigma \log(\frac{36SAN^2}{\delta})}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}} + 4\sqrt{\sigma\varepsilon_{\text{opt}}}}{1-\gamma} \\
 &\leq 10\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2N}} + 6\sqrt{\frac{2\sigma \log(\frac{36SAN^2}{\delta})}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}} + 4\sqrt{\sigma\varepsilon_{\text{opt}}}}{1-\gamma}. \tag{245}
 \end{aligned}$$

Finally, recalling the matrix form in (40), taking $N \geq \log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)$, we complete the proof:

$$\begin{aligned} & \left| \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right| \\ & \leq 2\sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N}} \sqrt{\text{Var}_{P^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} 1 + \frac{8\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N(1-\gamma)} 1 + 6\sqrt{\frac{2\sigma\log\left(\frac{36SAN^2}{\delta}\right)}{(1-\gamma)^2N}} 1 + \frac{2\gamma\varepsilon_{\text{opt}} + 4\sqrt{\sigma\varepsilon_{\text{opt}}}}{1-\gamma} 1 \\ & \leq 10\sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{(1-\gamma)^2N}} 1 + 6\sqrt{\frac{2\sigma\log\left(\frac{36SAN^2}{\delta}\right)}{(1-\gamma)^2N}} 1 + \frac{2\gamma\varepsilon_{\text{opt}} + 4\sqrt{\sigma\varepsilon_{\text{opt}}}}{1-\gamma} 1. \end{aligned} \quad (246)$$

D.2.3. Proof of Lemma 17 Following the proof pipeline of Lemma 12 in Appendix B.3.6, we first recall (131)

$$\left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}}\right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} \leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left(\underline{P}^{\hat{\pi}, \hat{V}}\right)^t \text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})}. \quad (247)$$

Recall that we denote the minimum value of $\hat{V}^{\hat{\pi}, \sigma}$ as $V_{\min} = \min_{s \in \mathcal{S}} \hat{V}^{\hat{\pi}, \sigma}(s)$ and $V' := \hat{V}^{\hat{\pi}, \sigma} - V_{\min} 1$. By the same argument as (155), we arrive at

$$\begin{aligned} & \text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma}) \\ & \leq \underline{P}^{\hat{\pi}, \hat{V}} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} 1 + \frac{2}{\gamma} \|V'\|_{\infty} \left| \left(\underline{P}^{\hat{\pi}, \hat{V}} - \underline{P}^{\hat{\pi}, \hat{V}}\right) \hat{V}^{\hat{\pi}, \sigma} \right| \\ & \leq \underline{P}^{\hat{\pi}, \hat{V}} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} 1 \\ & \quad + \frac{2}{\gamma} \|V'\|_{\infty} \left(10\sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{(1-\gamma)^2N}} + 6\sqrt{\frac{2\sigma\log\left(\frac{36SAN^2}{\delta}\right)}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}} + 4\sqrt{\sigma\varepsilon_{\text{opt}}}}{1-\gamma} \right) 1, \end{aligned} \quad (248)$$

where the last inequality follows by Lemma 16. Plugging (248) back into (247) and following the routine of (156) leads to

$$\begin{aligned} & \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}}\right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} \\ & \stackrel{(i)}{\leq} \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left(\underline{P}^{\hat{\pi}, \hat{V}}\right)^t \left(\underline{P}^{\hat{\pi}, \hat{V}} (V' \circ V') - \frac{1}{\gamma} V' \circ V'\right)} \\ & \quad + \sqrt{\frac{1}{(1-\gamma)^2\gamma^2} \left(2 + 20\sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{(1-\gamma)^2N}} + 12\sqrt{\frac{2\sigma\log\left(\frac{36SAN^2}{\delta}\right)}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}} + 8\sqrt{\sigma\varepsilon_{\text{opt}}}}{1-\gamma} \right) \|V'\|_{\infty} 1} \\ & \stackrel{(ii)}{\leq} \sqrt{\frac{\|V'\|_{\infty}^2}{\gamma(1-\gamma)}} 1 \\ & \quad + \sqrt{\frac{1}{(1-\gamma)^2\gamma^2} \left(2 + 20\sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{(1-\gamma)^2N}} + 12\sqrt{\frac{2\sigma\log\left(\frac{36SAN^2}{\delta}\right)}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}} + 8\sqrt{\sigma\varepsilon_{\text{opt}}}}{1-\gamma} \right) \|V'\|_{\infty} 1} \\ & \stackrel{(iii)}{\leq} \sqrt{\frac{\|V'\|_{\infty}^2}{\gamma(1-\gamma)}} 1 + \sqrt{\frac{32(1+\sqrt{\sigma})\|V'\|_{\infty}}{(1-\gamma)^2\gamma^2}} 1 \leq 7\sqrt{\frac{1}{(1-\gamma)^3\gamma^2}} + 6\sqrt{\frac{\sqrt{\sigma}}{(1-\gamma)^3\gamma^2}} 1, \end{aligned} \quad (249)$$

where (i) arises from following the routine of (135), (ii) holds by repeating the argument of (136), (iii) follows by taking $N \geq \frac{\log(\frac{54SA^2N^2}{(1-\gamma)^\delta})}{(1-\gamma)^2}$ and $\varepsilon_{\text{opt}} \leq \frac{(1-\gamma)^2}{\gamma}$, and the last inequality holds by $\|V'\|_\infty \leq \|V^{*,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$.

