

495 **Appendix****Table of Contents**

498	A Other related works	14
499	B Preliminaries	15
500	B.1 Basic facts	16
501	B.2 Properties of the robust Bellman operator	17
502	B.3 Additional facts of the empirical robust MDP	18
503	C Proof of the upper bound with TV distance: Theorem 1	20
504	C.1 Technical lemmas	20
505	C.2 Proof of Theorem 1	21
506	C.3 Proof of the auxiliary lemmas	30
507	D Proof of the lower bound with TV distance: Theorem 2	39
508	D.1 Construction of the hard problem instances	39
509	D.2 Establishing the minimax lower bound	40
510	D.3 Proof of the auxiliary facts	42
511	E Proof of the upper bound with χ^2 divergence: Theorem 3	46
512	E.1 Proof of Theorem 3	46
513	E.2 Proof of the auxiliary lemmas	47
514	F Proof of the lower bound with χ^2 divergence: Theorem 4	51
515	F.1 Construction of the hard problem instances	51
516	F.2 Establishing the minimax lower bound	53
517	F.3 Proof of the auxiliary facts	55

521 **A Other related works**

522 We limit our discussions primarily to provable RL algorithms in the tabular setting with finite state
523 and action spaces, which are most related to our work.

524 **Finite-sample guarantees for standard RL.** A surge of recent research has utilized the toolkit of
525 concentration inequalities to investigate the performance of standard RL algorithms in non-asymptotic
526 settings. There has been a considerable amount of research into non-asymptotic sample analysis of
527 standard RL for a variety of settings; a small set of samples include, but are not limited to, the works
528 via probably approximately correct (PAC) bounds for the generative model setting (Kearns and Singh,
529 1999; Beck and Srikant, 2012; ?; Chen et al., 2020; Azar et al., 2013; Sidford et al., 2018; Agarwal
530 et al., 2020; Li et al., 2023, 2020; Wainwright, 2019), the offline setting (Rashidinejad et al., 2021;
531 Xie et al., 2021; Yin et al., 2021; Shi et al., 2022; Li et al., 2022a; Jin et al., 2021; Yan et al., 2022),
532 and the online setting via regret analysis (Jin et al., 2018; Bai et al., 2019; Li et al., 2021; Zhang et al.,
533 2020b; Dong et al., 2019; Jafarnia-Jahromi et al., 2020; Yang et al., 2021).

534 **Robustness in RL.** Although standard RL has achieved remarkable success, current RL algorithms
535 are still limited since the agent may fail catastrophically if the deployed environment is subject to
536 perturbation, uncertainty, and even structural changes. To address these challenges, an emerging line
537 of works begin to address robustness of RL algorithms with respect to the uncertainty or perturbation
538 over different components of MDPs — state, action, reward, and the transition kernel; see Moos
539 et al. (2022) for a recent review. Besides the framework of distributionally robust MDPs (RMDPs)

(Iyengar, 2005) adopted by this work, to promote robustness in RL, there exist various other works including but not limited to Zhang et al. (2020a, 2021); Han et al. (2022); Qiaoben et al. (2021); Sun et al. (2021); Xiong et al. (2022) investigating the robustness w.r.t. state uncertainty, where the agent’s policy is chosen based on a perturbed observation generated from the state by adding restricted noise or adversarial attack. Besides, Tessler et al. (2019); Tan et al. (2020) considered the robustness to the uncertainty of the action, namely, the action is possibly distorted by an adversarial agent abruptly or smoothly.

Distributionally robust RL. Rooted in the literature of distributionally robust optimization, which has primarily been investigated in the context of supervised learning (Rahimian and Mehrotra, 2019; Gao, 2020; Bertsimas et al., 2018; Duchi and Namkoong, 2018; Blanchet and Murthy, 2019), distributionally robust dynamic programming and RMDPs have attracted considerable attention recently (Iyengar, 2005; Xu and Mannor, 2012; Wolff et al., 2012; Kaufman and Schaefer, 2013; Ho et al., 2018; Smirnova et al., 2019; Ho et al., 2021; Goyal and Grand-Clement, 2022; Derman and Mannor, 2020; Tamar et al., 2014; Badrinath and Kalathil, 2021). In the context of RMDPs, both empirical and theoretical studies have been widely conducted, although most prior theoretical analyses focus on planning with an exact knowledge of the uncertainty set (Iyengar, 2005; Xu and Mannor, 2012; Tamar et al., 2014), or are asymptotic in nature (Roy et al., 2017).

Resorting to the tools of high-dimensional statistics, various recent works begin to shift attention to understand the finite-sample performance of provable robust RL algorithms, under diverse data generating mechanisms and forms of the uncertainty set over the transition kernel. Besides the infinite-horizon setting, finite-sample complexity bounds for RMDPs with the TV distance and the χ^2 divergence are also developed for the finite-horizon setting in Xu et al. (2023); Dong et al. (2022). In addition, many other forms of uncertainty sets have been considered. For example, Wang and Zou (2021) considered a R-contamination uncertain set and proposed a provable robust Q-learning algorithm for the online setting with similar guarantees as standard MDPs. The KL divergence is another popular choice widely considered, where Yang et al. (2022); Panaganti and Kalathil (2022); Zhou et al. (2021); Shi and Chi (2022); Xu et al. (2023); Wang et al. (2023); ? investigated the sample complexity of both model-based and model-free algorithms under the simulator or offline settings. Xu et al. (2023) considered a variety of uncertainty sets including one associated with Wasserstein distance. Badrinath and Kalathil (2021) considered a general (s, a) -rectangular form of the uncertainty set and proposed a model-free algorithm for the online setting with linear function approximation to cope with large state spaces. Moreover, various other related issues have been explored such as the iteration complexity of the policy-based methods (Li et al., 2022b; Kumar et al., 2023), and regularization-based robust RL (Yang et al., 2023).

574 B Preliminaries

575 For convenience, we introduce the notation $[T] := \{1, \dots, T\}$ for any positive integer $T > 0$.
 576 Moreover, for any two vectors $x = [x_i]_{1 \leq i \leq n}$ and $y = [y_i]_{1 \leq i \leq n}$, the notation $x \leq y$ (resp. $x \geq y$)
 577 means $x_i \leq y_i$ (resp. $x_i \geq y_i$) for all $1 \leq i \leq n$. And for any vector x , we overload the notation
 578 by letting $x^{\circ 2} = [x(s, a)^2]_{(s, a) \in \mathcal{S} \times \mathcal{A}}$ (resp. $x^{\circ 2} = [x(s)^2]_{s \in \mathcal{S}}$). With slight abuse of notation, we
 579 denote 0 (resp. 1) as the all-zero (resp. all-one) vector, and drop the subscript ρ to write $\mathcal{U}^\sigma(\cdot) = \mathcal{U}_\rho^\sigma(\cdot)$
 580 whenever the argument holds for all divergence ρ .

581 **Matrix notation.** To continue, we recall or introduce some additional matrix notation that is useful
 582 throughout the analysis.

- 583 • $P^0 \in \mathbb{R}^{SA \times S}$: the matrix of the nominal transition kernel with $P_{s,a}^0$ as the (s, a) -th row.
- 584 • $\hat{P}^0 \in \mathbb{R}^{SA \times S}$: the matrix of the estimated nominal transition kernel with $\hat{P}_{s,a}^0$ as the
 585 (s, a) -th row.
- 586 • $r \in \mathbb{R}^{SA}$: a vector representing the reward function r (so that $r_{(s,a)} = r(s, a)$ for all
 587 $(s, a) \in \mathcal{S} \times \mathcal{A}$).

588
589

- $\Pi^\pi \in \{0, 1\}^{S \times SA}$: a projection matrix associated with a given deterministic policy π taking the following form

$$\Pi^\pi = \begin{pmatrix} \mathbf{e}_{\pi(1)}^\top & 0^\top & \cdots & 0^\top \\ 0^\top & \mathbf{e}_{\pi(2)}^\top & \cdots & 0^\top \\ \vdots & \vdots & \ddots & \vdots \\ 0^\top & 0^\top & \cdots & \mathbf{e}_{\pi(S)}^\top \end{pmatrix}, \quad (18)$$

590

where $\mathbf{e}_{\pi(1)}^\top, \mathbf{e}_{\pi(2)}^\top, \dots, \mathbf{e}_{\pi(S)}^\top \in \mathbb{R}^A$ are standard basis vectors.

591
592
593
594

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

$$\text{Var}_P(s, a) := \text{Var}_{P_{s,a}}(V). \quad (19)$$

595
596
597
598

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

$$P_{s,a}^V = \operatorname{argmin}_{P \in \mathcal{U}^\sigma(P_{s,a}^0)} PV, \quad \text{and} \quad \hat{P}_{s,a}^V = \operatorname{argmin}_{P \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} PV. \quad (20a)$$

599

Furthermore, we make use of the following short-hand notation:

$$\begin{aligned} P_{s,a}^{\pi,V} &:= P_{s,a}^{V,\pi,\sigma} = \operatorname{argmin}_{P \in \mathcal{U}^\sigma(P_{s,a}^0)} PV^{\pi,\sigma}, \\ P_{s,a}^{\pi,\hat{V}} &:= P_{s,a}^{\hat{V},\pi,\sigma} = \operatorname{argmin}_{P \in \mathcal{U}^\sigma(P_{s,a}^0)} P\hat{V}^{\pi,\sigma}, \end{aligned} \quad (20b)$$

$$\begin{aligned} \hat{P}_{s,a}^{\pi,V} &:= \hat{P}_{s,a}^{V,\pi,\sigma} = \operatorname{argmin}_{P \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} PV^{\pi,\sigma}, \\ \hat{P}_{s,a}^{\pi,\hat{V}} &:= \hat{P}_{s,a}^{\hat{V},\pi,\sigma} = \operatorname{argmin}_{P \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} P\hat{V}^{\pi,\sigma}. \end{aligned} \quad (20c)$$

600
601

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

602
603

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

$$\begin{aligned} P^\pi &:= \Pi^\pi P^0, & \hat{P}^\pi &:= \Pi^\pi \hat{P}^0, & \underline{P}^{\pi,V} &:= \Pi^\pi P^{\pi,V}, & \underline{P}^{\pi,\hat{V}} &:= \Pi^\pi P^{\pi,\hat{V}}, \\ \hat{\underline{P}}^{\pi,V} &:= \Pi^\pi \hat{P}^{\pi,V}, & \text{and} & & \hat{\underline{P}}^{\pi,\hat{V}} &:= \Pi^\pi \hat{P}^{\pi,\hat{V}}. \end{aligned} \quad (21)$$

604
605

We denote P_s^π as the s -th row of the transition matrix P^π ; similar quantities can be defined for the other matrices as well.

606

B.1 Basic facts

607
608
609

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

$$\text{KL}(p \parallel q) := p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q} \quad \text{and} \quad \chi^2(p \parallel q) := \frac{(p-q)^2}{q} + \frac{(p-q)^2}{1-q} = \frac{(p-q)^2}{q(1-q)}, \quad (22)$$

610
611
612

which represent respectively the KL divergence and the χ^2 divergence of $\text{Ber}(p)$ from $\text{Ber}(q)$ (Tsybakov and Zaiats, 2009). We make note of the following useful property about the KL divergence in Tsybakov and Zaiats (2009, Lemma 2.7).

613

Lemma 1. For any $p, q \in (0, 1)$, it holds that

$$\text{KL}(p \parallel q) \leq \frac{(p-q)^2}{q(1-q)}. \quad (23)$$

614 **Variance.** For any probability vector $P \in \mathbb{R}^{1 \times S}$ and vector $V \in \mathbb{R}^S$, we denote the variance

$$\text{Var}_P(V) := P(V \circ V) - (PV) \circ (PV). \quad (24)$$

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

616 **Lemma 2.** 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
617 $P \in \Delta(S)$, one has

$$|\text{Var}_P(V_1) - \text{Var}_P(V_2)| \leq \frac{2x}{(1-\gamma)}. \quad (25)$$

618 *Proof.* It is immediate to check that

$$\begin{aligned} |\text{Var}_P(V_1) - \text{Var}_P(V_2)| &= |P(V_1 \circ V_1) - (PV_1) \circ (PV_1) - P(V_2 \circ V_2) + (PV_2) \circ (PV_2)| \\ &\leq |P(V_1 \circ V_1 - V_2 \circ V_2)| + |(PV_1 + PV_2)P(V_1 - V_2)| \\ &\leq 2\|V_1 + V_2\|_\infty \|V_1 - V_2\|_\infty \leq \frac{2x}{(1-\gamma)}. \end{aligned} \quad (26)$$

619 where the penultimate inequality holds by the triangle inequality. \square

620 B.2 Properties of the robust Bellman operator

621 **γ -contraction of the robust Bellman operator.** It is worth noting that the robust Bellman operator
622 (cf. (6)) shares the nice γ -contraction property of the standard Bellman operator, stated as below.

623 **Lemma 3** (γ -Contraction). (Iyengar, 2005, Theorem 3.2) For any $\gamma \in [0, 1)$, the robust Bellman
624 operator $\mathcal{T}^\sigma(\cdot)$ (cf. (6)) is a γ -contraction w.r.t. $\|\cdot\|_\infty$. Namely, for any $Q_1, Q_2 \in \mathbb{R}^{S \times \mathcal{A}}$ s.t.
625 $Q_1(s, a), Q_2(s, a) \in [0, \frac{1}{1-\gamma}]$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, one has

$$\|\mathcal{T}^\sigma(Q_1) - \mathcal{T}^\sigma(Q_2)\|_\infty \leq \gamma \|Q_1 - Q_2\|_\infty. \quad (27)$$

626 Additionally, $Q^{*,\sigma}$ is the unique fixed point of $\mathcal{T}^\sigma(\cdot)$ obeying $0 \leq Q^{*,\sigma}(s, a) \leq \frac{1}{1-\gamma}$ for all $(s, a) \in$
627 $\mathcal{S} \times \mathcal{A}$.

628 **Dual equivalence of the robust Bellman operator.** Fortunately, the robust Bellman operator can
629 be evaluated efficiently by resorting to its dual formulation (Iyengar, 2005). In what follows, we shall
630 illustrate this for the two choices of the divergence ρ of interest. Before continuing, for any $V \in \mathbb{R}^S$,
631 we denote $[V]_\alpha$ as its clipped version by some non-negative value α , namely,

$$[V]_\alpha(s) := \begin{cases} \alpha, & \text{if } V(s) > \alpha, \\ V(s), & \text{otherwise.} \end{cases} \quad (28)$$

632 • TV distance, where the uncertainty set is $U_\rho^\sigma(\hat{P}_{s,a}^0) := U_{\rho_{\text{TV}}}^\sigma(\hat{P}_{s,a}^0) := U_{\rho_{\text{TV}}}^\sigma(\hat{P}_{s,a}^0)$ w.r.t. the
633 TV distance $\rho = \rho_{\text{TV}}$ defined in (7). In particular, we have the following lemma due to
634 strong duality, which is a direct consequence of Iyengar (2005, Lemma 4.3).

635 **Lemma 4** (Strong duality for TV). Consider any probability vector $P \in \Delta(S)$, any fixed
636 uncertainty level σ and the uncertainty set $U^\sigma(P) := U_{\rho_{\text{TV}}}^\sigma(P)$. For any vector $V \in \mathbb{R}^S$
637 obeying $V \geq 0$, recalling the definition of $[V]_\alpha$ in (28), one has

$$\inf_{P \in U^\sigma(P)} \mathcal{P}V = \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P[V]_\alpha - \sigma \left(\alpha - \min_{s'} [V]_\alpha(s') \right) \right\}. \quad (29)$$

638 In view of the above lemma, the following dual update rule is equivalent to (13) in DRVI:

$$\begin{aligned} \hat{Q}_t(s, a) &= r(s, a) \\ &+ \gamma \max_{\alpha \in [\min_s \hat{V}_{t-1}(s), \max_s \hat{V}_{t-1}(s)]} \left\{ \hat{P}_{s,a}^0 [\hat{V}_{t-1}]_\alpha - \sigma \left(\alpha - \min_{s'} [\hat{V}_{t-1}]_\alpha(s') \right) \right\}. \end{aligned} \quad (30)$$

639 • χ^2 divergence, where the uncertainty set is $U_\rho^\sigma(\hat{P}_{s,a}^0) := U_{\chi^2}^\sigma(\hat{P}_{s,a}^0) := U_{\rho_{\chi^2}}^\sigma(\hat{P}_{s,a}^0)$ w.r.t. the
640 χ^2 divergence $\rho = \rho_{\chi^2}$ defined in (8). We introduce the following lemma which directly
641 follows from (Iyengar, 2005, Lemma 4.2).