D.2.4. Proof of Lemma 18 For any $0 \leq \alpha_1, \alpha_2 \leq 1/(1-\gamma)$, one has

$$\begin{aligned}
& |J_{s,a}(\alpha_1, V) - J_{s,a}(\alpha_2, V)| \\
&= \left| \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_{\alpha_1})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_1})} \right| - \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_{\alpha_2})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_2})} \right| \right| \\
&\stackrel{(i)}{\leq} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_{\alpha_1})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_1})} - \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_{\alpha_2})} + \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_2})} \right| \\
&\leq \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_{\alpha_1})} - \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_{\alpha_2})} \right| + \left| \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_1})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_2})} \right| \\
&\stackrel{(ii)}{\leq} \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_{\alpha_2}) - \text{Var}_{\hat{P}_{s,a}^0}([V]_{\alpha_1})} + \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_2}) - \text{Var}_{P_{s,a}^0}([V]_{\alpha_1})} \\
&\stackrel{(iii)}{\leq} \sqrt{\left| \hat{P}_{s,a}^0([V]_{\alpha_1}) \circ ([V]_{\alpha_1}) - ([V]_{\alpha_2}) \circ ([V]_{\alpha_2}) \right|} + \sqrt{\left| \hat{P}_{s,a}^0([V]_{\alpha_1} + [V]_{\alpha_2}) \cdot \hat{P}_{s,a}^0([V]_{\alpha_1} - [V]_{\alpha_2}) \right|} \\
&\quad + \sqrt{\left| P_{s,a}^0([V]_{\alpha_1}) \circ ([V]_{\alpha_1}) - ([V]_{\alpha_2}) \circ ([V]_{\alpha_2}) \right|} + \sqrt{\left| P_{s,a}^0([V]_{\alpha_1} + [V]_{\alpha_2}) \cdot P_{s,a}^0([V]_{\alpha_1} - [V]_{\alpha_2}) \right|} \\
&\leq 2\sqrt{2(\alpha_1 + \alpha_2)|\alpha_1 - \alpha_2|} \leq 4\sqrt{\frac{|\alpha_1 - \alpha_2|}{1-\gamma}}. \tag{250}
\end{aligned}$$

Here, (i) holds by the fact $\left| |x| - |y| \right| \leq |x - y|$ for all $x, y \in \mathbb{R}$, (ii) follows from the fact that $\sqrt{x} - \sqrt{y} \leq \sqrt{x - y}$ for any $x \geq y \geq 0$ and $\text{Var}_P([V]_{\alpha_2}) \geq \text{Var}_P([V]_{\alpha_1})$ for any transition kernel $P \in \Delta(\mathcal{S})$, (iii) holds by the definition of $\text{Var}_P(\cdot)$ defined in (38), and the last inequality arises from $0 \leq \alpha_1, \alpha_2 \leq 1/(1-\gamma)$.

E. Proof of the lower bound with χ^2 divergence: Theorem 4

To prove Theorem 4, we shall first construct some hard instances and then characterize the sample complexity requirements over these instances. The structure of the hard instances is the same as the ones used in the proof of Theorem 2.

E.1. Construction of the hard problem instances

First, note that we shall use the same MDPs defined in Section 5.3 as follows

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

In particular, we shall keep the structure of the transition kernel in (56), reward function in (61) and initial state distribution in (62), while p and Δ shall be tailored to the χ^2 divergence case.

Uncertainty set of the transition kernels. Recalling the uncertainty set associated with χ^2 divergence in (204), for any uncertainty level σ , the uncertainty set throughout this section is defined as $\mathcal{U}^\sigma(P^\phi)$:

$$\mathcal{U}^\sigma(P^\phi) := \otimes \mathcal{U}_{\chi^2}^\sigma(P_{s,a}^\phi), \quad \mathcal{U}_{\chi^2}^\sigma(P_{s,a}^\phi) := \left\{ P_{s,a} \in \Delta(\mathcal{S}) : \sum_{s' \in \mathcal{S}} \frac{(P(s' | s, a) - P^\phi(s' | s, a))^2}{P^\phi(s' | s, a)} \leq \sigma \right\}, \quad (251)$$

where $\Delta(\mathcal{S})$ denotes the simplex over the state space \mathcal{S} . Clearly, $\mathcal{U}^\sigma(P_{s,a}^\phi) = \{P_{s,a}^\phi\}$ whenever the state transition is deterministic for χ^2 divergence. Here, q and Δ (recall (59)) which determine the instances are specified as

$$0 \leq q = \begin{cases} 1 - \gamma & \text{if } \sigma \in (0, \frac{1-\gamma}{4}) \\ \frac{\sigma}{1+\sigma} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty) \end{cases}, \quad p = q + \Delta, \quad (252)$$

and

$$0 < \Delta \leq \begin{cases} \frac{1}{4}(1 - \gamma) & \text{if } \sigma \in (0, \frac{1-\gamma}{4}) \\ \min \left\{ \frac{1-\gamma}{\gamma(3+\sigma)}, \frac{1}{2(1+\sigma)} \right\} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty) \end{cases}. \quad (253)$$

This directly ensures that

$$p = \Delta + q \leq \max \left\{ \frac{\frac{1}{2} + \sigma}{1 + \sigma}, \frac{5}{4}(1 - \gamma) \right\} \leq 1$$

since $\gamma \in [\frac{3}{4}, 1)$.

To continue, for any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, we denote the infimum probability of moving to the next state s' associated with any perturbed transition kernel $P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^\phi)$ as

$$\underline{P}^\phi(s' | s, a) := \inf_{P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^\phi)} P(s' | s, a). \quad (254)$$

In addition, we denote the transition from state 0 to state 1 as follows, which plays an important role in the analysis,

$$\underline{p} := \underline{P}^\phi(1 | 0, \phi), \quad \underline{q} := \underline{P}^\phi(1 | 0, 1 - \phi). \quad (255)$$

Before continuing, we introduce some facts about \underline{p} and \underline{q} which are summarized as the following lemma; the proof is postponed to Appendix E.3.1.

LEMMA 19. Consider any $\sigma \in (0, \infty)$ and any p, q, Δ obeying (252) and (253), the following properties hold

$$\begin{cases} \frac{1-\gamma}{2} < \underline{q} < 1 - \gamma, & \underline{q} + \frac{3}{4}\Delta \leq \underline{p} \leq \underline{q} + \Delta \leq \frac{5(1-\gamma)}{4} & \text{if } \sigma \in (0, \frac{1-\gamma}{4}), \\ \underline{q} = 0, & \frac{\sigma+1}{2}\Delta \leq \underline{p} \leq (3 + \sigma)\Delta & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty). \end{cases} \quad (256)$$

Value functions and optimal policies. Armed with the above facts, we are positioned to derive the corresponding robust value functions, the optimal policies, and their corresponding optimal robust value functions. For any RMDP \mathcal{M}_ϕ with the uncertainty set defined in (251), we denote the robust optimal policy as π_ϕ^* , the robust value function of any policy π (resp. the optimal policy π_ϕ^*) as $V_\phi^{\pi, \sigma}$ (resp. $V_\phi^{*, \sigma}$). The following lemma describes some key properties of the robust (optimal) value functions and optimal policies whose proof is postponed to Appendix E.3.2.

LEMMA 20. For any $\phi \in \{0, 1\}$ and any policy π , one has

$$V_\phi^{\pi, \sigma}(0) = \frac{\gamma z_\phi^\pi}{(1 - \gamma)(1 - \gamma(1 - z_\phi^\pi))}, \quad (257)$$

where z_ϕ^π is defined as

$$z_\phi^\pi := \underline{p}\pi(\phi | 0) + \underline{q}\pi(1 - \phi | 0). \quad (258)$$

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

$$V_\phi^{*, \sigma}(0) = \frac{\gamma \underline{p}}{(1 - \gamma)(1 - \gamma(1 - \underline{p}))}, \quad (259a)$$

$$\pi_\phi^*(\phi | s) = 1, \quad \text{for } s \in \mathcal{S}. \quad (259b)$$

E.2. Establishing the minimax lower bound

Our goal is to control the performance gap w.r.t. any policy estimator $\hat{\pi}$ based on the generated dataset and the chosen initial distribution φ in (62), which gives

$$\langle \varphi, V_\phi^{*, \sigma} - V_\phi^{\hat{\pi}, \sigma} \rangle = V_\phi^{*, \sigma}(0) - V_\phi^{\hat{\pi}, \sigma}(0). \quad (260)$$

Step 1: converting the goal to estimate ϕ . To achieve the goal, we first introduce the following fact which shall be verified in Appendix E.3.3: given

$$\varepsilon \leq \frac{1}{768(1 - \gamma)}, \quad (261)$$

and choosing

$$\Delta = \begin{cases} 18(1 - \gamma)^2 \varepsilon & \text{if } \sigma \in (0, \frac{1 - \gamma}{4}), \\ \frac{64\varepsilon(1 - \gamma)^2}{3(1 + \sigma)} & \text{if } \sigma \in [\frac{1 - \gamma}{4}, \infty), \end{cases} \quad (262)$$

which satisfies the requirement of Δ in (253), it holds that for any policy $\hat{\pi}$,

$$\langle \varphi, V_\phi^{*, \sigma} - V_\phi^{\hat{\pi}, \sigma} \rangle \geq 2\varepsilon(1 - \hat{\pi}(\phi | 0)). \quad (263)$$

Step 2: arriving at the final results. To continue, following the same definitions and argument in Appendix C.2, we recall the minimax probability of the error and its property as follows:

$$p_e \geq \frac{1}{4} \exp \left\{ -N \left(\text{KL}(P^0(\cdot|0,0) \| P^1(\cdot|0,0)) + \text{KL}(P^0(\cdot|0,1) \| P^1(\cdot|0,1)) \right) \right\}, \quad (264)$$

then we can complete the proof by showing $p_e \geq \frac{1}{8}$ given the bound for the sample size N . In the following, we shall control the KL divergence terms in (264) in three different cases.

- Case 1: $\sigma \in (0, \frac{1-\gamma}{4})$. In this case, applying $\gamma \in [\frac{3}{4}, 1)$ yields

$$\begin{aligned} 1 - q > 1 - p = 1 - q - \Delta > \gamma - \frac{1-\gamma}{4} > \frac{3}{4} - \frac{1}{16} > \frac{1}{2}, \\ p \geq q = 1 - \gamma. \end{aligned} \quad (265)$$

Armed with the above facts, applying Tsybakov (2009, Lemma 2.7) yields

$$\begin{aligned} \text{KL}(P^0(\cdot|0,1) \| P^1(\cdot|0,1)) &= \text{KL}(p \| q) \leq \frac{(p-q)^2}{(1-p)p} \stackrel{(i)}{=} \frac{\Delta^2}{p(1-p)} \\ &\stackrel{(ii)}{=} \frac{324(1-\gamma)^4 \varepsilon^2}{p(1-p)} \\ &\stackrel{(iii)}{\leq} 648(1-\gamma)^3 \varepsilon^2, \end{aligned} \quad (266)$$

where (i) follows from the definition in (252), (ii) holds by plugging in the expression of Δ in (262), and (iii) arises from (265). The same bound can be established for $\text{KL}(P^0(\cdot|0,0) \| P^1(\cdot|0,0))$. Substituting (266) back into (264) demonstrates that: if the sample size is chosen as

$$N \leq \frac{\log 2}{1296(1-\gamma)^3 \varepsilon^2}, \quad (267)$$

then one necessarily has

$$p_e \geq \frac{1}{4} \exp \left\{ -N \cdot 1296(1-\gamma)^3 \varepsilon^2 \right\} \geq \frac{1}{8}. \quad (268)$$

- Case 2: $\sigma \in [\frac{1-\gamma}{4}, \infty)$. Applying the facts of Δ in (253), one has

$$\begin{aligned} 1 - q > 1 - p = 1 - q - \Delta &\geq \frac{1}{1+\sigma} - \frac{1}{2(1+\sigma)} = \frac{1}{2(1+\sigma)}, \\ p \geq q &= \frac{\sigma}{1+\sigma}. \end{aligned} \quad (269)$$

Given (269), applying Tsybakov (2009, Lemma 2.7) yields

$$\begin{aligned} \text{KL}(P^0(\cdot|0,1) \| P^1(\cdot|0,1)) &= \text{KL}(p \| q) \leq \frac{(p-q)^2}{(1-p)p} \stackrel{(i)}{=} \frac{\Delta^2}{p(1-p)} \\ &\stackrel{(ii)}{=} \frac{\frac{4096\varepsilon^2(1-\gamma)^4}{9(1+\sigma)^2}}{p(1-p)} \\ &\stackrel{(iii)}{\leq} \frac{\frac{4096\varepsilon^2(1-\gamma)^4}{9(1+\sigma)^2}}{\frac{\sigma}{2(1+\sigma)^2}} \leq \frac{8192(1-\gamma)^4 \varepsilon^2}{\sigma}, \end{aligned} \quad (270)$$

where (i) follows from the definition in (252), (ii) holds by plugging in the expression of Δ in (262), and (iii) arises from (269). The same bound can be established for $\text{KL}(P^0(\cdot | 0, 0) \| P^1(\cdot | 0, 0))$.