642
643
644

Lemma 5 (Strong duality for χ^2). *Consider any probability vector $P \in \Delta(\mathcal{S})$, any fixed uncertainty level σ and the uncertainty set $\mathcal{U}^\sigma(P) := \mathcal{U}_{\chi^2}^\sigma(P)$. For any vector $V \in \mathbb{R}^{\mathcal{S}}$ obeying $V \geq 0$, one has*

$$\inf_{\mathcal{P} \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P[V]_\alpha - \sqrt{\sigma \text{Var}_P([V]_\alpha)} \right\}, \quad (31)$$

645
646

where $\text{Var}_P(\cdot)$ is defined as (24).

In view of the above lemma, the update rule (13) in DRVI can be equivalently written as:

$$\begin{aligned} \widehat{Q}_t(s, a) &= r(s, a) \\ &+ \gamma \max_{\alpha \in [\min_s \widehat{V}_{t-1}(s), \max_s \widehat{V}_{t-1}(s)]} \left\{ \widehat{P}_{s,a}^0 [\widehat{V}_{t-1}]_\alpha - \sqrt{\sigma \text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}_{t-1}]_\alpha)} \right\}. \end{aligned} \quad (32)$$

647

The proofs of Lemma 4 and Lemma 5 are provided as follows.

648

Proof of Lemma 4. To begin with, applying (Iyengar, 2005, Lemma 4.3), the term of interest obeys

$$\inf_{\mathcal{P} \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\mu \in \mathbb{R}^{\mathcal{S}}, \mu \geq 0} \left\{ P(V - \mu) - \sigma \left(\max_{s'} \{V(s') - \mu(s')\} - \min_{s'} \{V(s') - \mu(s')\} \right) \right\}, \quad (33)$$

649
650

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

$$\mu^*(s) = \begin{cases} V(s) - \alpha, & \text{if } V(s) > \alpha \\ 0, & \text{otherwise.} \end{cases} \quad (34)$$

651

Therefore, (33) can be solved by optimizing α as below (Iyengar, 2005, Lemma 4.3):

$$\inf_{\mathcal{P} \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P[V]_\alpha - \sigma \left(\alpha - \min_{s'} [V]_\alpha(s') \right) \right\}. \quad (35)$$

652

□

653

Proof of Lemma 5. Due to strong duality (Iyengar, 2005, Lemma 4.2), it holds that

$$\inf_{\mathcal{P} \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\mu \in \mathbb{R}^{\mathcal{S}}, \mu \geq 0} \left\{ P(V - \mu) - \sqrt{\sigma \text{Var}_P(V - \mu)} \right\}, \quad (36)$$

654

and the optimal μ^* obeys

$$\mu^*(s) = \begin{cases} V(s) - \alpha, & \text{if } V(s) > \alpha \\ 0, & \text{otherwise.} \end{cases} \quad (37)$$

655
656

for some $\alpha \in [\min_s V(s), \max_s V(s)]$. As a result, solving (36) is equivalent to optimizing the scalar α as below:

$$\inf_{\mathcal{P} \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P[V]_\alpha - \sqrt{\sigma \text{Var}_P([V]_\alpha)} \right\}. \quad (38)$$

657

□

658

B.3 Additional facts of the empirical robust MDP

659

Bellman equations of the empirical robust MDP $\widehat{\mathcal{M}}_{\text{rob}}$. To begin with, recall that the empirical robust MDP $\widehat{\mathcal{M}}_{\text{rob}} = \{\mathcal{S}, \mathcal{A}, \gamma, \mathcal{U}^\sigma(\widehat{P}^0), r\}$ based on the estimated nominal distribution \widehat{P}^0 constructed in (10) and its corresponding robust value function (resp. robust Q-function) $\widehat{V}^{\pi, \sigma}$ (resp. $\widehat{Q}^{\pi, \sigma}$).

660
661

662

Note that $\widehat{Q}^{*, \sigma}$ is the unique fixed point of $\widehat{\mathcal{T}}^\sigma(\cdot)$ (see Lemma 3), the empirical robust Bellman operator constructed using \widehat{P}^0 . Moreover, similar to (??), for $\widehat{\mathcal{M}}_{\text{rob}}$, the Bellman's optimality

663

664 principle gives the following *robust Bellman consistency equation* (resp. *robust Bellman optimality*
665 *equation*):

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \widehat{Q}^{\pi, \sigma}(s, a) = r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s,a}^0)} \mathcal{P} \widehat{V}^{\pi, \sigma}, \quad (39a)$$

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \widehat{Q}^{*, \sigma}(s, a) = r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s,a}^0)} \mathcal{P} \widehat{V}^{*, \sigma}. \quad (39b)$$

666 With these in mind, combined with the matrix notation, for any policy π , we can write the robust
667 Bellman consistency equations as

$$Q^{\pi, \sigma} = r + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P^0)} \mathcal{P} V^{\pi, \sigma} \quad \text{and} \quad \widehat{Q}^{\pi, \sigma} = r + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}^0)} \mathcal{P} \widehat{V}^{\pi, \sigma}, \quad (40)$$

668 which leads to

$$\begin{aligned} V^{\pi, \sigma} &= r_\pi + \gamma \Pi^\pi \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P^0)} \mathcal{P} V^{\pi, \sigma} \stackrel{(i)}{=} r_\pi + \gamma \underline{P}^{\pi, V} V^{\pi, \sigma}, \\ \widehat{V}^{\pi, \sigma} &= r_\pi + \gamma \Pi^\pi \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}^0)} \mathcal{P} \widehat{V}^{\pi, \sigma} \stackrel{(ii)}{=} r_\pi + \gamma \underline{\widehat{P}}^{\pi, \widehat{V}} \widehat{V}^{\pi, \sigma}, \end{aligned} \quad (41)$$

669 where (i) and (ii) holds by the definitions in (18), (20) and (21).

670 Encouragingly, the above property of the robust Bellman operator ensures the fast convergence of
671 DRVI. We collect this consequence in the following lemma.

672 **Lemma 6.** *Let $\widehat{Q}_0 = 0$. The iterates $\{\widehat{Q}_t\}, \{\widehat{V}_t\}$ of DRVI obey*

$$\forall t \geq 0 : \quad \|\widehat{Q}_t - \widehat{Q}^{*, \sigma}\|_\infty \leq \frac{\gamma^t}{1 - \gamma} \quad \text{and} \quad \|\widehat{V}_t - \widehat{V}^{*, \sigma}\|_\infty \leq \frac{\gamma^t}{1 - \gamma}. \quad (42)$$

673 Furthermore, the output policy $\widehat{\pi}$ obeys

$$\|\widehat{V}^{*, \sigma} - \widehat{V}^{\widehat{\pi}, \sigma}\|_\infty \leq \frac{2\gamma \varepsilon_{\text{opt}}}{1 - \gamma}, \quad \text{where} \quad \|\widehat{V}^{*, \sigma} - \widehat{V}_{T-1}\|_\infty =: \varepsilon_{\text{opt}}. \quad (43)$$

674 *Proof of Lemma 6.* Applying the γ -contraction property in Lemma 3 directly yields that for any
675 $t \geq 0$,

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

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

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

677 where the penultimate inequality holds by the maximum operator is 1-Lipschitz. This completes the
678 proof of (42).

679 We now move to establish (43). Note that there exists at least one state $s_0 \in \mathcal{S}$ that is associated with
680 the maximum of the value gap, i.e.,

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

681 Recall $\widehat{\pi}^*$ is the optimal robust policy for the empirical RMDP $\widehat{\mathcal{M}}_{\text{rob}}$. For convenience, we denote
682 $a_1 = \widehat{\pi}^*(s_0)$ and $a_2 = \widehat{\pi}(s_0)$. Then, since $\widehat{\pi}$ is the greedy policy w.r.t. \widehat{Q}_T , one has

$$r(s_0, a_1) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s_0, a_1}^0)} \mathcal{P} \widehat{V}_{T-1} = \widehat{Q}_T(s_0, a_1) \leq \widehat{Q}_T(s_0, a_2) = r(s_0, a_2) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s_0, a_2}^0)} \mathcal{P} \widehat{V}_{T-1}. \quad (44)$$

683 Recalling the notation in (20), the above fact and (43) altogether yield

$$\begin{aligned}
r(s_0, a_1) + \gamma \widehat{P}_{s_0, a_1}^{\widehat{V}_{T-1}} \left(\widehat{V}^{*, \sigma} - \varepsilon_{\text{opt}} \mathbf{1} \right) &\leq r(s_0, a_1) + \gamma \widehat{P}_{s_0, a_1}^{\widehat{V}_{T-1}} \widehat{V}_{T-1} \\
&\leq r(s_0, a_2) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s_0, a_2}^0)} \mathcal{P} \widehat{V}_{T-1} \\
&\stackrel{(i)}{\leq} r(s_0, a_2) + \gamma \widehat{P}_{s_0, a_2}^{\widehat{V}^{\widehat{\pi}, \sigma}} \widehat{V}_{T-1} \\
&\leq r(s_0, a_2) + \gamma \widehat{P}_{s_0, a_2}^{\widehat{V}^{\widehat{\pi}, \sigma}} \left(\widehat{V}^{*, \sigma} + \varepsilon_{\text{opt}} \mathbf{1} \right), \tag{45}
\end{aligned}$$

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

$$\begin{aligned}
&\left\| \widehat{V}^{*, \sigma} - \widehat{V}^{\widehat{\pi}, \sigma} \right\|_\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}_{T-1}} \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}_{T-1}} \widehat{V}^{*, \sigma} \right) \\
&\stackrel{(ii)}{\leq} 2\gamma \varepsilon_{\text{opt}} + \gamma \widehat{P}_{s_0, a_2}^{\widehat{V}^{\widehat{\pi}, \sigma}} \left(\widehat{V}^{*, \sigma} - \widehat{V}^{\widehat{\pi}, \sigma} \right) + \gamma \left(\widehat{P}_{s_0, a_1}^{\widehat{V}_{T-1}} \widehat{V}^{*, \sigma} - \widehat{P}_{s_0, a_1}^{\widehat{V}_{T-1}} \widehat{V}^{*, \sigma} \right) \\
&\leq 2\gamma \varepsilon_{\text{opt}} + \gamma \left\| \widehat{V}^{*, \sigma} - \widehat{V}^{\widehat{\pi}, \sigma} \right\|_\infty, \tag{46}
\end{aligned}$$

where (i) holds by plugging in (45), and (ii) follows from $\inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s_0, a_1}^0)} \mathcal{P} \widehat{V}^{*, \sigma} \leq \mathcal{P} \widehat{V}^{*, \sigma}$ for any $\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s_0, a_1}^0)$. Rearranging (46) leads to

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

685

□

686 C Proof of the upper bound with TV distance: Theorem 1

687 Throughout this section, for any transition kernel P , the uncertainty set is taken as (see (7))

$$\mathcal{U}^\sigma(P) := \mathcal{U}_{\text{TV}}^\sigma(P) = \otimes \mathcal{U}_{\text{TV}}^\sigma(P_{s,a}), \quad \mathcal{U}_{\text{TV}}^\sigma(P_{s,a}) := \left\{ P'_{s,a} \in \Delta(\mathcal{S}) : \frac{1}{2} \|P'_{s,a} - P_{s,a}\|_1 \leq \sigma \right\}. \tag{47}$$

688 C.1 Technical lemmas

689 We begin with a key lemma concerning the dynamic range of the robust value function $V^{\pi, \sigma}$ (cf. (??)),
690 which produces tighter control when σ is large; the proof is deferred to Appendix C.3.1.

691 **Lemma 7.** For any nominal transition kernel $P \in \mathbb{R}^{S \times A \times S}$, any fixed uncertainty level σ , and any
692 policy π , its corresponding robust value function $V^{\pi, \sigma}$ (cf. (??)) satisfies

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

693 Next, we introduce the following lemma, whose proof is postponed in Appendix C.3.2.

694 **Lemma 8.** Consider an MDP with transition kernel matrix P and reward function $0 \leq r \leq 1$. For any
 695 policy π and its associated state transition matrix $P_\pi := \Pi^\pi P$ and value function $0 \leq V^{\pi,P} \leq \frac{1}{1-\gamma}$
 696 (cf. (1)), one has

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

697 C.2 Proof of Theorem 1

698 The main proof idea of Theorem 1 is similar to that of Agarwal et al. (2020) and Li et al. (2020)
 699 while the argument needs essential adjustments in order to adapt to the robustness setting. Before
 700 proceeding, applying Lemma 6 yields that for any $\varepsilon_{\text{opt}} > 0$, as long as $T \geq \log(\frac{1}{(1-\gamma)\varepsilon_{\text{opt}}})$, one has

$$\|\widehat{V}^{*,\sigma} - \widehat{V}^{\widehat{\pi},\sigma}\|_\infty \leq \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}, \quad (48)$$

701 allowing us to justify the more general statement in Remark ???. To control the performance gap
 702 $\|V^{*,\sigma} - V^{\widehat{\pi},\sigma}\|_\infty$, the proof is divided into several key steps.

703 **Step 1: decomposing the error.** Recall the optimal robust policy π^* w.r.t. \mathcal{M}_{rob} and the optimal
 704 robust policy $\widehat{\pi}^*$, the optimal robust value function $\widehat{V}^{*,\sigma}$ (resp. robust value function $\widehat{Q}^{\pi^*,\sigma}$) w.r.t.
 705 $\widehat{\mathcal{M}}_{\text{rob}}$. The term of interest $V^{*,\sigma} - V^{\widehat{\pi},\sigma}$ can be decomposed as

$$\begin{aligned} V^{*,\sigma} - V^{\widehat{\pi},\sigma} &= (V^{\pi^*,\sigma} - \widehat{V}^{\pi^*,\sigma}) + (\widehat{V}^{\pi^*,\sigma} - \widehat{V}^{\widehat{\pi}^*,\sigma}) + (\widehat{V}^{\widehat{\pi}^*,\sigma} - \widehat{V}^{\widehat{\pi},\sigma}) + (\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma}) \\ &\stackrel{(i)}{\leq} (V^{\pi^*,\sigma} - \widehat{V}^{\pi^*,\sigma}) + (\widehat{V}^{\widehat{\pi}^*,\sigma} - \widehat{V}^{\widehat{\pi},\sigma}) + (\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma}) \\ &\stackrel{(ii)}{\leq} (V^{\pi^*,\sigma} - \widehat{V}^{\pi^*,\sigma}) + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} 1 + (\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma}) \end{aligned} \quad (49)$$

706 where (i) holds by $\widehat{V}^{\pi^*,\sigma} - \widehat{V}^{\widehat{\pi}^*,\sigma} \leq 0$ since $\widehat{\pi}^*$ is the robust optimal policy for $\widehat{\mathcal{M}}_{\text{rob}}$, and (ii) comes
 707 from the fact in (48).

708 To control the two important terms in (49), we first consider a more general term $\widehat{V}^{\pi,\sigma} - V^{\pi,\sigma}$ for
 709 any policy π . Towards this, plugging in (41) yields

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

710 where (i) holds by observing

$$\underline{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} \leq \underline{P}^{\pi,V} \widehat{V}^{\pi,\sigma}$$

711 due to the optimality of $\underline{P}^{\pi,\widehat{V}}$ (cf. (20)). Rearranging terms leads to

$$\widehat{V}^{\pi,\sigma} - V^{\pi,\sigma} \leq \gamma (I - \gamma \underline{P}^{\pi,V})^{-1} \left(\widehat{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - \underline{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} \right). \quad (50)$$

712 Similarly, we can also deduce

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

$$\geq \gamma \left(I - \gamma \underline{P}^{\pi, \hat{V}} \right)^{-1} \left(\hat{\underline{P}}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - \underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} \right). \quad (51)$$

713 Combining (50) and (51), we arrive at

$$\begin{aligned} \left\| \hat{V}^{\pi, \sigma} - V^{\pi, \sigma} \right\|_{\infty} &\leq \gamma \max \left\{ \left\| \left(I - \gamma \underline{P}^{\pi, V} \right)^{-1} \left(\hat{\underline{P}}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - \underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} \right) \right\|_{\infty}, \right. \\ &\quad \left. \left\| \left(I - \gamma \underline{P}^{\pi, \hat{V}} \right)^{-1} \left(\hat{\underline{P}}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - \underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} \right) \right\|_{\infty} \right\}. \end{aligned} \quad (52)$$

714 By decomposing the error in a symmetric way, we can similarly obtain

$$\begin{aligned} \left\| \hat{V}^{\pi, \sigma} - V^{\pi, \sigma} \right\|_{\infty} &\leq \gamma \max \left\{ \left\| \left(I - \gamma \hat{\underline{P}}^{\pi, V} \right)^{-1} \left(\hat{\underline{P}}^{\pi, V} V^{\pi, \sigma} - \underline{P}^{\pi, V} V^{\pi, \sigma} \right) \right\|_{\infty}, \right. \\ &\quad \left. \left\| \left(I - \gamma \hat{\underline{P}}^{\pi, \hat{V}} \right)^{-1} \left(\hat{\underline{P}}^{\pi, V} V^{\pi, \sigma} - \underline{P}^{\pi, V} V^{\pi, \sigma} \right) \right\|_{\infty} \right\}. \end{aligned} \quad (53)$$

715 With the above facts in mind, we are ready to control the two terms $\|\hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma}\|_{\infty}$ and
716 $\|\hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma}\|_{\infty}$ in (49) separately. More specifically, taking $\pi = \pi^*$, applying (53) leads to

$$\begin{aligned} \left\| \hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma} \right\|_{\infty} &\leq \gamma \max \left\{ \left\| \left(I - \gamma \hat{\underline{P}}^{\pi^*, V} \right)^{-1} \left(\hat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \right\|_{\infty}, \right. \\ &\quad \left. \left\| \left(I - \gamma \hat{\underline{P}}^{\pi^*, \hat{V}} \right)^{-1} \left(\hat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \right\|_{\infty} \right\}. \end{aligned} \quad (54)$$

717 Similarly, taking $\pi = \hat{\pi}$, applying (52) leads to

$$\begin{aligned} \left\| \hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma} \right\|_{\infty} &\leq \gamma \max \left\{ \left\| \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \left(\hat{\underline{P}}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right) \right\|_{\infty}, \right. \\ &\quad \left. \left\| \left(I - \gamma \underline{P}^{\hat{\pi}, V} \right)^{-1} \left(\hat{\underline{P}}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right) \right\|_{\infty} \right\}. \end{aligned} \quad (55)$$

718 **Step 2: controlling $\|\hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma}\|_{\infty}$: bounding the first term in (54).** To control the two
719 terms in (54), we first introduce the following lemma whose proof is postponed to Appendix C.3.3.

720 **Lemma 9.** Consider any $\delta \in (0, 1)$. Setting $N \geq \log\left(\frac{18SAN}{\delta}\right)$, with probability at least $1 - \delta$, one
721 has

$$\begin{aligned} \left| \hat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right| &\leq 2 \sqrt{\frac{\log\left(\frac{18SAN}{\delta}\right)}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{\pi^*, \sigma})} + \frac{\log\left(\frac{18SAN}{\delta}\right)}{N(1-\gamma)} 1 \\ &\leq 3 \sqrt{\frac{\log\left(\frac{18SAN}{\delta}\right)}{(1-\gamma)^2 N}} 1, \end{aligned} \quad (56)$$

722 where $\text{Var}_{P^{\pi^*}}(V^{\pi^*, \sigma})$ is defined in (19).

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

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

$$\begin{aligned}
&\leq \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \mathbf{1} + 2 \underbrace{\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, V}}(V^{*, \sigma})}}_{=: \mathcal{C}_1} \\
&+ 2 \underbrace{\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \sqrt{|\text{Var}_{\widehat{P}^{\pi^*}}(V^{*, \sigma}) - \text{Var}_{\widehat{P}^{\pi^*, V}}(V^{*, \sigma})|}}_{=: \mathcal{C}_2} \\
&+ 2 \underbrace{\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \left(\sqrt{\text{Var}_{P^{\pi^*}}(V^{*, \sigma})} - \sqrt{\text{Var}_{\widehat{P}^{\pi^*}}(V^{*, \sigma})} \right)}_{=: \mathcal{C}_3}, \tag{57}
\end{aligned}$$

724 where (i) holds by $\left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \geq 0$, (ii) follows from Lemma 9, and the last inequality arise
725 from

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

726 by applying the triangle inequality.

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

$$\left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \mathbf{1} = \left(I + \sum_{t=1}^{\infty} \gamma^t \left(\widehat{P}^{\pi^*, V} \right)^t \right) \mathbf{1} = \frac{1}{1-\gamma} \mathbf{1}. \tag{58}$$

729 Armed with this fact, we shall control the other three terms $\mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3$ in (57) separately.

730 • Consider \mathcal{C}_1 . We first introduce the following lemma, whose proof is postponed to Ap-
731 pendix C.3.4.

732 **Lemma 10.** Consider any $\delta \in (0, 1)$. With probability at least $1 - \delta$, one has

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

733 Applying Lemma 10 and inserting back to (57) leads to

$$\begin{aligned}
\mathcal{C}_1 &= 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{P}^{\pi^*, V} \right)^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, V}}(V^{*, \sigma})} \\
&\leq 8 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} \left(1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right) \mathbf{1}. \tag{59}
\end{aligned}$$

734 • Consider \mathcal{C}_2 . First, denote $V' := V^{*, \sigma} - \min_{s' \in \mathcal{S}} V^{*, \sigma}(s') \mathbf{1}$, by Lemma 7, it follows that

$$0 \leq V' \leq \frac{1}{\gamma \max\{1-\gamma, \sigma\}} \mathbf{1}. \tag{60}$$

735

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

$$\begin{aligned} \left| \text{Var}_{\tilde{P}_{s,a}}(V^{*,\sigma}) - \text{Var}_{P_{s,a}}(V^{*,\sigma}) \right| &= \left| \text{Var}_{\tilde{P}_{s,a}}(V') - \text{Var}_{P_{s,a}}(V') \right| \\ &\leq \|\tilde{P}_{s,a} - P_{s,a}\|_1 \|V'\|_\infty^2 \\ &\leq \frac{2\sigma}{\gamma^2(\max\{1-\gamma, \sigma\})^2} 1 \leq \frac{2}{\gamma^2 \max\{1-\gamma, \sigma\}} 1. \end{aligned} \quad (61)$$

736

Applying the above relation we obtain

$$\begin{aligned} \mathcal{C}_2 &= 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \sqrt{\left| \text{Var}_{\hat{P}^{\pi^*}}(V^{*,\sigma}) - \text{Var}_{\hat{P}^{\pi^*, V}}(V^{*,\sigma}) \right|} \\ &= 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \sqrt{\left| \Pi^{\pi^*} (\text{Var}_{\hat{P}_0}(V^{*,\sigma}) - \text{Var}_{\hat{P}^{\pi^*, V}}(V^{*,\sigma})) \right|} \\ &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \sqrt{\left\| \text{Var}_{\hat{P}_0}(V^{*,\sigma}) - \text{Var}_{\hat{P}^{\pi^*, V}}(V^{*,\sigma}) \right\|_\infty} 1 \\ &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \sqrt{\frac{2}{\gamma^2 \max\{1-\gamma, \sigma\}}} 1 \\ &= 2\sqrt{\frac{2 \log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1, \end{aligned} \quad (62)$$

737

where the last equality uses $\left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} 1 = \frac{1}{1-\gamma}$ (cf. (58)).

738

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

739

Lemma 11. (*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 $1 - \delta$,

740

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

741

Applying Lemma 11 with $\pi = \pi^*$ and $V = V^{*,\sigma}$ leads to

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

742

which can be plugged in (57) to verify

$$\begin{aligned} \mathcal{C}_3 &= 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \left(\sqrt{\text{Var}_{P^{\pi^*}}(V^{*,\sigma})} - \sqrt{\text{Var}_{\hat{P}^{\pi^*}}(V^{*,\sigma})} \right) \\ &\leq \frac{4}{(1-\gamma)} \frac{\log(\frac{SAN}{\delta}) \|V^{*,\sigma}\|_\infty}{N} 1 \leq \frac{4 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1, \end{aligned} \quad (63)$$

743

where the last line uses $\left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} 1 = \frac{1}{1-\gamma}$ (cf. (58)).

744

Finally, inserting the results of \mathcal{C}_1 in (59), \mathcal{C}_2 in (62), \mathcal{C}_3 in (63), and (58) back into (57) gives

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

$$\begin{aligned}
& + 2\sqrt{\frac{2\log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^2\max\{1-\gamma,\sigma\}N}}1 + \frac{4\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N}1 + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)^2}1 \\
& \leq 10\sqrt{\frac{2\log(\frac{18SAN}{\delta})}{\gamma^3(1-\gamma)^2\max\{1-\gamma,\sigma\}N}}\left(1 + \sqrt{\frac{\log(\frac{SAN}{\delta})}{(1-\gamma)^2N}}\right)1 + \frac{5\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N}1 \\
& \leq 160\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2\max\{1-\gamma,\sigma\}N}}1 + \frac{5\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N}1, \tag{64}
\end{aligned}$$

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

746 **Step 3: controlling $\|\widehat{V}^{\pi^*,\sigma} - V^{\pi^*,\sigma}\|_\infty$: bounding the second term in (54).** To proceed, applying
747 Lemma 9 on the second term of the right hand side of (54) leads to

$$\begin{aligned}
& (I - \gamma\widehat{P}^{\pi^*,\widehat{V}})^{-1}(\widehat{P}^{\pi^*,V}V^{\pi^*,\sigma} - \underline{P}^{\pi^*,V}V^{\pi^*,\sigma}) \\
& \leq 2(I - \gamma\widehat{P}^{\pi^*,\widehat{V}})^{-1}\left(\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}}\sqrt{\text{Var}_{P^{\pi^*}}(V^{*,\sigma})} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)}1\right) \\
& \leq \frac{2\log(\frac{18SAN}{\delta})}{N(1-\gamma)}(I - \gamma\widehat{P}^{\pi^*,\widehat{V}})^{-1}1 + 2\underbrace{\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}}(I - \gamma\widehat{P}^{\pi^*,\widehat{V}})^{-1}\sqrt{\text{Var}_{\widehat{P}^{\pi^*,\widehat{V}}}(\widehat{V}^{\pi^*,\sigma})}}_{=:C_4} \\
& \quad + 2\underbrace{\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}}(I - \gamma\widehat{P}^{\pi^*,\widehat{V}})^{-1}\left(\sqrt{\text{Var}_{\widehat{P}^{\pi^*,\widehat{V}}}(V^{\pi^*,\sigma} - \widehat{V}^{\pi^*,\sigma})}\right)}_{=:C_5} \\
& \quad + 2\underbrace{\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}}(I - \gamma\widehat{P}^{\pi^*,\widehat{V}})^{-1}\left(\sqrt{|\text{Var}_{\widehat{P}^{\pi^*}}(V^{*,\sigma}) - \text{Var}_{\widehat{P}^{\pi^*,\widehat{V}}}(V^{*,\sigma})|}\right)}_{=:C_6} \\
& \quad + 2\underbrace{\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}}(I - \gamma\widehat{P}^{\pi^*,\widehat{V}})^{-1}\left(\sqrt{\text{Var}_{P^{\pi^*}}(V^{*,\sigma})} - \sqrt{\text{Var}_{\widehat{P}^{\pi^*}}(V^{*,\sigma})}\right)}_{=:C_7}, \tag{65}
\end{aligned}$$

748 where the last term \widetilde{C}_3 can be controlled the same as C_3 in (63). We now bound the above terms
749 separately.

750 • Applying Lemma 8 with $P = \widehat{P}^{\pi^*,\widehat{V}}$, $\pi = \pi^*$ and taking $V = \widehat{V}^{\pi^*,\sigma}$ which obeys
751 $\widehat{V}^{\pi^*,\sigma} = r_{\pi^*} + \gamma\widehat{P}^{\pi^*,\widehat{V}}\widehat{V}^{\pi^*,\sigma}$, and in view of (58), the term C_4 in (65) can be controlled as
752 follows:

$$\begin{aligned}
C_4 & = 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}}(I - \gamma\widehat{P}^{\pi^*,\widehat{V}})^{-1}\sqrt{\text{Var}_{\widehat{P}^{\pi^*,\widehat{V}}}(\widehat{V}^{\pi^*,\sigma})} \\
& \leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}}\sqrt{\frac{8(\max_s \widehat{V}^{\pi^*,\sigma}(s) - \min_s \widehat{V}^{\pi^*,\sigma}(s))}{\gamma^2(1-\gamma)^2}}1 \\
& \leq 8\sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^3(1-\gamma)^2\max\{1-\gamma,\sigma\}N}}1, \tag{66}
\end{aligned}$$

753 where the last inequality holds by applying Lemma 7.

754

- To continue, considering \mathcal{C}_5 , we directly observe that (in view of (58))

$$\begin{aligned} \mathcal{C}_5 &= 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \widehat{P}^{\pi^*, \widehat{V}} \right)^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, \widehat{V}}} (V^{\pi^*, \sigma} - \widehat{V}^{\pi^*, \sigma})} \\ &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \left\| V^{*, \sigma} - \widehat{V}^{\pi^*, \sigma} \right\|_{\infty} 1. \end{aligned} \quad (67)$$

755

- Then, it is easily verified that \mathcal{C}_6 can be controlled similarly as (62) as follows:

$$\mathcal{C}_6 \leq 2\sqrt{\frac{2\log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1. \quad (68)$$

756

- Similarly, \mathcal{C}_7 can be controlled the same as (63) shown below:

$$\mathcal{C}_7 \leq \frac{4\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1. \quad (69)$$

757

Combining the results in (66), (67), (68), and (69) and inserting back to (65) leads to

$$\begin{aligned} &\left(I - \gamma \widehat{P}^{\pi^*, \widehat{V}} \right)^{-1} \left(\widehat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \leq 8\sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 \\ &\quad + 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \left\| V^{*, \sigma} - \widehat{V}^{\pi^*, \sigma} \right\|_{\infty} 1 + 2\sqrt{\frac{2\log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 + \frac{4\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1 \\ &\leq 80\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 + 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \left\| V^{*, \sigma} - \widehat{V}^{\pi^*, \sigma} \right\|_{\infty} 1 + \frac{4\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1, \end{aligned} \quad (70)$$

758

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

759

Finally, inserting (64) and (70) back to (54) yields

$$\begin{aligned} \left\| \widehat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma} \right\|_{\infty} &\leq \max \left\{ 160\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{5\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}, \right. \\ &\quad \left. 80\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \left\| V^{*, \sigma} - \widehat{V}^{\pi^*, \sigma} \right\|_{\infty} + \frac{4\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} \right\} \\ &\leq 160\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{8\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}, \end{aligned} \quad (71)$$

760

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

761

Step 4: controlling $\|\widehat{V}^{\widehat{\pi}, \sigma} - V^{\widehat{\pi}, \sigma}\|_{\infty}$: bounding the first term in (55). Unlike the earlier term, we now need to deal with the complicated statistical dependency between $\widehat{\pi}$ and the empirical RMDP. 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 C.3.5.