Substituting (270) back into (264) demonstrates that: if the sample size is chosen as

$$N \leq \frac{\sigma \log 2}{16384(1-\gamma)^4 \varepsilon^2}, \quad (271)$$

then one necessarily has

$$p_e \geq \frac{1}{4} \exp \left\{ -N \frac{16384(1-\gamma)^4 \varepsilon^2}{\sigma} \right\} \geq \frac{1}{8}. \quad (272)$$

Step 3: putting things together. Finally, summing up the results in (267) and (271), combined with the requirement in (261), one has when

$$\varepsilon \leq \frac{c_1}{1-\gamma}, \quad (273)$$

taking

$$N \leq c_2 \begin{cases} \frac{1}{(1-\gamma)^3 \varepsilon^2} & \text{if } \sigma \in (0, \frac{1-\gamma}{4}) \\ \frac{\sigma}{(1-\gamma)^4 \varepsilon^2} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty) \end{cases} \quad (274)$$

leads to $p_e \geq \frac{1}{8}$, for some universal constants $c_1, c_2 > 0$.

E.3. Proof of the auxiliary facts

We begin with some basic facts about the χ^2 divergence defined in (75) for any two Bernoulli distributions $\text{Ber}(w)$ and $\text{Ber}(x)$, denoted as

$$f(w, x) := \chi^2(x \| w) = \frac{(w-x)^2}{w} + \frac{(1-w-(1-x))^2}{1-w} = \frac{(w-x)^2}{w(1-w)}. \quad (275)$$

For $x \in [0, w)$, it is easily verified that the partial derivative w.r.t. x obeys $\frac{\partial f(w, x)}{\partial x} = \frac{2(x-w)}{w(1-w)} < 0$, implying that

$$\forall x_1 < x_2 \in [0, w), \quad f(w, x_1) > f(w, x_2). \quad (276)$$

In other words, the χ^2 divergence $f(w, x)$ increases as x decreases from w to 0.

Next, we introduce the following function for any fixed $\sigma \in (0, \infty)$ and any $x \in [\frac{\sigma}{1+\sigma}, 1)$:

$$f_\sigma(x) := \inf_{\{y: \chi^2(y \| x) \leq \sigma, y \in [0, x]\}} y \stackrel{(i)}{=} \max \left\{ 0, x - \sqrt{\sigma x(1-x)} \right\} = x - \sqrt{\sigma x(1-x)}, \quad (277)$$

where (i) has been verified in Yang et al. (2022, Corollary B.2), and the last equality holds since $x \geq \frac{\sigma}{1+\sigma}$. The next lemma summarizes some useful facts about $f_\sigma(\cdot)$, which again has been verified in Yang et al. (2022, Lemma B.12 and Corollary B.2).

LEMMA 21. Consider any $\sigma \in (0, \infty)$. For $x \in [\frac{\sigma}{1+\sigma}, 1)$, $f_\sigma(x)$ is convex and differentiable, which obeys

$$f'_\sigma(x) = 1 + \frac{\sqrt{\sigma}(2x-1)}{2\sqrt{x(1-x)}}.$$

E.3.1. Proof of Lemma 19 Let us control \underline{q} and \underline{p} respectively.

Step 1: controlling \underline{q} . We shall control \underline{q} in different cases w.r.t. the uncertainty level σ .

- Case 1: $\sigma \in (0, \frac{1-\gamma}{4})$. In this case, recall that $q = 1 - \gamma$ defined in (252), applying (277) with $x = q$ leads to

$$1 - \gamma = q > \underline{q} = f_\sigma(q) = 1 - \gamma - \sqrt{\sigma\gamma(1-\gamma)} \geq 1 - \gamma - \sqrt{\frac{1-\gamma}{4}\gamma(1-\gamma)} > \frac{1-\gamma}{2}. \quad (278)$$

- Case 2: $\sigma \in [1, \infty)$. Note that it suffices to treat $P_{0,1-\phi}^\phi$ as a Bernoulli distribution $\text{Ber}(q)$ over states 1 and 0, since we do not allow transition to other states. Recalling $q = \frac{\sigma}{1+\sigma}$ in (252) and noticing the fact that

$$f(q, 0) = \frac{q^2}{q} + \frac{(1 - (1-q))^2}{1-q} = \frac{q}{(1-q)} = \sigma, \quad (279)$$

one has the probability $\text{Ber}(0)$ falls into the uncertainty set of $\text{Ber}(q)$ of size σ . As a result, recalling the definition (255) leads to

$$\underline{q} = \underline{P}^\phi(1 | 0, 1 - \phi) = 0, \quad (280)$$

since $\underline{q} \geq 0$.

Step 2: controlling \underline{p} . To characterize the value of \underline{p} , we also divide into several cases separately.