763

764

765

Lemma 12. Consider any $\delta \in (0, 1)$. Taking $N \geq \log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)$, with probability at least $1 - \delta$,

766

one has

$$\begin{aligned} \left| \widehat{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - \underline{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right| &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*, \sigma})} 1 + \frac{8\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} 1 + \frac{2\gamma\epsilon_{\text{opt}}}{1-\gamma} \\ &\leq 10\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2 N}} 1 + \frac{2\gamma\epsilon_{\text{opt}}}{1-\gamma} 1. \end{aligned} \quad (72)$$

767 With Lemma 12 in hand, we have

$$\begin{aligned}
& \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \left(\hat{\underline{P}}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right) \\
& \stackrel{(i)}{\leq} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \left| \hat{\underline{P}}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right| \\
& \leq 2 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \sqrt{\text{Var}_{P^{\hat{\pi}}}(\hat{V}^{\star, \sigma})} \\
& \quad + \left(I - \gamma P_Q^{\hat{\pi}, V^{\hat{\pi}}} \right)^{-1} \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \right) 1 \\
& \stackrel{(ii)}{\leq} \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} + \frac{2\gamma\varepsilon_{\text{opt}}}{(1-\gamma)^2} \right) 1 + \underbrace{2 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})}}_{=: \mathcal{D}_1} \\
& \quad + \underbrace{2 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \sqrt{|\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\star, \sigma}) - \text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})|}}_{=: \mathcal{D}_2} \\
& \quad + \underbrace{2 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \sqrt{|\text{Var}_{P^{\hat{\pi}}}(\hat{V}^{\star, \sigma}) - \text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\star, \sigma})|}}_{=: \mathcal{D}_3}, \tag{73}
\end{aligned}$$

768 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
769 that falls into $\Delta(\mathcal{S})$.

770 The remainder of the proof will focus on controlling the three terms in (73) separately.

771 • For \mathcal{D}_1 , we introduce the following lemma, whose proof is postponed to C.3.6.

772 **Lemma 13.** 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
773 with probability at least $1 - \delta$,

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

774 Applying Lemma 13 and (58) to (73) leads to

$$\begin{aligned}
\mathcal{D}_1 &= 2 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} \\
&\leq 12 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}} N} 1. \tag{74}
\end{aligned}$$

775 • Applying Lemma 2 with $\|\hat{V}^{\star, \sigma} - \hat{V}^{\hat{\pi}, \sigma}\|_{\infty} \leq \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}$ and (58), \mathcal{D}_2 can be controlled as

$$\begin{aligned}
\mathcal{D}_2 &= 2 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \sqrt{|\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\star, \sigma}) - \text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})|} \\
&\leq 4 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \frac{\sqrt{\gamma\varepsilon_{\text{opt}}}}{1-\gamma} \leq 4 \sqrt{\frac{\gamma\varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1. \tag{75}
\end{aligned}$$

776 • \mathcal{D}_3 can be controlled similar to \mathcal{C}_2 in (62) as follows:

$$\mathcal{D}_3 = 2 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \sqrt{|\text{Var}_{P^{\hat{\pi}}}(\hat{V}^{\star, \sigma}) - \text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\star, \sigma})|}$$

$$\begin{aligned}
&\leq 4\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}}\right)^{-1} \sqrt{\frac{1}{\gamma^2 \max\{1-\gamma, \sigma\}}} \\
&\leq 4\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} \quad (76)
\end{aligned}$$

777 Finally, summing up the results in (74), (75), and (76) and inserting them back to (73) yields: taking
778 $N \geq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2}$ and $\varepsilon_{\text{opt}} \leq \frac{1-\gamma}{\gamma}$, with probability at least $1 - \delta$,

$$\begin{aligned}
&\left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}}\right)^{-1} \left(\hat{\underline{P}}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma}\right) \leq \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} + \frac{2\gamma\varepsilon_{\text{opt}}}{(1-\gamma)^2}\right) 1 \\
&\quad + 12\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 + 4\sqrt{\frac{\gamma\varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4N}} 1 \\
&\quad + 4\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 \\
&\leq 16\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 + \frac{14 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} 1, \quad (77)
\end{aligned}$$

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

780 **Step 5: controlling $\|\hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma}\|_{\infty}$: bounding the second term in (55).** Towards this, applying
781 Lemma 12 leads to

$$\begin{aligned}
&\left(I - \gamma \underline{P}^{\hat{\pi}, V}\right)^{-1} \left(\hat{\underline{P}}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma}\right) \leq \left(I - \gamma \underline{P}^{\hat{\pi}, V}\right)^{-1} \left|\hat{\underline{P}}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma}\right| \\
&\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, V}\right)^{-1} \sqrt{\text{Var}_{P^{\hat{\pi}}}(\hat{V}^{*, \sigma})} \\
&\quad + \left(I - \gamma \underline{P}^{\hat{\pi}, 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 \\
&\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 + \underbrace{2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, V}\right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, V}}(V^{\hat{\pi}, \sigma})}}_{=: \mathcal{D}_4} \\
&\quad + \underbrace{2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, V}\right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, V}}(\hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma})}}_{=: \mathcal{D}_5} \\
&\quad + \underbrace{2\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}, V}}(\hat{V}^{*, \sigma}) - \text{Var}_{\underline{P}^{\hat{\pi}, V}}(\hat{V}^{\hat{\pi}, \sigma})\right|}}_{=: \mathcal{D}_6} \\
&\quad + \underbrace{2\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}^{*, \sigma}) - \text{Var}_{\underline{P}^{\hat{\pi}, V}}(\hat{V}^{\hat{\pi}, \sigma})\right|}}_{=: \mathcal{D}_7}. \quad (78)
\end{aligned}$$

782 We shall bound each of the terms separately.

783
784

- Applying Lemma 8 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 similar to (66) as follows:

$$\mathcal{D}_4 \leq 8 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1. \quad (79)$$

785

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

$$\begin{aligned} \mathcal{D}_5 &= 2 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\hat{\pi}, V} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, V}}(\widehat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma})} \\ &\leq 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\hat{\pi}, \sigma} - \widehat{V}^{\hat{\pi}, \sigma} \right\|_{\infty} 1. \end{aligned} \quad (80)$$

786
787

- Next, observing that \mathcal{D}_6 and \mathcal{D}_7 are almost the same as the terms \mathcal{D}_2 (controlled in (75)) and \mathcal{D}_3 (controlled in (76)) in (73), it is easily verified that they can be controlled as follows

$$\mathcal{D}_6 \leq 4 \sqrt{\frac{\gamma \varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1, \quad \mathcal{D}_7 \leq 4 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1. \quad (81)$$

788 Then inserting the results in (79), (80), and (81) back to (78) leads to

$$\begin{aligned} &\left(I - \gamma \underline{P}^{\hat{\pi}, V} \right)^{-1} \left(\widehat{\underline{P}}^{\hat{\pi}, \widehat{V}} \widehat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \widehat{V}} \widehat{V}^{\hat{\pi}, \sigma} \right) \\ &\leq \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} + \frac{2\gamma \varepsilon_{\text{opt}}}{(1-\gamma)^2} \right) 1 + 8 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 \\ &\quad + 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\hat{\pi}, \sigma} - \widehat{V}^{\hat{\pi}, \sigma} \right\|_{\infty} 1 + 4 \sqrt{\frac{\gamma \varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1 \\ &\quad + 4 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 \\ &\leq 12 \sqrt{\frac{2 \log(\frac{8SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 + \frac{14 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} 1 + 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\hat{\pi}, \sigma} - \widehat{V}^{\hat{\pi}, \sigma} \right\|_{\infty} 1, \end{aligned} \quad (82)$$

789 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
790 letting $N \geq \frac{\log(\frac{54SAN^2}{\delta})}{1-\gamma}$.

791 Finally, inserting (77) and (82) back to (55) yields: taking $\varepsilon_{\text{opt}} \leq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}$ and $N \geq$
792 $\frac{16 \log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2}$, with probability at least $1 - \delta$, one has

$$\begin{aligned} \left\| \widehat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma} \right\|_{\infty} &\leq \max \left\{ 16 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{14 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2}, \right. \\ &\quad \left. 12 \sqrt{\frac{2 \log(\frac{8SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{14 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} + 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\hat{\pi}, \sigma} - \widehat{V}^{\hat{\pi}, \sigma} \right\|_{\infty} \right\} \\ &\leq 24 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{28 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2}. \end{aligned} \quad (83)$$

793 **Step 6: summing up the results.** Summing up the results in (71) and (83) and inserting back
 794 to (49) complete the proof as follows: taking $\varepsilon_{\text{opt}} \leq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}$ and $N \geq \frac{16 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2}$, with
 795 probability at least $1 - \delta$,

$$\begin{aligned}
 \|V^{*,\sigma} - V^{\hat{\pi},\sigma}\|_{\infty} &\leq \|V^{\pi^*,\sigma} - \hat{V}^{\pi^*,\sigma}\|_{\infty} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} + \|\hat{V}^{\hat{\pi},\sigma} - V^{\hat{\pi},\sigma}\|_{\infty} \\
 &\leq \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} + 160\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma,\sigma\}N}} + \frac{8 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} \\
 &\quad + 24\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma,\sigma\}N}} + \frac{28 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} \\
 &\leq 184\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma,\sigma\}N}} + \frac{36 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} \\
 &\leq 1508\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2 \max\{1-\gamma,\sigma\}N}}, \tag{84}
 \end{aligned}$$

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

797 C.3 Proof of the auxiliary lemmas

798 C.3.1 Proof of Lemma 7

799 To begin, note that there at leasts exist one state s_0 for any $V^{\pi,\sigma}$ such that $V^{\pi,\sigma}(s_0) =$
 800 $\min_{s \in \mathcal{S}} V^{\pi,\sigma}(s)$. With this in mind, for any policy π , one has by the definition in (??) and the
 801 Bellman's equation (??),

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

802 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
 803 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
 804 the values of some elements of $P_{s,a}$ to obey $P_{s,a} \geq \tilde{P}_{s,a} \geq 0$ and $\sum_{s'} (P_{s,a}(s') - \tilde{P}_{s,a}(s')) = \sigma$.
 805 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
 806 $\frac{1}{2} \|\tilde{P}_{s,a} + \sigma e_{s_0}^{\top} - P_{s,a}\|_1 \leq \frac{1}{2} \|\tilde{P}_{s,a} - P_{s,a}\|_1 + \frac{\sigma}{2} = \sigma$. Consequently,

$$\begin{aligned}
 \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s,a})} \mathcal{P}V^{\pi,\sigma} &\leq \left(\tilde{P}_{s,a} + \sigma e_{s_0}^{\top} \right) V^{\pi,\sigma} \leq \|\tilde{P}_{s,a}\|_1 \|V^{\pi,\sigma}\|_{\infty} + \sigma V^{\pi,\sigma}(s_0) \\
 &\leq (1-\sigma) \max_{s \in \mathcal{S}} V^{\pi,\sigma}(s) + \sigma \min_{s \in \mathcal{S}} V^{\pi,\sigma}(s), \tag{85}
 \end{aligned}$$

807 where the second inequality holds by $\|\tilde{P}_{s,a}\|_1 = \sum_{s'} \tilde{P}_{s,a}(s') = -\sum_{s'} (P_{s,a}(s') - \tilde{P}_{s,a}(s')) +$
 808 $\sum_{s'} P_{s,a}(s') = 1 - \sigma$. Plugging this back to the previous relation gives

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

809 which, by rearranging terms, immediately yields

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

810 **C.3.2 Proof of Lemma 8**

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

$$\begin{aligned} (I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} &= \frac{1}{1 - \gamma} (1 - \gamma) (I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} \\ &\stackrel{(i)}{\leq} \frac{1}{1 - \gamma} \sqrt{(1 - \gamma) (I - \gamma P_\pi)^{-1} \text{Var}_{P_\pi}(V^{\pi,P})} \\ &= \sqrt{\frac{1}{1 - \gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \text{Var}_{P_\pi}(V^{\pi,P})}, \end{aligned} \quad (86)$$

813 where (i) holds by Jensen's inequality.

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

$$\begin{aligned} &\text{Var}_{P_\pi}(V^{\pi,P}) \\ &\stackrel{(i)}{=} \text{Var}_{P_\pi}(V') = P_\pi(V' \circ V') - (P_\pi V') \circ (P_\pi V') \\ &\stackrel{(ii)}{=} P_\pi(V' \circ V') - \frac{1}{\gamma^2} (V' - r_\pi + (1 - \gamma)V_{\min} \mathbf{1}) \circ (V' - r_\pi + (1 - \gamma)V_{\min} \mathbf{1}) \\ &= P_\pi(V' \circ V') - \frac{1}{\gamma^2} V' \circ V' + \frac{2}{\gamma^2} V' \circ (r_\pi - (1 - \gamma)V_{\min} \mathbf{1}) \\ &\quad - \frac{1}{\gamma^2} (r_\pi - (1 - \gamma)V_{\min} \mathbf{1}) \circ (r_\pi - (1 - \gamma)V_{\min} \mathbf{1}) \\ &\leq P_\pi(V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_\infty \mathbf{1}, \end{aligned} \quad (87)$$

816 where (i) holds by the fact that $\text{Var}_{P_\pi}(V^{\pi,P} - b \mathbf{1}) = \text{Var}_{P_\pi}(V^{\pi,P})$ for any scalar b and $V^{\pi,P} \in \mathbb{R}^S$,
817 (ii) follows from $V' = r_\pi + \gamma P_\pi V^{\pi,P} - V_{\min} \mathbf{1} = r_\pi - (1 - \gamma)V_{\min} \mathbf{1} + \gamma P_\pi V'$, and the last line
818 arises from $\frac{1}{\gamma^2} V' \circ V' \geq \frac{1}{\gamma} V' \circ V'$ and $\|r_\pi - (1 - \gamma)V_{\min} \mathbf{1}\|_\infty \leq 1$. Plugging (87) back to (86)
819 leads to

$$\begin{aligned} (I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} &\leq \sqrt{\frac{1}{1 - \gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \left(P_\pi(V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_\infty \mathbf{1} \right)} \\ &\stackrel{(i)}{\leq} \sqrt{\frac{1}{1 - \gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \left(P_\pi(V' \circ V') - \frac{1}{\gamma} V' \circ V' \right) \right|} + \sqrt{\frac{1}{1 - \gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \frac{2}{\gamma^2} \|V'\|_\infty \mathbf{1}} \\ &\leq \sqrt{\frac{1}{1 - \gamma}} \sqrt{\left| \left(\sum_{t=0}^{\infty} \gamma^t (P_\pi)^{t+1} - \sum_{t=0}^{\infty} \gamma^{t-1} (P_\pi)^t \right) (V' \circ V') \right|} + \sqrt{\frac{2 \|V'\|_\infty \mathbf{1}}{\gamma^2 (1 - \gamma)^2}} \\ &\stackrel{(ii)}{\leq} \sqrt{\frac{\|V'\|_\infty^2 \mathbf{1}}{\gamma (1 - \gamma)}} + \sqrt{\frac{2 \|V'\|_\infty \mathbf{1}}{\gamma^2 (1 - \gamma)^2}} \\ &\leq \sqrt{\frac{8 \|V'\|_\infty \mathbf{1}}{\gamma^2 (1 - \gamma)^2}}, \end{aligned} \quad (88)$$

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

822 **C.3.3 Proof of Lemma 9**

823 **Step 1: controlling the point-wise concentration.** We first consider a more general term w.r.t. any
824 fixed (independent from \hat{P}^0) value vector V obeying $0 \leq V \leq \frac{1}{1 - \gamma} \mathbf{1}$ and any policy π . Invoking

825 Lemma 4 leads to that for any $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned} \left| \widehat{P}_{s,a}^{\pi,V} V - P_{s,a}^{\pi,V} V \right| &\leq \left| \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ \widehat{P}_{s,a}^0 [V]_\alpha - \sigma \left(\alpha - \min_{s'} [V]_\alpha (s') \right) \right\} \right. \\ &\quad \left. - \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P_{s,a}^0 [V]_\alpha - w\sigma \left(\alpha - \min_{s'} [V]_\alpha (s') \right) \right\} \right| \\ &\leq \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \underbrace{\left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_\alpha \right|}_{=: g_{s,a}(\alpha, V)}, \end{aligned} \quad (89)$$