- Case 1: $\sigma \in (0, \frac{1-\gamma}{4})$. In this case, note that $p > q = 1 - \gamma \geq \frac{\sigma}{1+\sigma}$. Therefore, applying that $f_\sigma(\cdot)$ is convex and the form of its derivative in Lemma 21, one has

$$\begin{aligned} \underline{p} &= f_\sigma(p) \geq f_\sigma(q) + f'_\sigma(q)(p - q) \\ &= \underline{q} + \left(1 + \frac{\sqrt{\sigma}(2q - 1)}{2\sqrt{q(1-q)}}\right)\Delta \geq \underline{q} + \left(1 - \frac{\sqrt{\frac{1-\gamma}{4}}(1 - 2(1-\gamma))}{2\sqrt{(1-\gamma)\gamma}}\right)\Delta \geq \underline{q} + \frac{3\Delta}{4}. \end{aligned} \quad (281)$$

Similarly, applying Lemma 21 leads to

$$\begin{aligned} \underline{p} &= f_\sigma(p) \leq f_\sigma(q) + f'_\sigma(p)(p - q) \\ &= \underline{q} + \left(1 - \frac{\sqrt{\sigma}(1 - 2p)}{2\sqrt{p(1-p)}}\right)\Delta \leq \underline{q} + \Delta, \end{aligned} \quad (282)$$

where the last inequality holds by $1 - 2p > 0$ due to the fact $p = q + \Delta \leq \frac{5}{4}(1 - \gamma) \leq \frac{5}{16} < \frac{1}{2}$ (cf. (253) and $\gamma \in [\frac{3}{4}, 1)$). To sum up, given $\sigma \in (0, \frac{1-\gamma}{4})$, combined with (278), we arrive at

$$\underline{q} + \frac{3}{4}\Delta \leq \underline{p} \leq \underline{q} + \Delta \leq \frac{5(1-\gamma)}{4}, \quad (283)$$

where the last inequality holds by $\Delta \leq \frac{1}{4}(1 - \gamma)$ (see (253)).

• Case 2: $\sigma \in [\frac{1-\gamma}{4}, \infty)$. We recall that $p = q + \Delta > q = \frac{\sigma}{1+\sigma}$ in (252). To derive the lower bound for \underline{p} in (255), similar to (281), one has

$$\begin{aligned} \underline{p} &= f_\sigma(p) \geq f_\sigma(q) + f'_\sigma(q)(p - q) \\ &= \underline{q} + \left(1 + \frac{\sqrt{\sigma}(2q - 1)}{2\sqrt{q(1 - q)}}\right) \Delta \\ &\stackrel{(i)}{=} 0 + \left(1 + \frac{\sqrt{\sigma} \frac{\sigma - 1}{1 + \sigma}}{2\sqrt{\frac{\sigma}{1 + \sigma} \frac{1}{1 + \sigma}}}\right) \Delta = \left(1 + \frac{\sigma - 1}{2}\right) \Delta = \left(\frac{\sigma + 1}{2}\right) \Delta, \end{aligned} \quad (284)$$

where (i) follows from $q = \frac{\sigma}{1+\sigma}$ and $\underline{q} = 0$ (see (280)). For the other direction, similar to (282), we have

$$\begin{aligned} \underline{p} &= f_\sigma(p) \leq f_\sigma(q) + f'_\sigma(p)(p - q) = \underline{q} + \left(1 + \frac{\sqrt{\sigma}(2p - 1)}{2\sqrt{p(1 - p)}}\right) \Delta \\ &\stackrel{(i)}{=} \left(1 + \frac{\sqrt{\sigma}(2p - 1)}{2\sqrt{p(1 - p)}}\right) \Delta \stackrel{(ii)}{=} \left(1 + \frac{\sqrt{\sigma} \left(\frac{\sigma - 1}{1 + \sigma} + 2\Delta\right)}{2\sqrt{\left(\frac{\sigma}{1 + \sigma} + \Delta\right) \left(\frac{1}{1 + \sigma} - \Delta\right)}}\right) \Delta \\ &\stackrel{(iii)}{\leq} \left(1 + \frac{\sqrt{\sigma}(1 + 2\Delta)}{2\sqrt{\frac{\sigma}{1 + \sigma} \cdot \frac{1}{2(1 + \sigma)}}}\right) \Delta \stackrel{(iv)}{\leq} \left(1 + (1 + \sigma) \left(1 + \frac{1}{1 + \sigma}\right)\right) \Delta = (3 + \sigma)\Delta, \end{aligned} \quad (285)$$

where (i) holds by $\underline{q} = 0$ (see (280)), (ii) follows from plugging in $p = q + \Delta = \frac{\sigma}{1+\sigma} + \Delta$, and (iii) and (iv) arise from $\Delta = \min\left\{\frac{1}{4}(1 - \gamma), \frac{1}{2(1 + \sigma)}\right\} \leq 1$ in (253). Combining (284) and (285) yields

$$\frac{\sigma + 1}{2} \Delta \leq \underline{p} \leq (3 + \sigma)\Delta. \quad (286)$$

Step 3: combining all the results. Finally, summing up the results for both \underline{q} (in (278) and (280)) and \underline{p} (in (283) and (286)), we arrive at the advertised bound.

E.3.2. Proof of Lemma 20

The robust value function for any policy π . For any \mathcal{M}_ϕ with $\phi \in \{0, 1\}$, we first characterize the robust value function of any policy π over different states.

Towards this, it is easily observed that for any policy π , the robust value functions at state $s = 1$ or any $s \in \{2, 3, \dots, S - 1\}$ obey

$$V_\phi^{\pi, \sigma}(1) \stackrel{(i)}{=} 1 + \gamma V_\phi^{\pi, \sigma}(1) = \frac{1}{1 - \gamma} \quad (287a)$$

and

$$\forall s \in \{2, 3, \dots, S - 1\}: \quad V_\phi^{\pi, \sigma}(s) \stackrel{(ii)}{=} 0 + \gamma V_\phi^{\pi, \sigma}(1) = \frac{\gamma}{1 - \gamma}, \quad (287b)$$

where (i) and (ii) are according to the facts that the transitions defined over states $s \geq 1$ in (56) give only one possible next state 1, leading to a non-random transition in the uncertainty set associated with χ^2 divergence, and $r(1, a) = 1$ for all $a \in \mathcal{A}$ and $r(s, a) = 0$ holds for all $(s, a) \in \{2, 3, \dots, S-1\} \times \mathcal{A}$.

To continue, the robust value function at state 0 with policy π satisfies

$$\begin{aligned} V_\phi^{\pi, \sigma}(0) &= \mathbb{E}_{a \sim \pi(\cdot | 0)} \left[r(0, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,a}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} \right] \\ &= 0 + \gamma \pi(\phi | 0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,\phi}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} + \gamma \pi(1 - \phi | 0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,1-\phi}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} \end{aligned} \quad (288)$$

$$\stackrel{(i)}{\leq} \frac{\gamma}{1 - \gamma}, \quad (289)$$

where (i) holds by that $\|V_\phi^{\pi, \sigma}\|_\infty \leq \frac{1}{1-\gamma}$. Summing up the results in (287b) and (289) leads to

$$\forall s \in \{2, 3, \dots, S-1\}, \quad V_\phi^{\pi, \sigma}(1) > V_\phi^{\pi, \sigma}(s) \geq V_\phi^{\pi, \sigma}(0). \quad (290)$$

With the transition kernel in (56) over state 0 and the fact in (290), (288) can be rewritten as

$$\begin{aligned} V_\phi^{\pi, \sigma}(0) &= \gamma \pi(\phi | 0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,\phi}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} + \gamma \pi(1 - \phi | 0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,1-\phi}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} \\ &\stackrel{(i)}{=} \gamma \pi(\phi | 0) \left[\underline{p} V_\phi^{\pi, \sigma}(1) + (1 - \underline{p}) V_\phi^{\pi, \sigma}(0) \right] + \gamma \pi(1 - \phi | 0) \left[\underline{q} V_\phi^{\pi, \sigma}(1) + (1 - \underline{q}) V_\phi^{\pi, \sigma}(0) \right] \\ &\stackrel{(ii)}{=} \gamma z_\phi^\pi V_\phi^{\pi, \sigma}(1) + \gamma (1 - z_\phi^\pi) V_\phi^{\pi, \sigma}(0) \\ &= \frac{\gamma z_\phi^\pi}{(1 - \gamma)(1 - \gamma(1 - z_\phi^\pi))}, \end{aligned} \quad (291)$$

where (i) holds by the definition of \underline{p} and \underline{q} in (255), (ii) follows from the definition of z_ϕ^π in (258), and the last line holds by applying (287a) and solving the resulting linear equation for $V_\phi^{\pi, \sigma}(0)$.

Optimal policy and its optimal value function. To continue, observing that $V_\phi^{\pi, \sigma}(0) =: f(z_\phi^\pi)$ is increasing in z_ϕ^π since the derivative of $f(z_\phi^\pi)$ w.r.t. z_ϕ^π obeys

$$f'(z_\phi^\pi) = \frac{\gamma(1 - \gamma)(1 - \gamma(1 - z_\phi^\pi)) - \gamma^2 z_\phi^\pi(1 - \gamma)}{(1 - \gamma)^2 (1 - \gamma(1 - z_\phi^\pi))^2} = \frac{\gamma}{(1 - \gamma(1 - z_\phi^\pi))^2} > 0,$$

where the last inequality holds by $0 \leq z_\phi^\pi \leq 1$. Further, z_ϕ^π is also increasing in $\pi(\phi | 0)$ (see the fact $\underline{p} \geq \underline{q}$ from Lemma 19), the optimal robust policy in state 0 thus obeys

$$\pi_\phi^*(\phi | 0) = 1. \quad (292)$$

Considering that the action does not influence the state transition for all states $s > 0$, without loss of generality, we choose the optimal robust policy to obey

$$\forall s > 0: \quad \pi_\phi^*(\phi | s) = 1. \quad (293)$$

Taking $\pi = \pi_\phi^*$ and $z_\phi^{\pi_\phi^*} = \underline{p}$ in (291), we complete the proof by showing the corresponding optimal robust value function at state 0 as follows:

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

E.3.3. Proof of the claim (263) Plugging in the definition of φ , we arrive at that for any policy π ,

$$\begin{aligned} \langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\pi,\sigma} \rangle &= V_\phi^{*,\sigma}(0) - V_\phi^{\pi,\sigma}(0) \\ &\stackrel{(i)}{=} \frac{\gamma \underline{p}}{(1-\gamma)(1-\gamma(1-\underline{p}))} - \frac{\gamma z_\phi^\pi}{(1-\gamma)(1-\gamma(1-z_\phi^\pi))} \\ &= \frac{\gamma(\underline{p} - z_\phi^\pi)}{(1-\gamma(1-\underline{p}))(1-\gamma(1-z_\phi^\pi))} \stackrel{(ii)}{\geq} \frac{\gamma(\underline{p} - z_\phi^\pi)}{(1-\gamma(1-\underline{p}))^2} \stackrel{(iii)}{=} \frac{\gamma(\underline{p} - \underline{q})(1 - \pi(\phi|0))}{(1-\gamma(1-\underline{p}))^2}, \end{aligned} \quad (294)$$

where (i) holds by applying Lemma 20, (ii) arises from $z_\phi^\pi \leq \underline{p}$ (see the definition of z_ϕ^π in (258) and the fact $\underline{p} \geq \underline{q} + \frac{3\Delta}{4}$ from Lemma 19), and (iii) follows from the definition of z_ϕ^π in (258).

To further control (294), we consider it in two cases separately:

- Case 1: $\sigma \in (0, \frac{1-\gamma}{4})$. In this case, applying Lemma 19 to (294) yields

$$\begin{aligned} \langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\pi,\sigma} \rangle &\geq \frac{\gamma(\underline{p} - \underline{q})(1 - \pi(\phi|0))}{(1-\gamma(1-\underline{p}))^2} \geq \frac{\gamma \frac{3\Delta}{4}(1 - \pi(\phi|0))}{\left(1-\gamma\left(1-\frac{5(1-\gamma)}{4}\right)\right)^2} \\ &\geq \frac{\Delta(1 - \pi(\phi|0))}{9(1-\gamma)^2} = 2\varepsilon(1 - \pi(\phi|0)), \end{aligned} \quad (295)$$

where the penultimate inequality follows from $\gamma \geq 3/4$, and the last inequality holds by taking the specification of Δ in (262) as follows:

$$\Delta = 18(1-\gamma)^2\varepsilon. \quad (296)$$

It is easily verified that taking $\varepsilon \leq \frac{1}{72(1-\gamma)}$ as in (261) directly leads to meeting the requirement in (253), i.e., $\Delta \leq \frac{1}{4}(1-\gamma)$.