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

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

$$\begin{aligned} g_{s,a}(\alpha, V) &= \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_\alpha \right| \leq \sqrt{\frac{2 \log(\frac{2}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} + \frac{2 \log(\frac{2}{\delta})}{3N(1-\gamma)} \\ &\leq \sqrt{\frac{2 \log(\frac{2}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{2 \log(\frac{2}{\delta})}{3N(1-\gamma)}. \end{aligned} \quad (90)$$

829 **Step 2: deriving the uniform concentration.** To obtain the union bound, we first notice that
830 $g_{s,a}(\alpha, V)$ is 1-Lipschitz w.r.t. α for any V obeying $\|V\|_\infty \leq \frac{1}{1-\gamma}$. In addition, we can construct
831 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
832 bound and (90), it holds with probability at least $1 - \frac{\delta}{SA}$ that for all $\alpha \in N_{\varepsilon_1}$,

$$g_{s,a}(\alpha, V) \leq \sqrt{\frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{3N(1-\gamma)}. \quad (91)$$

833 Combined with (89), it yields that,

$$\begin{aligned} \left| \widehat{P}_{s,a}^{\pi,V} V - P_{s,a}^{\pi,V} V \right| &\leq \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_\alpha \right| \\ &\stackrel{(i)}{\leq} \varepsilon_1 + \sup_{\alpha \in N_{\varepsilon_1}} \left| \left(P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_\alpha \right| \\ &\stackrel{(ii)}{\leq} \varepsilon_1 + \sqrt{\frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{3N(1-\gamma)} \end{aligned} \quad (92)$$

$$\begin{aligned} &\stackrel{(iii)}{\leq} \sqrt{\frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{\log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N(1-\gamma)} \\ &\stackrel{(iv)}{\leq} 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} \end{aligned} \quad (93)$$

$$\begin{aligned} &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \|V\|_\infty + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} \\ &\leq 3\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \end{aligned} \quad (94)$$

834 where (i) follows from that the optimal α^* falls into the ε_1 -ball centered around some point inside
835 N_{ε_1} and $g_{s,a}(\alpha, V)$ is 1-Lipschitz, (ii) holds by (91), (iii) arises from taking $\varepsilon_1 = \frac{\log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{3N(1-\gamma)}$, (iv)
836 is verified by $|N_{\varepsilon_1}| \leq \frac{3}{\varepsilon_1(1-\gamma)} \leq 9N$, and the last inequality is due to the fact $\|V^{\star, \sigma}\|_\infty \leq \frac{1}{1-\gamma}$ and
837 letting $N \geq \log(\frac{18SAN}{\delta})$.

838 To continue, applying (93) and (94) with $\pi = \pi^*$ and $V = V^{*,\sigma}$ (independent with \widehat{P}^0) and taking
839 the union bound over $(s, a) \in \mathcal{S} \times \mathcal{A}$ gives that with probability at least $1 - \delta$, it holds simultaneously
840 for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ that

$$\begin{aligned} \left| \widehat{P}_{s,a}^{\pi^*,V} V^{*,\sigma} - P_{s,a}^{\pi^*,V} V^{*,\sigma} \right| &\leq 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V^{*,\sigma})} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} \\ &\leq 3 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}}. \end{aligned} \quad (95)$$

841 By converting (95) to the matrix form, one has with probability at least $1 - \delta$,

$$\begin{aligned} \left| \widehat{\underline{P}}^{\pi^*,V} V^{\pi^*,\sigma} - \underline{P}^{\pi^*,V} V^{\pi^*,\sigma} \right| &\leq 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{*,\sigma})} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} \mathbf{1} \\ &\leq 3 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \mathbf{1}. \end{aligned} \quad (96)$$

842 C.3.4 Proof of Lemma 10

843 Following the same argument as (86), it follows

$$\left(I - \gamma \widehat{\underline{P}}^{\pi^*,V} \right)^{-1} \sqrt{\text{Var}_{\widehat{\underline{P}}^{\pi^*,V}}(V^{*,\sigma})} = \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left(\widehat{\underline{P}}^{\pi^*,V} \right)^t \text{Var}_{\widehat{\underline{P}}^{\pi^*,V}}(V^{*,\sigma})}. \quad (97)$$

844 To continue, we first focus on controlling $\text{Var}_{\widehat{\underline{P}}^{\pi^*,V}}(V^{*,\sigma})$. Towards this, denoting the minimum
845 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
846 Bellman's consistency equation in (41))

$$\begin{aligned} V' &= V^{*,\sigma} - V_{\min} \mathbf{1} = r_{\pi^*} + \gamma \underline{P}^{\pi^*,V} V^{*,\sigma} - V_{\min} \mathbf{1} \\ &= r_{\pi^*} + \gamma \widehat{\underline{P}}^{\pi^*,V} V^{*,\sigma} + \gamma \left(\underline{P}^{\pi^*,V} - \widehat{\underline{P}}^{\pi^*,V} \right) V^{*,\sigma} - V_{\min} \mathbf{1} \\ &= r_{\pi^*} - (1-\gamma) V_{\min} \mathbf{1} + \gamma \widehat{\underline{P}}^{\pi^*,V} V' + \gamma \left(\underline{P}^{\pi^*,V} - \widehat{\underline{P}}^{\pi^*,V} \right) V^{*,\sigma} \\ &= r'_{\pi^*} + \gamma \widehat{\underline{P}}^{\pi^*,V} V' + \gamma \left(\underline{P}^{\pi^*,V} - \widehat{\underline{P}}^{\pi^*,V} \right) V^{*,\sigma}, \end{aligned} \quad (98)$$

847 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,
848 we control $\text{Var}_{\widehat{\underline{P}}^{\pi^*,V}}(V^{*,\sigma})$ as follows:

$$\begin{aligned} \text{Var}_{\widehat{\underline{P}}^{\pi^*,V}}(V^{*,\sigma}) &\stackrel{(i)}{=} \text{Var}_{\widehat{\underline{P}}^{\pi^*,V}}(V') = \widehat{\underline{P}}^{\pi^*,V} (V' \circ V') - \left(\widehat{\underline{P}}^{\pi^*,V} V' \right) \circ \left(\widehat{\underline{P}}^{\pi^*,V} V' \right) \\ &\stackrel{(ii)}{=} \widehat{\underline{P}}^{\pi^*,V} (V' \circ V') - \frac{1}{\gamma^2} \left(V' - r'_{\pi^*} - \gamma \left(\underline{P}^{\pi^*,V} - \widehat{\underline{P}}^{\pi^*,V} \right) V^{*,\sigma} \right)^{\circ 2} \\ &= \widehat{\underline{P}}^{\pi^*,V} (V' \circ V') - \frac{1}{\gamma^2} V' \circ V' + \frac{2}{\gamma^2} V' \circ \left(r'_{\pi^*} + \gamma \left(\underline{P}^{\pi^*,V} - \widehat{\underline{P}}^{\pi^*,V} \right) V^{*,\sigma} \right) \\ &\quad - \frac{1}{\gamma^2} \left(r'_{\pi^*} + \gamma \left(\underline{P}^{\pi^*,V} - \widehat{\underline{P}}^{\pi^*,V} \right) V^{*,\sigma} \right)^{\circ 2} \\ &\stackrel{(iii)}{\leq} \widehat{\underline{P}}^{\pi^*,V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} \mathbf{1} + \frac{2}{\gamma} \|V'\|_{\infty} \left| \left(\underline{P}^{\pi^*,V} - \widehat{\underline{P}}^{\pi^*,V} \right) V^{*,\sigma} \right| \end{aligned} \quad (99)$$

$$\leq \widehat{\underline{P}}^{\pi^*,V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} \mathbf{1} + \frac{6}{\gamma} \|V'\|_{\infty} \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \mathbf{1}, \quad (100)$$

849 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
850 from (98), (iii) arises from $\frac{1}{\gamma^2} V' \circ V' \geq \frac{1}{\gamma} V' \circ V'$ and $-1 \leq r_{\pi^*} - (1-\gamma) V_{\min} \mathbf{1} = r'_{\pi^*} \leq r_{\pi^*} \leq 1$,
851 and the last inequality holds by Lemma 9.

852 Plugging (100) into (97) leads to

$$\begin{aligned}
& (I - \gamma \widehat{P}^{\pi^*, V})^{-1} \sqrt{\text{Var}_{\widehat{P}^{\pi^*, V}}(V^*, \sigma)} \\
& \leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (\widehat{P}^{\pi^*, V})^t \left(\widehat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} 1 + \frac{6}{\gamma} \|V'\|_{\infty} \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} 1 \right)} \\
& \stackrel{(i)}{\leq} \sqrt{\frac{1}{1-\gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t (\widehat{P}^{\pi^*, V})^t \left(\widehat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' \right) \right|} \\
& \quad + \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (\widehat{P}^{\pi^*, V})^t \left(\frac{2}{\gamma^2} \|V'\|_{\infty} 1 + \frac{6}{\gamma} \|V'\|_{\infty} \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} 1 \right)} \\
& \leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t (\widehat{P}^{\pi^*, V})^t \left[\widehat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' \right] \right|} + \sqrt{\frac{\left(2 + 6 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right) \|V'\|_{\infty}}{(1-\gamma)^2 \gamma^2}} 1, \tag{101}
\end{aligned}$$

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

$$\begin{aligned}
& \left| \sum_{t=0}^{\infty} \gamma^t (\widehat{P}^{\pi^*, V})^t \left(\widehat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' \right) \right| \\
& = \left| \left(\sum_{t=0}^{\infty} \gamma^t (\widehat{P}^{\pi^*, V})^{t+1} - \sum_{t=0}^{\infty} \gamma^{t-1} (\widehat{P}^{\pi^*, V})^t \right) (V' \circ V') \right| \leq \frac{1}{\gamma} \|V'\|_{\infty}^2 1 \tag{102}
\end{aligned}$$

855 by recursion. Inserting (102) back to (101) leads to

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

856 where the penultimate inequality follows from applying Lemma 7 with $P = P^0$ and $\pi = \pi^*$:

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

857 C.3.5 Proof of Lemma 12

858 To begin with, for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, invoking the results in (89), we have

$$\begin{aligned}
& \left| \widehat{P}_{s,a}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - P_{s,a}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right| \leq \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha} \right| \\
& \stackrel{(i)}{\leq} \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left(\left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}^{*, \sigma}]_{\alpha} \right| + \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) \left([\widehat{V}^{\widehat{\pi}, \sigma}]_{\alpha} - [\widehat{V}^{*, \sigma}]_{\alpha} \right) \right| \right)
\end{aligned}$$

$$\begin{aligned}
&\leq \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left(\left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}^{*, \sigma}]_\alpha \right| + \|P_{s,a}^0 - \widehat{P}_{s,a}^0\|_1 \left\| [\widehat{V}^{\widehat{\pi}, \sigma}]_\alpha - [\widehat{V}^{*, \sigma}]_\alpha \right\|_\infty \right) \\
&\stackrel{(ii)}{\leq} \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}^{*, \sigma}]_\alpha \right| + 2 \left\| \widehat{V}^{\widehat{\pi}, \sigma} - \widehat{V}^{*, \sigma} \right\|_\infty \\
&\stackrel{(iii)}{\leq} \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}^{*, \sigma}]_\alpha \right| + \frac{2\gamma \varepsilon_{\text{opt}}}{1-\gamma}, \tag{104}
\end{aligned}$$

859 where (i) holds by the triangle inequality, and (ii) follows from $\|P_{s,a}^0 - \widehat{P}_{s,a}^0\|_1 \leq 2$ and
860 $\left\| [\widehat{V}^{\widehat{\pi}, \sigma}]_\alpha - [\widehat{V}^{*, \sigma}]_\alpha \right\|_\infty \leq \left\| \widehat{V}^{\widehat{\pi}, \sigma} - \widehat{V}^{*, \sigma} \right\|_\infty$, and (iii) follows from (48).

861 To control $\left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}^{*, \sigma}]_\alpha \right|$ in (104) for any given $\alpha \in [0, \frac{1}{1-\gamma}]$, and tame the dependency
862 between $\widehat{V}^{*, \sigma}$ and \widehat{P}^0 , we resort to the following leave-one-out argument motivated by (Agarwal et al.,
863 2020; Li et al., 2022a; Shi and Chi, 2022). Specifically, we first construct a set of auxiliary RMDPs
864 which simultaneously have the desired statistical independence between robust value functions and
865 the estimated nominal transition kernel, and are minimally different from the original RMDPs under
866 consideration. Then we control the term of interest associated with these auxiliary RMDPs and show
867 the value is close to the target quantity for the desired RMDP. The process is divided into several
868 steps as below.

869 **Step 1: construction of auxiliary RMDPs with deterministic empirical nominal transitions.**

870 Recall that we target the empirical infinite-horizon robust MDP $\widehat{\mathcal{M}}_{\text{rob}}$ with the nominal transition
871 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
872 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 .
873 In particular, we define the nominal transition kernel and reward function of $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ as $P^{s,u}$ and $r^{s,u}$,
874 which are expressed as follows

$$\begin{cases} P^{s,u}(s' | s, a) = \mathbb{1}(s' = s) & \text{for all } (s', a) \in \mathcal{S} \times \mathcal{A}, \\ P^{s,u}(\cdot | \tilde{s}, a) = \widehat{P}^0(\cdot | \tilde{s}, a) & \text{for all } (\tilde{s}, a) \in \mathcal{S} \times \mathcal{A} \text{ and } \tilde{s} \neq s, \end{cases} \tag{105}$$

875 and

$$\begin{cases} r^{s,u}(s, a) = u & \text{for all } a \in \mathcal{A}, \\ r^{s,u}(\tilde{s}, a) = r(\tilde{s}, a) & \text{for all } (\tilde{s}, a) \in \mathcal{S} \times \mathcal{A} \text{ and } \tilde{s} \neq s. \end{cases} \tag{106}$$

876 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
877 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 ,
878 which will be leveraged later.

879 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

$$\forall (\tilde{s}, a) \in \mathcal{S} \times \mathcal{A}: \quad \widehat{\mathcal{T}}_{s,u}^\sigma(Q)(\tilde{s}, a) = r^{s,u}(\tilde{s}, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{\tilde{s},a}^{s,u})} \mathcal{P}V, \quad \text{with } V(\tilde{s}) = \max_a Q(\tilde{s}, a). \tag{107}$$

880 **Step 2: fixed-point equivalence between $\widehat{\mathcal{M}}_{\text{rob}}$ and the auxiliary RMDP $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$.** Recall that
881 $\widehat{Q}^{*, \sigma}$ is the unique fixed point of $\widehat{\mathcal{T}}^\sigma(\cdot)$ with the corresponding robust value $\widehat{V}^{*, \sigma}$. We assert that the
882 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
883 robust value function $\widehat{V}^{*, \sigma}$ derived from $\widehat{\mathcal{T}}^\sigma(\cdot)$, as long as we choose u in the following manner:

$$u^* := u^*(s) = \widehat{V}^{*, \sigma}(s) - \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(e_s)} \mathcal{P}\widehat{V}^{*, \sigma}. \tag{108}$$

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

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

$$r^{s,u^*}(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^{s,u^*})} \mathcal{P}\widehat{V}^{*, \sigma} = u^* + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(e_s)} \mathcal{P}\widehat{V}^{*, \sigma}$$

$$= \widehat{V}^{*,\sigma}(s) - \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(e_s)} \mathcal{P} \widehat{V}^{*,\sigma} + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(e_s)} \mathcal{P} \widehat{V}^{*,\sigma} = \widehat{V}^{*,\sigma}(s), \quad (109)$$

887 where the first equality follows from the definition of $P_{s,a}^{s,u^*}$ in (105), and the second equality
888 follows from plugging in the definition of u^* in (108).

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

$$\begin{aligned} r^{s,u^*}(s', a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s',a}^{s,u^*})} \mathcal{P} \widehat{V}^{*,\sigma} &= r(s', a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s',a}^0)} \mathcal{P} \widehat{V}^{*,\sigma} \\ &= \widehat{\mathcal{T}}^\sigma(\widehat{Q}^{*,\sigma})(s', a) = \widehat{Q}^{*,\sigma}(s', a), \end{aligned} \quad (110)$$

890 where the first equality follows from the definitions in (106) and (105), and the last line
891 arises from the definition of the robust Bellman operator in (12), and that $\widehat{Q}^{*,\sigma}$ is the fixed
892 point of $\widehat{\mathcal{T}}^\sigma(\cdot)$ (see Lemma 3).