- Case 2: $\sigma \in [\frac{1-\gamma}{4}, \infty)$. Similarly, applying Lemma 19 to (294) gives

$$\langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\pi,\sigma} \rangle \geq \frac{\gamma(\underline{p} - \underline{q})(1 - \pi(\phi|0))}{(1-\gamma(1-\underline{p}))^2} \geq \frac{\gamma \frac{\sigma+1}{2} \Delta(1 - \pi(\phi|0))}{\min\left\{1, (1-\gamma(1-(3+\sigma)\Delta))^2\right\}} \quad (297)$$

Before continuing, it can be verified that

$$1 - \gamma(1 - (3 + \sigma)\Delta) = 1 - \gamma + \gamma(3 + \sigma)\Delta \stackrel{(i)}{\leq} 2(1 - \gamma), \quad (298)$$

where (i) is obtained by $\Delta \leq \frac{1-\gamma}{\gamma(3+\sigma)}$ (see (253)). Applying the above fact to (297) gives

$$\begin{aligned} \langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\pi,\sigma} \rangle &\geq \frac{\gamma^{\frac{\sigma+1}{2}} \Delta (1 - \pi(\phi|0))}{\min \left\{ 1, (1 - \gamma(1 - (3 + \sigma)\Delta))^2 \right\}} \stackrel{(i)}{\geq} \frac{3(\sigma + 1)\Delta (1 - \pi(\phi|0))}{32(1 - \gamma)^2} \\ &= 2\varepsilon(1 - \pi(\phi|0)), \end{aligned} \quad (299)$$

where (i) holds by $\gamma \geq \frac{3}{4}$ and (297), and the last equality holds by the specification in (262):

$$\Delta = \frac{64\varepsilon(1 - \gamma)^2}{3(1 + \sigma)}. \quad (300)$$

As a result, it is easily verified that the requirement in (253)

$$\Delta \leq \frac{1 - \gamma}{\gamma(3 + \sigma)} \quad (301)$$

is met if we let

$$\varepsilon \leq \frac{1}{768(1 - \gamma)} \quad (302)$$

as in (261).

The proof is then completed by summing up the results in the above two cases.

F. Proof for the offline setting

F.1. Proof of the upper bounds: Corollary 1 and Corollary 3

As the proofs of Corollary 1 and Corollary 3 are similar, without loss of generality, we first focus on Corollary 1 in the case of TV distance.

To begin with, suppose we have access to in total N_b independent sample tuples $\{s_i, a_i, s'_i, r_i\}_{i=1}^{N_b}$ from either the generative model or a historical dataset. We denote the number of samples generated based on the state-action pair (s, a) as $N(s, a)$, i.e.,

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad N(s, a) = \sum_{i=1}^{N_b} \mathbb{1}\{s_i = s, a_i = a\}. \quad (303)$$

Then according to (13), we can construct an empirical nominal transition for DRVI (Algorithm 1).

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \hat{P}^0(s' | s, a) := \frac{1}{N(s, a)} \sum_{i=1}^{N(s, a)} \mathbb{1}\{s_i = s, a_i = a, s'_i = s'\}. \quad (304)$$

Armed with the above estimate of nominal transition kernel, we introduce a slightly more general version of Theorem 1, which follows directly from the same proof routine in Appendix B.2.

THEOREM 5 (Upper bound under TV distance). *Let the uncertainty set be $\mathcal{U}_p^\sigma(\cdot) = \mathcal{U}_{TV}^\sigma(\cdot)$, as specified by the TV distance (9). Consider any discount factor $\gamma \in [\frac{1}{4}, 1)$, uncertainty level $\sigma \in (0, 1)$, and $\delta \in (0, 1)$. Based on the empirical nominal transition kernel in (304), let $\hat{\pi}$ be the output policy of Algorithm 1 after $T = C_1 \log\left(\frac{N_b}{1-\gamma}\right)$ iterations. Then with probability at least $1 - \delta$, one has*

$$\forall s \in \mathcal{S}: \quad V^{*,\sigma}(s) - V^{\hat{\pi},\sigma}(s) \leq \varepsilon \quad (305)$$

for any $\varepsilon \in \left(0, \sqrt{1/\max\{1-\gamma, \sigma\}}\right]$, as long as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad N(s, a) \geq \frac{C_2}{(1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2} \log\left(\frac{SAN_b}{(1-\gamma)\delta}\right). \quad (306)$$

Here, $C_1, C_2 > 0$ are some large enough universal constants.

Furthermore, we invoke a fact derived from basic concentration inequalities (Li et al. 2024a) as below.

LEMMA 22. *Consider any $\delta \in (0, 1)$ and a dataset with N_b independent samples satisfying Assumption 1. With probability at least $1 - \delta$, the quantities $\{N(s, a)\}$ obey*

$$\max\left\{N(s, a), \frac{2}{3} \log \frac{N_b}{\delta}\right\} \geq \frac{N_b \mu^b(s, a)}{12} \quad (307)$$

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

Now we are ready to verify Corollary 1. Armed with a historical dataset \mathcal{D}^b with N_b independent samples that obeys Assumption 1, one has with probability at least $1 - \delta$,

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad N(s, a) \geq \frac{N_b \mu^b(s, a)}{12} \geq \frac{N_b \mu_{\min}}{12} \quad (308)$$

as long as $N_b \geq \frac{8 \log \frac{N_b}{\delta}}{\mu_{\min}} \geq \frac{8 \log \frac{N_b}{\delta}}{\mu^b(s, a)}$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. Consequently, given $N_b \geq \frac{8 \log \frac{N_b}{\delta}}{\mu_{\min}}$, applying Theorem 5 with the fact $N(s, a) \geq \frac{N_b \mu_{\min}}{12}$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ (see (308)) directly leads to: DRVI can achieve an ε -optimal policy as long as

$$N(s, a) \geq \frac{N_b \mu_{\min}}{12} \geq \frac{C_2}{(1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2} \log\left(\frac{SAN_b}{(1-\gamma)\delta}\right), \quad (309)$$

namely

$$N_b \geq \frac{C_3}{\mu_{\min}(1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2} \log\left(\frac{SAN_b}{(1-\gamma)\delta}\right), \quad (310)$$

where C_3 is some large enough universal constant. Note that the above inequality directly implies $N_b \geq \frac{8 \log \frac{N_b}{\delta}}{\mu_{\min}}$. This completes the proof of Corollary 1. The same argument holds for Corollary 3.