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

$$\begin{cases} \widehat{Q}_{s,u^*}^{*,\sigma}(s, a) = \widehat{V}^{*,\sigma}(s) & \text{for all } a \in \mathcal{A}, \\ \widehat{Q}_{s,u^*}^{*,\sigma}(s', a) = \widehat{Q}^{*,\sigma}(s', a) & \text{for all } s' \neq s \text{ and } a \in \mathcal{A}. \end{cases} \quad (111)$$

895 Consequently, we confirm the existence of a fixed point of the operator $\widehat{\mathcal{T}}_{s,u^*}^\sigma(\cdot)$. In addition, its
896 corresponding value function $\widehat{V}_{s,u^*}^{*,\sigma}$ also coincides with $\widehat{V}^{*,\sigma}$. Note that the corresponding facts
897 between $\widehat{\mathcal{M}}_{\text{rob}}$ and $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ in Step 1 and step 2 holds in fact for any uncertainty set.

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

$$0 \leq u^* \leq \widehat{V}^{*,\sigma}(s) \leq \frac{1}{1-\gamma}. \quad (112)$$

899 We can construct a N_{ε_2} -net over the interval $[0, \frac{1}{1-\gamma}]$, where the size is bounded by $|N_{\varepsilon_2}| \leq \frac{3}{\varepsilon_2(1-\gamma)}$
900 (Vershynin, 2018). Following the same arguments in the proof of Lemma 3, we can demonstrate
901 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
902 $0 \leq \widehat{Q}_{s,u}^{*,\sigma} \leq \frac{1}{1-\gamma} \cdot 1$. Consequently, the corresponding robust value function also satisfies $\|\widehat{V}_{s,u}^{*,\sigma}\|_\infty \leq$
903 $\frac{1}{1-\gamma}$.

904 By the definitions in (105) and (106), we observe that for all $u \in N_{\varepsilon_2}$, $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ is statistically in-
905 dependent 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
906 fixed α . With this in mind, invoking the fact in (93) and (94) and taking the union bound over
907 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
908 $(s, a, u) \in \mathcal{S} \times \mathcal{A} \times N_{\varepsilon_2}$ that

$$\begin{aligned} & \max_{\alpha \in [0, 1/(1-\gamma)]} \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}_{s,u}^{*,\sigma}]_\alpha \right| \\ & \leq \varepsilon_2 + 2 \sqrt{\frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}_{s,u}^{*,\sigma})} + \frac{2 \log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{3N(1-\gamma)} \\ & \leq \varepsilon_2 + 3 \sqrt{\frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{(1-\gamma)^2 N}}, \end{aligned} \quad (113)$$

909 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 $N \geq$
910 $\log\left(\frac{18SAN|N_{\varepsilon_2}|}{\delta}\right)$.

911 **Step 4: uniform concentration.** Recalling that $u^* \in [0, \frac{1}{1-\gamma}]$ (see (112)), we can always find
 912 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 (107) yields

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

913 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

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

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

$$\left\| \widehat{Q}_{s,\bar{u}}^{*,\sigma} - \widehat{Q}_{s,u^*}^{*,\sigma} \right\|_\infty \leq \frac{\varepsilon_2}{(1-\gamma)} \quad \text{and} \quad \left\| \widehat{V}_{s,\bar{u}}^{*,\sigma} - \widehat{V}_{s,u^*}^{*,\sigma} \right\|_\infty \leq \left\| \widehat{Q}_{s,\bar{u}}^{*,\sigma} - \widehat{Q}_{s,u^*}^{*,\sigma} \right\|_\infty \leq \frac{\varepsilon_2}{(1-\gamma)}. \quad (114)$$

915 Armed with the above facts, to control the first term in (104), invoking the identity $\widehat{V}^{*,\sigma} = \widehat{V}_{s,u^*}^{*,\sigma}$
 916 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}^{\hat{\pi},\sigma}(s), \max_s \widehat{V}^{\hat{\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{2\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,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})} \\ &\quad + \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)} \quad (115) \\ &\leq 10\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2 N}}, \quad (116) \end{aligned}$$

917 where (i) holds by the triangle inequality, (ii) arises from (the last inequality holds by (114))

$$\left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) \left([\widehat{V}_{s,\bar{u}}^{*,\sigma}]_\alpha - [\widehat{V}_{s,u^*}^{*,\sigma}]_\alpha \right) \right| \leq \left\| P_{s,a}^0 - \widehat{P}_{s,a}^0 \right\|_1 \left\| [\widehat{V}_{s,\bar{u}}^{*,\sigma}]_\alpha - [\widehat{V}_{s,u^*}^{*,\sigma}]_\alpha \right\|_\infty$$

$$\leq 2 \left\| \widehat{V}_{s,\bar{u}}^{*,\sigma} - \widehat{V}_{s,u^*}^{*,\sigma} \right\|_{\infty} \leq \frac{2\varepsilon_2}{(1-\gamma)}, \quad (117)$$

918 (iii) follows from (113), (iv) can be verified by applying Lemma 2 with (114). Here, the penultimate
 919 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
 920 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 $N \geq \log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)$.

921 **Step 5: finishing up.** Inserting (115) and (116) back into (104) and combining with (116) give that
 922 with probability at least $1 - \delta$,

$$\begin{aligned} \left| \widehat{P}_{s,a}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - P_{s,a}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right| &\leq \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi},\sigma}(s), \max_s \widehat{V}^{\widehat{\pi},\sigma}(s)]} \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}^{*,\sigma}]_{\alpha} \right| + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \\ &\leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| (P_{s,a}^0 - \widehat{P}_{s,a}^0) [\widehat{V}^{*,\sigma}]_{\alpha} \right| + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \\ &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma})} + \frac{8\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \\ &\leq 10\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \end{aligned} \quad (118)$$

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

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

$$\begin{aligned} \left| \widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} \widehat{\underline{V}}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{\underline{V}}^{\widehat{\pi},\sigma} \right| &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{*,\sigma})} \mathbf{1} + \frac{8\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} \mathbf{1} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \mathbf{1} \\ &\leq 10\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2N}} \mathbf{1} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \mathbf{1}. \end{aligned} \quad (119)$$

925 C.3.6 Proof of Lemma 13

926 The proof can be achieved by directly applying the same routine as Appendix C.3.4. Towards this,
 927 similar to (97), we arrive at

$$\left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{\underline{V}}^{\widehat{\pi},\sigma})} \leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left(\underline{P}^{\widehat{\pi},\widehat{V}} \right)^t \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{\underline{V}}^{\widehat{\pi},\sigma})}. \quad (120)$$

928 To control $\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{\underline{V}}^{\widehat{\pi},\sigma})$, 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
 929 $V' := \widehat{V}^{\widehat{\pi},\sigma} - V_{\min} \mathbf{1}$. By the same argument as (99), we arrive at

$$\begin{aligned} \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{\underline{V}}^{\widehat{\pi},\sigma}) &\leq \underline{P}^{\widehat{\pi},\widehat{V}} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} \mathbf{1} + \frac{2}{\gamma} \|V'\|_{\infty} \left| \left(\widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} - \underline{P}^{\widehat{\pi},\widehat{V}} \right) \widehat{\underline{V}}^{\widehat{\pi},\sigma} \right| \\ &\leq \underline{P}^{\widehat{\pi},\widehat{V}} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} \mathbf{1} + \frac{2}{\gamma} \|V'\|_{\infty} \left(10\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \right) \mathbf{1}, \end{aligned} \quad (121)$$

930 where the last inequality makes use of Lemma 12. Plugging (121) back into (120) leads to

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

$$\begin{aligned}
&\stackrel{\text{(ii)}}{\leq} \sqrt{\frac{\|V'\|_\infty^2}{\gamma(1-\gamma)}} 1 + \sqrt{\frac{\left(2 + 20\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)^\delta})}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}\right) \|V'\|_\infty}{(1-\gamma)^2\gamma^2}} 1 \\
&\stackrel{\text{(iii)}}{\leq} \sqrt{\frac{\|V'\|_\infty^2}{\gamma(1-\gamma)}} 1 + \sqrt{\frac{24\|V'\|_\infty}{(1-\gamma)^2\gamma^2}} 1 \leq 6\sqrt{\frac{\|V'\|_\infty}{(1-\gamma)^2\gamma^2}} 1,
\end{aligned} \tag{122}$$

931 where (i) arises from following the routine of (101), (ii) holds by repeating the argument of (102),
932 (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 $\|V'\|_\infty \leq$
933 $\|V^{*,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$.

934 Finally, applying Lemma 7 with $P = \widehat{P}^0$ and $\pi = \widehat{\pi}$ yields

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

935 which can be inserted into (122) and gives

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

936 D Proof of the lower bound with TV distance: Theorem 2

937 To prove Theorem 2, we shall first construct some hard instances and then characterize the sample
938 complexity requirements over these instances. Note that the hard instances for robust MDPs are
939 different from those for standard MDPs, due to the asymmetric structure induced by the robust RL
940 problem formulation to consider the worst-case performance. By constructing a new class of hard
941 instances inspired by the asymmetric structure of the RMDP, we develop a new lower bound in
942 Theorem 2 that is tighter than prior art (Yang et al., 2022).

943 D.1 Construction of the hard problem instances

944 **Construction of two hard MDPs.** Suppose there are two standard MDPs defined as below:

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

945 Here, γ is the discount parameter, $\mathcal{S} = \{0, 1, \dots, S-1\}$ is the state space. Given any state
946 $s \in \{2, 3, \dots, S-1\}$, the corresponding action space are $\mathcal{A} = \{0, 1, 2, \dots, A-1\}$. While for
947 states $s = 0$ or $s = 1$, the action space is only $\mathcal{A}' = \{0, 1\}$. For any $\phi \in \{0, 1\}$, the transition kernel
948 P^ϕ of the constructed MDP \mathcal{M}_ϕ is defined as

$$P^\phi(s' \mid s, a) = \begin{cases} p\mathbb{1}(s' = 1) + (1-p)\mathbb{1}(s' = 0) & \text{if } (s, a) = (0, \phi) \\ q\mathbb{1}(s' = 1) + (1-q)\mathbb{1}(s' = 0) & \text{if } (s, a) = (0, 1-\phi) \\ \mathbb{1}(s' = 1) & \text{if } s \geq 1 \end{cases}, \tag{123}$$

949 where p and q are set to satisfy

$$0 \leq p \leq 1 \quad \text{and} \quad 0 \leq q = p - \Delta \tag{124}$$

950 for some p and $\Delta > 0$ that shall be introduced later. The above transition kernel P^ϕ implies that state
951 1 is an absorbing state, namely, the MDP will always stay after it arrives at 1.

952 Then, we define the reward function as

$$r(s, a) = \begin{cases} 1 & \text{if } s = 1 \\ 0 & \text{otherwise} \end{cases}. \tag{125}$$

953 Additionally, we choose the following initial state distribution:

$$\varphi(s) = \begin{cases} 1, & \text{if } s = 0 \\ 0, & \text{otherwise} \end{cases}. \tag{126}$$

954 Here, the constructed two instances are set with different probability transition from state 0 with
955 reward 0 but not state 1 with reward 1 (which were used in standard MDPs (Li et al., 2022a)), yielding
956 a larger gap between the value functions of the two instances.

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

$$\mathcal{U}^\sigma(P) := \mathcal{U}_{\text{TV}}^\sigma(P) = \otimes \mathcal{U}_{\text{TV}}^\sigma(P_{s,a}), \quad \mathcal{U}_{\text{TV}}^\sigma(P_{s,a}) := \left\{ P'_{s,a} \in \Delta(\mathcal{S}) : \frac{1}{2} \|P'_{s,a} - P_{s,a}\|_1 \leq \sigma \right\}, \quad (127)$$

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

$$p = (1 + c_1) \max\{1 - \gamma, \sigma\} \quad \text{and} \quad \Delta \leq c_1 \max\{1 - \gamma, \sigma\}, \quad (128)$$

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

$$(1 + c_1) \sigma \leq 1 - c_0 + c_1 \sigma \leq 1 - \frac{c_0}{2} < 1, \quad (1 + c_1) (1 - \gamma) \leq \frac{3}{2} (1 - \gamma) \leq \frac{3}{4} < 1. \quad (129)$$

963 Consequently, applying (124) directly leads to

$$p \geq q \geq \max\{1 - \gamma, \sigma\}. \quad (130)$$

964 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
 965 next state s' associated with any perturbed transition kernel $P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^\phi)$ as

$$\underline{P}^\phi(s' | s, a) := \inf_{P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^\phi)} P(s' | s, a) = \max\{P(s' | s, a) - \sigma, 0\}, \quad (131)$$

966 where the last equation can be easily verified by the definition of $\mathcal{U}^\sigma(P^\phi)$ in (127). As shall be seen,
 967 the transition from state 0 to state 1 plays an important role in the analysis, for convenience, we
 968 denote

$$\underline{p} := \underline{P}^\phi(1 | 0, \phi) = p - \sigma, \quad \underline{q} := \underline{P}^\phi(1 | 0, 1 - \phi) = q - \sigma, \quad (132)$$

969 which follows from the fact that $p \geq q \geq \sigma$ in (130).

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

976 **Lemma 14.** *For any $\phi = \{0, 1\}$ and any policy π , the robust value function obeys*

$$V_\phi^{\pi, \sigma}(0) = \frac{\gamma \left(z_\phi^\pi - \sigma \right)}{(1 - \gamma) \left(1 + \frac{\gamma(z_\phi^\pi - \sigma)}{1 - \gamma(1 - \sigma)} \right) (1 - \gamma(1 - \sigma))}, \quad (133)$$

977 where z_ϕ^π is defined as

$$z_\phi^\pi := p\pi(\phi | 0) + q\pi(1 - \phi | 0). \quad (134)$$

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

$$V_\phi^{*, \sigma}(0) = \frac{\gamma(p - \sigma)}{(1 - \gamma) \left(1 + \frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)} \right) (1 - \gamma(1 - \sigma))}, \quad (135a)$$

$$\pi_\phi^*(\phi | s) = 1, \quad \text{for } s \in \mathcal{S}. \quad (135b)$$

979 D.2 Establishing the minimax lower bound

980 Note that our goal is to control the quantity w.r.t. any policy estimator $\hat{\pi}$ based on the chosen initial
 981 distribution φ in (126) and the dataset consisting of N samples over each state-action pair generated
 982 from the nominal transition kernel P^ϕ , which gives

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

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

$$\Delta = 32(1-\gamma) \max\{1-\gamma, \sigma\} \varepsilon \leq c_1 \max\{1-\gamma, \sigma\} \quad (136)$$

985 which satisfies (128), it leads to that for any policy $\hat{\pi}$,

$$\langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\hat{\pi},\sigma} \rangle \geq 2\varepsilon(1 - \hat{\pi}(\phi|0)). \quad (137)$$

986 With this connection established between the policy $\hat{\pi}$ and its sub-optimality gap as depicted in (137),
 987 we can now proceed to build an estimate for ϕ . Here, we denote \mathbb{P}_ϕ as the probability distribution
 988 when the MDP is \mathcal{M}_ϕ , where ϕ can take on values in the set $\{0, 1\}$.

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

$$\mathbb{P}_\phi \left\{ \langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\hat{\pi},\sigma} \rangle \leq \varepsilon \right\} \geq \frac{7}{8}, \quad (138)$$

990 then in view of (137), we necessarily have $\hat{\pi}(\phi|0) \geq \frac{1}{2}$ with probability at least $\frac{7}{8}$. With this in mind,
 991 we are motivated to construct the following estimate $\hat{\phi}$ for $\phi \in \{0, 1\}$:

$$\hat{\phi} = \arg \max_{a \in \{0,1\}} \hat{\pi}(a|0), \quad (139)$$

992 which obeys

$$\mathbb{P}_\phi \{ \hat{\phi} = \phi \} \geq \mathbb{P}_\phi \{ \hat{\pi}(\phi|0) > 1/2 \} \geq \frac{7}{8}. \quad (140)$$

993 Subsequently, our aim is to demonstrate that (140) cannot occur without an adequate number of
 994 samples, which would in turn contradict (137).