F.2. Proof of the lower bounds: Corollary 2 and Corollary 4

Analogous to Appendix F.1, without loss of generality, we first focus on verifying Corollary 2 where we use the TV distance to measure the uncertainty set.

We stick to the two hard instances \mathcal{M}_0 and \mathcal{M}_1 (i.e., \mathcal{M}_ϕ with $\phi \in \{0, 1\}$) constructed in the proof for Theorem 2 (Appendix C.1). Recall that the state space is defined as $\mathcal{S} = \{0, 1, 2, \dots, S-1\}$, where the corresponding action space for any state $s \in \{2, 3, \dots, S-1\}$ is $\mathcal{A} = \{0, 1, 2, \dots, A-1\}$. For states $s = 0$ or $s = 1$, the action space is only $\mathcal{A}' = \{0, 1\}$. Hence, for a given factor $\mu_{\min} \in (0, \frac{1}{SA}]$, we can construct a historical dataset \mathcal{D}^b with N_b samples such that the data coverage becomes the smallest over the state-action pairs $(0, 0)$ and $(0, 1)$, i.e.,

$$\mu^b(0, 0) = \mu^b(0, 1) = \mu_{\min} \quad \text{and} \quad \mu^b(s, a) = \frac{1 - 2\mu_{\min}}{(S-2)A + 2}, \quad \forall s \in \{1, 2, \dots, S-1\}. \quad (311)$$

Armed with the above hard instance and historical dataset, we follow the proof procedure in Appendix C.2 to verify the corollary. Our goal is to distinguish between the two hypotheses $\phi \in \{0, 1\}$ by considering the minimax probability of error as follows:

$$p_e := \inf_{\psi} \max \{ \mathbb{P}_0(\psi \neq 0), \mathbb{P}_1(\psi \neq 1) \}, \quad (312)$$

where the infimum is taken over all possible tests ψ constructed from the samples in \mathcal{D}^b .

Recall that we denote μ_ϕ (resp. $\mu_\phi(s)$) as the distribution of a sample tuple (s_i, a_i, s'_i) under the nominal transition kernel P^ϕ associated with \mathcal{M}_ϕ and the samples are generated independently. Analogous to (172), one has

$$\begin{aligned} p_e &\geq \frac{1}{4} \exp \left(-N_b \text{KL}(\mu_0 \parallel \mu_1) \right) \\ &= \frac{1}{4} \exp \left\{ -N_b \mu_{\min} \left(\text{KL}(P^0(\cdot | 0, 0) \parallel P^1(\cdot | 0, 0)) + \text{KL}(P^0(\cdot | 0, 1) \parallel P^1(\cdot | 0, 1)) \right) \right\}, \end{aligned} \quad (313)$$

where the last inequality holds by observing that

$$\begin{aligned} \text{KL}(\mu_0 \parallel \mu_1) &= \sum_{s, a, s'} \mu^b(s, a) \text{KL}(P^0(s' | s, a) \parallel P^1(s' | s, a)) \\ &= \sum_{a \in \{0, 1\}} \mu^b(0, a) \text{KL}(P^0(\cdot | 0, a) \parallel P^1(\cdot | 0, a)) = \mu_{\min} \sum_{a \in \{0, 1\}} \text{KL}(P^0(\cdot | 0, a) \parallel P^1(\cdot | 0, a)). \end{aligned} \quad (314)$$

Here, the last line holds by the fact that $P^0(\cdot | s, a)$ and $P^1(\cdot | s, a)$ (associated with \mathcal{M}_0 and \mathcal{M}_1) only differ from each other in state-action pairs $(0, 0)$ and $(0, 1)$, each has a visitation density of μ_{\min} . Consequently, following the same routine from (173) to the end of Appendix C.2, we apply (174) and (175) with $N = N_b \mu_{\min}$ and complete the proof by showing: if the sample size is selected as

$$N_b \mu_{\min} = N \leq \frac{c_1 \log 2}{8192(1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2}, \quad (315)$$

then one necessarily has

$$p_\varepsilon = \inf_{\hat{\pi}} \max \left\{ \mathbb{P}_0 \left(V^{*,\sigma}(\varphi) - V^{\hat{\pi},\sigma}(\varphi) > \varepsilon \right), \mathbb{P}_1 \left(V^{*,\sigma}(\varphi) - V^{\hat{\pi},\sigma}(\varphi) > \varepsilon \right) \right\} \geq \frac{1}{8}. \quad (316)$$

We can follow the same argument to complete the proof of Corollary 4.

References

- Agarwal, A., Kakade, S., and Yang, L. F. (2020). Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pages 67–83. PMLR.
- Iyengar, G. N. (2005). Robust dynamic programming. *Mathematics of Operations Research*, 30(2):257–280.
- Li, G., Shi, L., Chen, Y., Chi, Y., and Wei, Y. (2024a). Settling the sample complexity of model-based offline reinforcement learning. *The Annals of Statistics*, 52(1):233–260.
- Li, G., Wei, Y., Chi, Y., and Chen, Y. (2024b). Breaking the sample size barrier in model-based reinforcement learning with a generative model. *Operations Research*, 72(1):203–221.
- Panaganti, K. and Kalathil, D. (2022). Sample complexity of robust reinforcement learning with a generative model. In *International Conference on Artificial Intelligence and Statistics*, pages 9582–9602. PMLR.
- Shi, L. and Chi, Y. (2024). Distributionally robust model-based offline reinforcement learning with near-optimal sample complexity. *Journal of Machine Learning Research*, 25(200):1–91.
- Tsybakov, A. B. (2009). *Introduction to nonparametric estimation*, volume 11. Springer.
- Vershynin, R. (2018). *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press.
- Yang, W., Zhang, L., and Zhang, Z. (2022). Toward theoretical understandings of robust Markov decision processes: Sample complexity and asymptotics. *The Annals of Statistics*, 50(6):3223–3248.