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

$$p_e := \inf_{\psi} \max \{ \mathbb{P}_0(\psi \neq 0), \mathbb{P}_1(\psi \neq 1) \}. \quad (141)$$

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

1000 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
 1001 the nominal transition kernel P^ϕ associated with \mathcal{M}_ϕ and the samples are generated independently.
 1002 Applying standard results from [Tsybakov and Zaiats \(2009, Theorem 2.2\)](#) and the additivity of the
 1003 KL divergence (cf. [Tsybakov and Zaiats \(2009, Page 85\)](#)), we obtain

$$\begin{aligned} p_e &\geq \frac{1}{4} \exp \left(-NSAKL(\mu_0 \parallel \mu_1) \right) \\ &= \frac{1}{4} \exp \left\{ -N \left(\text{KL}(P^0(\cdot|0,0) \parallel P^1(\cdot|0,0)) + \text{KL}(P^0(\cdot|0,1) \parallel P^1(\cdot|0,1)) \right) \right\}, \end{aligned} \quad (142)$$

1004 where the last inequality holds by observing that

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

1005 Here, the last equality holds by the fact that $P^0(\cdot|s,a)$ and $P^1(\cdot|s,a)$ only differ when $s = 0$.

1006 Now, our focus shifts towards bounding the terms involving the KL divergence in (142). Given
 1007 $p \geq q \geq \max\{1-\gamma, \sigma\}$ (cf. (130)), applying Lemma 1 (cf. (23)) gives

$$\text{KL}(P^0(\cdot|0,1) \parallel P^1(\cdot|0,1)) = \text{KL}(p \parallel q) \leq \frac{(p-q)^2}{(1-p)p} \stackrel{(i)}{=} \frac{\Delta^2}{p(1-p)}$$

$$\begin{aligned}
& \underline{\text{(ii)}} \frac{1024(1-\gamma)^2 \max\{1-\gamma, \sigma\}^2 \varepsilon^2}{p(1-p)} \\
& \leq \frac{1024(1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2}{1-p} \leq \frac{4096}{c_1} (1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2, \tag{143}
\end{aligned}$$

1008 where (i) stems from the definition in (124), (ii) follows by the expression of Δ in (136), and the last
1009 inequality arises from $1-q \geq 1-p \geq \frac{c_0}{4}$ (see (129)).

1010 Note that it can be shown that $\text{KL}(P^0(\cdot|0,0) \| P^1(\cdot|0,0))$ can be upper bounded in a same manner.
1011 Substituting (143) back into (142) demonstrates that: if the sample size is selected as

$$N \leq \frac{c_1 \log 2}{8192(1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2}, \tag{144}$$

1012 then one necessarily has

$$p_e \geq \frac{1}{4} \exp \left\{ -N \frac{8192}{c_1} (1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2 \right\} \geq \frac{1}{8}, \tag{145}$$

1013 **Step 3: putting the results together.** Lastly, suppose that there exists an estimator $\hat{\pi}$ such that

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

1014 According to Step 1, the estimator $\hat{\phi}$ defined in (139) must satisfy

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

1015 However, this cannot occur under the sample size condition (144) to avoid contradiction with (145).
1016 Thus, we have completed the proof.

1017 D.3 Proof of the auxiliary facts

1018 D.3.1 Proof of Lemma 14

1019 **Deriving the robust value function over different states.** For any \mathcal{M}_ϕ with $\phi \in \{0, 1\}$, we first
1020 characterize the robust value function of any policy π over different states. Before proceeding, we
1021 denote the minimum of the robust value function over states as below:

$$V_{\phi, \min}^{\pi, \sigma} := \min_{s \in \mathcal{S}} V_\phi^{\pi, \sigma}(s). \tag{146}$$

1022 Clearly, there exists at least one state $s_{\phi, \min}^\pi$ that satisfies $V_\phi^{\pi, \sigma}(s_{\phi, \min}^\pi) = V_{\phi, \min}^{\pi, \sigma}$.

1023 With this in mind, it is easily observed that for any policy π , the robust value function at state $s = 1$
1024 obeys

$$\begin{aligned}
V_\phi^{\pi, \sigma}(1) &= \mathbb{E}_{a \sim \pi(\cdot|1)} \left[r(1, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{1,a}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} \right] \\
&\stackrel{\text{(i)}}{=} 1 + \gamma \mathbb{E}_{a \sim \pi(\cdot|1)} \left[\underline{P}^\phi(1|1, a) V_\phi^{\pi, \sigma}(1) \right] + \gamma \sigma V_{\phi, \min}^{\pi, \sigma} \stackrel{\text{(ii)}}{=} 1 + \gamma(1-\sigma) V_\phi^{\pi, \sigma}(1) + \gamma \sigma V_{\phi, \min}^{\pi, \sigma}, \tag{147}
\end{aligned}$$

1025 where (i) holds by $r(1, a) = 1$ for all $a \in \mathcal{A}'$ and (131), and (ii) follows from $P^\phi(1|1, a) = 1$ for all
1026 $a \in \mathcal{A}'$.

1027 Similarly, for any $s \in \{2, 3, \dots, S-1\}$, we have

$$\begin{aligned}
V_\phi^{\pi, \sigma}(s) &= 0 + \gamma \mathbb{E}_{a \sim \pi(\cdot|s)} \left[\underline{P}^\phi(1|s, a) V_\phi^{\pi, \sigma}(1) \right] + \gamma \sigma V_{\phi, \min}^{\pi, \sigma} \\
&= \gamma(1-\sigma) V_\phi^{\pi, \sigma}(1) + \gamma \sigma V_{\phi, \min}^{\pi, \sigma}, \tag{148}
\end{aligned}$$

1028 since $r(s, a) = 0$ for all $s \in \{2, 3, \dots, S-1\}$ and the definition in (131).

1029 Finally, we move onto compute $V_\phi^{\pi,\sigma}(0)$, the robust value function at state 0 associated with any
 1030 policy π . First, it obeys

$$\begin{aligned} V_\phi^{\pi,\sigma}(0) &= \mathbb{E}_{a \sim \pi(\cdot|0)} \left[r(0, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,a}^\phi)} \mathcal{P} V_\phi^{\pi,\sigma} \right] \\ &= 0 + \gamma \pi(\phi|0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,\phi}^\phi)} \mathcal{P} V_\phi^{\pi,\sigma} + \gamma \pi(1-\phi|0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,1-\phi}^\phi)} \mathcal{P} V_\phi^{\pi,\sigma}. \end{aligned} \quad (149)$$

1031 Recall the transition kernel defined in (123) and the fact about the uncertainty set over state 0 in
 1032 (132), it is easily verified that the following probability vector $P_1 \in \Delta(\mathcal{S})$ obeys $P_1 \in \mathcal{U}^\sigma(P_{0,\phi}^\phi)$,
 1033 which is defined as

$$\begin{aligned} P_1(0) &= 1 - p + \sigma \mathbb{1}(0 = s_{\phi,\min}^\pi), & P_1(1) &= \underline{p} = p - \sigma, \\ P_1(s) &= \sigma \mathbb{1}(s = s_{\phi,\min}^\pi), & \forall s \in \{2, 3, \dots, S-1\}, \end{aligned} \quad (150)$$

1034 where $\underline{p} = p - \sigma$ due to (132). Similarly, the following probability vector $P_2 \in \Delta(\mathcal{S})$ also falls into
 1035 the uncertainty set $\mathcal{U}^\sigma(P_{0,1-\phi}^\phi)$:

$$\begin{aligned} P_2(0) &= 1 - q + \sigma \mathbb{1}(0 = s_{\phi,\min}^\pi), & P_2(1) &= \underline{q} = q - \sigma, \\ P_2(s) &= \sigma \mathbb{1}(0 = s_{\phi,\min}^\pi) & \forall s \in \{2, 3, \dots, S-1\}. \end{aligned} \quad (151)$$

1036 It is noticed that P_0 and P_1 defined above are the worst-case perturbations, since the probability
 1037 mass at state 1 will be moved to the state with the least value. Plugging the above facts about
 1038 $P_1 \in \mathcal{U}^\sigma(P_{0,\phi}^\phi)$ and $P_2 \in \mathcal{U}^\sigma(P_{0,1-\phi}^\phi)$ into (149), we arrive at

$$\begin{aligned} V_\phi^{\pi,\sigma}(0) &\leq \gamma \pi(\phi|0) P_1 V_\phi^{\pi,\sigma} + \gamma \pi(1-\phi|0) P_2 V_\phi^{\pi,\sigma} \\ &= \gamma \pi(\phi|0) \left[(p - \sigma) V_\phi^{\pi,\sigma}(1) + (1 - p) V_\phi^{\pi,\sigma}(0) + \sigma V_{\phi,\min}^{\pi,\sigma} \right] \\ &\quad + \gamma \pi(1-\phi|0) \left[(q - \sigma) V_\phi^{\pi,\sigma}(1) + (1 - q) V_\phi^{\pi,\sigma}(0) + \sigma V_{\phi,\min}^{\pi,\sigma} \right] \\ &\stackrel{(i)}{=} \gamma (z_\phi^\pi - \sigma) V_\phi^{\pi,\sigma}(1) + \gamma \sigma V_{\phi,\min}^{\pi,\sigma} + \gamma (1 - z_\phi^\pi) V_\phi^{\pi,\sigma}(0), \end{aligned} \quad (152)$$

1039 where the last equality holds by the definition of z_ϕ^π in (134). To continue, recursively applying (152)
 1040 yields

$$\begin{aligned} &V_\phi^{\pi,\sigma}(0) \\ &\leq \gamma (z_\phi^\pi - \sigma) V_\phi^{\pi,\sigma}(1) + \gamma \sigma V_{\phi,\min}^{\pi,\sigma} + \gamma (1 - z_\phi^\pi) \left[\gamma (z_\phi^\pi - \sigma) V_\phi^{\pi,\sigma}(1) \right. \\ &\quad \left. + \gamma \sigma V_{\phi,\min}^{\pi,\sigma} + \gamma (1 - z_\phi^\pi) V_\phi^{\pi,\sigma}(0) \right] \\ &\stackrel{(i)}{\leq} \gamma (z_\phi^\pi - \sigma) V_\phi^{\pi,\sigma}(1) + \gamma \sigma V_{\phi,\min}^{\pi,\sigma} + \gamma (1 - z_\phi^\pi) \left[\gamma z_\phi^\pi V_\phi^{\pi,\sigma}(1) + \gamma (1 - z_\phi^\pi) V_\phi^{\pi,\sigma}(0) \right] \\ &\leq \dots \\ &\leq \gamma (z_\phi^\pi - \sigma) V_\phi^{\pi,\sigma}(1) + \gamma \sigma V_{\phi,\min}^{\pi,\sigma} + \gamma z_\phi^\pi \sum_{t=1}^{\infty} \gamma^t (1 - z_\phi^\pi)^t V_\phi^{\pi,\sigma}(1) + \lim_{t \rightarrow \infty} \gamma^t (1 - z_\phi^\pi)^t V_\phi^{\pi,\sigma}(0) \\ &\stackrel{(ii)}{\leq} \gamma (z_\phi^\pi - \sigma) V_\phi^{\pi,\sigma}(1) + \gamma \sigma V_{\phi,\min}^{\pi,\sigma} + \gamma (1 - z_\phi^\pi) \frac{\gamma z_\phi^\pi}{1 - \gamma(1 - z_\phi^\pi)} V_\phi^{\pi,\sigma}(1) + 0 \\ &< \gamma (z_\phi^\pi - \sigma) V_\phi^{\pi,\sigma}(1) + \gamma \sigma V_{\phi,\min}^{\pi,\sigma} + \gamma (1 - z_\phi^\pi) V_\phi^{\pi,\sigma}(1) \\ &= \gamma (1 - \sigma) V_\phi^{\pi,\sigma}(1) + \gamma \sigma V_{\phi,\min}^{\pi,\sigma}, \end{aligned} \quad (153)$$

1041 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
 1042 from the trivial fact that $\frac{\gamma z_\phi^\pi}{1 - \gamma(1 - z_\phi^\pi)} < 1$.

1043 Combining (147), (148), and (153), we have that for any policy π ,

$$V_\phi^{\pi,\sigma}(0) = V_{\phi,\min}^{\pi,\sigma}, \quad (154)$$

1044 which directly leads to

$$V_{\phi}^{\pi,\sigma}(1) = 1 + \gamma(1 - \sigma)V_{\phi}^{\pi,\sigma}(1) + \gamma\sigma V_{\phi,\min}^{\pi,\sigma} = \frac{1 + \gamma\sigma V_{\phi}^{\pi,\sigma}(0)}{1 - \gamma(1 - \sigma)}. \quad (155)$$

1045 Let's now return to the characterization of $V_{\phi}^{\pi,\sigma}(0)$. In view of (154), the equality in (152) holds, and
1046 we have

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

1047 where (i) arises from (155). Solving this relation gives

$$V_{\phi}^{\pi,\sigma}(0) = \frac{\frac{\gamma(z_{\phi}^{\pi} - \sigma)}{1 - \gamma(1 - \sigma)}}{(1 - \gamma)\left(1 + \frac{\gamma(z_{\phi}^{\pi} - \sigma)}{1 - \gamma(1 - \sigma)}\right)}. \quad (156)$$

1048 **The optimal robust policy and optimal robust value function.** We move on to characterize the
1049 robust optimal policy and its corresponding robust value function. To begin with, denoting

$$z := \frac{\gamma(z_{\phi}^{\pi} - \sigma)}{1 - \gamma(1 - \sigma)}, \quad (157)$$

1050 we rewrite (156) as

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

1051 Plugging in the fact that $z_{\phi}^{\pi} \geq q \geq \sigma > 0$ in (130), it follows that $z > 0$. So for any $z > 0$, the
1052 derivative of $f(z)$ w.r.t. z obeys

$$\frac{(1 - \gamma)(1 + z) - (1 - \gamma)z}{(1 - \gamma)^2(1 + z)^2} = \frac{1}{(1 - \gamma)(1 + z)^2} > 0. \quad (158)$$

1053 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
1054 the fact $p \geq q$ in (130)), the optimal policy in state 0 thus obeys

$$\pi_{\phi}^*(\phi | 0) = 1. \quad (159)$$

1055 Considering that the action does not influence the state transition for all states $s > 0$, without loss of
1056 generality, we choose the robust optimal policy to obey

$$\forall s > 0: \quad \pi_{\phi}^*(\phi | s) = 1. \quad (160)$$

1057 Taking $\pi = \pi_{\phi}^*$, we complete the proof by showing that the corresponding robust optimal robust value
1058 function at state 0 as follows:

$$V_{\phi}^{\star,\sigma}(0) = \frac{\frac{\gamma(z_{\phi}^{\pi^*} - \sigma)}{1 - \gamma(1 - \sigma)}}{(1 - \gamma)\left(1 + \frac{\gamma(z_{\phi}^{\pi^*} - \sigma)}{1 - \gamma(1 - \sigma)}\right)} = \frac{\frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)}}{(1 - \gamma)\left(1 + \frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)}\right)}. \quad (161)$$

1059 **D.3.2 Proof of the claim (137)**

1060 Plugging in the definition of φ , we arrive at that for any policy π ,

$$\langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\pi,\sigma} \rangle = V_\phi^{*,\sigma}(0) - V_\phi^{\pi,\sigma}(0) = \frac{\frac{\gamma(p-z_\phi^\pi)}{1-\gamma(1-\sigma)}}{(1-\gamma)\left(1 + \frac{\gamma(p-\sigma)}{1-\gamma(1-\sigma)}\right)\left(1 + \frac{\gamma(z_\phi^\pi-\sigma)}{1-\gamma(1-\sigma)}\right)}, \quad (162)$$

1061 which follows from applying (133) and basic calculus. Then, we proceed to control the above term in
1062 two cases separately in terms of the uncertainty level σ .

1063 • When $\sigma \in (0, 1-\gamma]$. Then regarding the important terms in (162), we observe that

$$1-\gamma < 1-\gamma(1-\sigma) \leq 1-\gamma(1-(1-\gamma)) = (1-\gamma)(1+\gamma) \leq 2(1-\gamma), \quad (163)$$

1064 which directly leads to

$$\frac{\gamma(z_\phi^\pi - \sigma)}{1-\gamma(1-\sigma)} \stackrel{(i)}{\leq} \frac{\gamma(p-\sigma)}{1-\gamma(1-\sigma)} \leq \frac{\gamma c_1(1-\gamma)}{1-\gamma(1-\sigma)} \stackrel{(ii)}{<} c_1\gamma, \quad (164)$$

1065 where (i) holds by $z_\phi^\pi < p$, and (ii) is due to (163). Inserting (163) and (164) back into
1066 (162), we arrive at

$$\begin{aligned} \langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\pi,\sigma} \rangle &\geq \frac{\frac{\gamma(p-z_\phi^\pi)}{2(1-\gamma)}}{(1-\gamma)(1+c_1\gamma)^2} \geq \frac{\gamma(p-z_\phi^\pi)}{8(1-\gamma)^2} \\ &= \frac{\gamma(p-q)(1-\pi(\phi|0))}{8(1-\gamma)^2} = \frac{\gamma\Delta(1-\pi(\phi|0))}{8(1-\gamma)^2} \geq 2\varepsilon(1-\pi(\phi|0)), \end{aligned} \quad (165)$$

1067 where the last inequality holds by setting $(\gamma \geq 1/2)$

$$\Delta = 32(1-\gamma)^2\varepsilon. \quad (166)$$

1068 Finally, it is easily verified that

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

1069 • When $\sigma \in (1-\gamma, 1-c_1]$. Regarding (162), we observe that

$$\gamma\sigma < 1-\gamma(1-\sigma) = 1-\gamma+\gamma\sigma \leq (1+\gamma)\sigma \leq 2\sigma, \quad (167)$$

1070 which directly leads to

$$\frac{\gamma(z_\phi^\pi - \sigma)}{1-\gamma(1-\sigma)} \leq \frac{\gamma(p-\sigma)}{1-\gamma(1-\sigma)} \leq \frac{\gamma c_1\sigma}{1-\gamma(1-\sigma)} \stackrel{(i)}{<} c_1, \quad (168)$$

1071 where (i) holds by (167). Inserting (167) and (168) back into (162), we arrive at

$$\begin{aligned} \langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\pi,\sigma} \rangle &\geq \frac{\frac{\gamma(p-z_\phi^\pi)}{2\sigma}}{(1-\gamma)(1+c_1)^2} \geq \frac{\gamma(p-z_\phi^\pi)}{8(1-\gamma)\sigma} = \frac{\gamma(p-q)(1-\pi(\phi|0))}{8(1-\gamma)\sigma} \\ &= \frac{\gamma\Delta(1-\pi(\phi|0))}{8(1-\gamma)\sigma} \geq 2\varepsilon(1-\pi(\phi|0)), \end{aligned} \quad (169)$$

1072 where the last inequality holds by letting $(\gamma \geq 1/2)$

$$\Delta = 32(1-\gamma)\sigma\varepsilon. \quad (170)$$

1073 Finally, it is easily verified that

$$\varepsilon \leq \frac{c_1}{32(1-\gamma)} \implies \Delta \leq c_1\sigma. \quad (171)$$

1074 **E Proof of the upper bound with χ^2 divergence: Theorem 3**

1075 The proof of Theorem 3 mainly follows the structure of the proof of Theorem 1 in Appendix C.
 1076 Throughout this section, for any nominal transition kernel P , the uncertainty set is taken as (see (8))

$$\begin{aligned} \mathcal{U}^\sigma(P) &= \mathcal{U}_{\chi^2}^\sigma(P) := \otimes \mathcal{U}_{\chi^2}^\sigma(P_{s,a}), \\ \mathcal{U}_{\chi^2}^\sigma(P_{s,a}) &:= \left\{ P'_{s,a} \in \Delta(\mathcal{S}) : \sum_{s' \in \mathcal{S}} \frac{(P'(s' | s, a) - P(s' | s, a))^2}{P(s' | s, a)} \leq \sigma \right\}. \end{aligned} \quad (172)$$

1077 **E.1 Proof of Theorem 3**

1078 In order to control the performance gap $\|V^{*\sigma} - V^{\hat{\pi}, \sigma}\|_\infty$, recall the error decomposition in (49):

$$V^{*\sigma} - V^{\hat{\pi}, \sigma} \leq \left(V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma} \right) + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \mathbf{1} + \left(\hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma} \right), \quad (173)$$

1079 where ε_{opt} (cf. (48)) shall be specified later (which justifies Remark ??). To further control (173), we
 1080 bound the remaining two terms separately.

1081 **Step 1: controlling** $\|\hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma}\|_\infty$. Towards this, recall the bound in (54) which holds for
 1082 any uncertainty set:

$$\begin{aligned} \|\hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma}\|_\infty &\leq \gamma \max \left\{ \left\| \left(I - \gamma \hat{P}^{\pi^*, \hat{V}} \right)^{-1} \left(\hat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \right\|_\infty, \right. \\ &\quad \left. \left\| \left(I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \left(\hat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \right\|_\infty \right\}. \end{aligned} \quad (174)$$

1083 To control the main term $\hat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma}$ in (174), we first introduce an important lemma
 1084 whose proof is postponed to Appendix E.2.1.

1085 **Lemma 15.** Consider any $\sigma > 0$ and the uncertainty set $\mathcal{U}^\sigma(\cdot) := \mathcal{U}_{\chi^2}^\sigma(\cdot)$. For any $\delta \in (0, 1)$ and
 1086 any fixed policy π , one has with probability at least $1 - \delta$,

$$\left\| \hat{P}^{\pi, V} V^{\pi, \sigma} - \underline{P}^{\pi, V} V^{\pi, \sigma} \right\|_\infty \leq 4 \sqrt{\frac{2(1+\sigma) \log\left(\frac{24SAN}{\delta}\right)}{(1-\gamma)^2 N}}.$$

1087 Applying Lemma 15 by taking $\pi = \pi^*$ gives

$$\left\| \hat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right\|_\infty \leq 4 \sqrt{\frac{2(1+\sigma) \log\left(\frac{24SAN}{\delta}\right)}{(1-\gamma)^2 N}}, \quad (175)$$

1088 which directly leads to

$$\begin{aligned} &\left\| \left(I - \gamma \hat{P}^{\pi^*, \hat{V}} \right)^{-1} \left(\hat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \right\|_\infty \\ &\leq \left\| \hat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right\|_\infty \cdot \left\| \left(I - \gamma \hat{P}^{\pi^*, \hat{V}} \right)^{-1} \mathbf{1} \right\|_\infty \leq 4 \sqrt{\frac{2(1+\sigma) \log\left(\frac{24SAN}{\delta}\right)}{(1-\gamma)^4 N}}. \end{aligned} \quad (176)$$

1089 Similarly, we have

$$\left\| \left(I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \left(\hat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \right\|_\infty \leq 4 \sqrt{\frac{2(1+\sigma) \log\left(\frac{24SAN}{\delta}\right)}{(1-\gamma)^4 N}}. \quad (177)$$

1090 Inserting (176) and (177) back to (174) yields

$$\left\| \hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma} \right\|_\infty \leq 4 \sqrt{\frac{2(1+\sigma) \log\left(\frac{24SAN}{\delta}\right)}{(1-\gamma)^4 N}}. \quad (178)$$

1091 **Step 2: controlling** $\left\| \widehat{V}^{\widehat{\pi}, \sigma} - V^{\widehat{\pi}, \sigma} \right\|_{\infty}$. Recall the bound in (55) which holds for any uncertainty
 1092 set:

$$\left\| \widehat{V}^{\widehat{\pi}, \sigma} - V^{\widehat{\pi}, \sigma} \right\|_{\infty} \leq \gamma \max \left\{ \left\| \left(I - \gamma \underline{P}^{\widehat{\pi}, V} \right)^{-1} \left(\widehat{\underline{P}}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - \underline{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right) \right\|_{\infty}, \right. \\ \left. \left\| \left(I - \gamma \underline{P}^{\widehat{\pi}, \widehat{V}} \right)^{-1} \left(\widehat{\underline{P}}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - \underline{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right) \right\|_{\infty} \right\}. \quad (179)$$

1093 We introduce the following lemma which controls $\widehat{\underline{P}}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - \underline{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma}$ in (179); the proof is
 1094 deferred to Appendix E.2.2.

1095 **Lemma 16.** Consider the uncertainty set $\mathcal{U}^{\sigma}(\cdot) := \mathcal{U}_{\chi^2}^{\sigma}(\cdot)$ and any $\delta \in (0, 1)$. With probability at
 1096 least $1 - \delta$, one has

$$\left\| \widehat{\underline{P}}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - \underline{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right\|_{\infty} \leq 12 \sqrt{\frac{2(1+\sigma) \log\left(\frac{36SAN^2}{\delta}\right)}{(1-\gamma)^2 N}} + \frac{2\gamma \varepsilon_{\text{opt}}}{1-\gamma} + 4 \sqrt{\frac{\sigma \varepsilon_{\text{opt}}}{(1-\gamma)^2}}. \quad (180)$$

1097 Repeating the arguments from (175) to (178) yields

$$\left\| \widehat{V}^{\widehat{\pi}, \sigma} - V^{\widehat{\pi}, \sigma} \right\|_{\infty} \leq 12 \sqrt{\frac{2(1+\sigma) \log\left(\frac{36SAN^2}{\delta}\right)}{(1-\gamma)^4 N}} + \frac{2\gamma \varepsilon_{\text{opt}}}{(1-\gamma)^2} + 4 \sqrt{\frac{\sigma \varepsilon_{\text{opt}}}{(1-\gamma)^4}}. \quad (181)$$

1098 Finally, inserting (178) and (181) back to (173) complete the proof

$$\begin{aligned} \left\| V^{*, \sigma} - V^{\widehat{\pi}, \sigma} \right\|_{\infty} &\leq \left\| V^{\pi^*, \sigma} - \widehat{V}^{\pi^*, \sigma} \right\|_{\infty} + \frac{2\gamma \varepsilon_{\text{opt}}}{1-\gamma} + \left\| \widehat{V}^{\widehat{\pi}, \sigma} - V^{\widehat{\pi}, \sigma} \right\|_{\infty} \\ &\leq 4 \sqrt{\frac{2(1+\sigma) \log\left(\frac{24SAN}{\delta}\right)}{(1-\gamma)^4 N}} + \frac{2\gamma \varepsilon_{\text{opt}}}{1-\gamma} + 12 \sqrt{\frac{2(1+\sigma) \log\left(\frac{36SAN^2}{\delta}\right)}{(1-\gamma)^4 N}} \\ &\quad + \frac{2\gamma \varepsilon_{\text{opt}}}{(1-\gamma)^2} + 4 \sqrt{\frac{\sigma \varepsilon_{\text{opt}}}{(1-\gamma)^4}} \\ &\leq 24 \sqrt{\frac{2(1+\sigma) \log\left(\frac{36SAN^2}{\delta}\right)}{(1-\gamma)^4 N}}, \end{aligned} \quad (182)$$

1099 where the last line holds by taking $\varepsilon_{\text{opt}} \leq \min \left\{ \sqrt{\frac{32(1+\sigma) \log\left(\frac{36SAN^2}{\delta}\right)}{N}}, \frac{4 \log\left(\frac{36SAN^2}{\delta}\right)}{N} \right\}$.

1100 E.2 Proof of the auxiliary lemmas

1101 E.2.1 Proof of Lemma 15

1102 **Step 1: controlling the point-wise concentration.** Consider any fixed policy π and the correspond-
 1103 ing robust value vector $V := V^{\pi, \sigma}$ (independent from \widehat{P}^0). Invoking Lemma 5 leads to that for any
 1104 $(s, a) \in \mathcal{S} \times \mathcal{A}$,

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

$$+ \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \sqrt{\sigma} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right|, \quad (183)$$

1105 where the first inequality follows by that the maximum operator is 1-Lipschitz, and the second
 1106 inequality follows from the triangle inequality. Observing that the first term in (183) is exactly the
 1107 same as (89), recalling the fact in (94) directly leads to: with probability at least $1 - \delta$,

$$\max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| (P_{s,a}^0 - \hat{P}_{s,a}^0) [V]_\alpha \right| \leq 2 \sqrt{\frac{\log(\frac{2SAN}{\delta})}{(1-\gamma)^2 N}} \quad (184)$$

1108 holds for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. Then the remainder of the proof focuses on controlling the second term
 1109 in (183).

1110 **Step 2: controlling the second term in (183).** For any given $(s, a) \in \mathcal{S} \times \mathcal{A}$ and fixed $\alpha \in [0, \frac{1}{1-\gamma}]$,
 1111 applying the concentration inequality (Panaganti and Kalathil, 2022, Lemma 6) with $\|[V]_\alpha\|_\infty \leq \frac{1}{1-\gamma}$,
 1112 we arrive at

$$\left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right| \leq \sqrt{\frac{2 \log(\frac{2}{\delta})}{(1-\gamma)^2 N}} \quad (185)$$

1113 holds with probability at least $1 - \delta$. To obtain a uniform bound, we first observe the follow lemma
 1114 proven in Appendix E.2.3.

1115 **Lemma 17.** For any V obeying $\|V\|_\infty \leq \frac{1}{1-\gamma}$, the function $J_{s,a}(\alpha, V) := \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_\alpha)} - \right.$
 1116 $\left. \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right|$ w.r.t. α obeys

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

1117 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,
 1118 2018). Armed with the above, we can derive the uniform bound over $\alpha \in [\min_s V(s), \max_s V(s)] \subset$
 1119 $[0, 1/(1-\gamma)]$: with probability at least $1 - \frac{\delta}{SA}$, it holds that for any $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned} & \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right| \\ & \leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right| \\ & \stackrel{(i)}{\leq} 4 \sqrt{\frac{\varepsilon_3}{1-\gamma}} + \sup_{\alpha \in N_{\varepsilon_3}} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right| \\ & \stackrel{(ii)}{\leq} 4 \sqrt{\frac{\varepsilon_3}{1-\gamma}} + \sqrt{\frac{2 \log(\frac{2SA|N_{\varepsilon_3}|}{\delta})}{(1-\gamma)^2 N}} \\ & \stackrel{(iii)}{\leq} 2 \sqrt{\frac{2 \log(\frac{2SA|N_{\varepsilon_3}|}{\delta})}{(1-\gamma)^2 N}} \leq 2 \sqrt{\frac{2 \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}}, \end{aligned} \quad (186)$$

1120 where (i) holds by the property of N_{ε_3} , (ii) follows from (185), (iii) arises from taking $\varepsilon_3 =$
 1121 $\frac{\log(\frac{2SA|N_{\varepsilon_3}|}{\delta})}{8N(1-\gamma)}$, and the last inequality is verified by $|N_{\varepsilon_3}| \leq \frac{3}{\varepsilon_3(1-\gamma)} \leq 24N$.

1122 Inserting (184) and (186) back to (183) and taking the union bound over $(s, a) \in \mathcal{S} \times \mathcal{A}$, we arrive at
 1123 that for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, with probability at least $1 - \delta$,

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

$$\leq \sqrt{\frac{2 \log(\frac{2SAN}{\delta})}{(1-\gamma)^2 N}} + 2\sqrt{\frac{2\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}} \leq 4\sqrt{\frac{2(1+\sigma) \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}}.$$

1124 Finally, we complete the proof by recalling the matrix form as below:

$$\left\| \widehat{\underline{P}}^{\pi, V} V^{\pi, \sigma} - \underline{P}^{\pi, V} V^{\pi, \sigma} \right\|_{\infty} \leq \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} \left| \widehat{P}_{s, a}^{\pi, V} V - P_{s, a}^{\pi, V} V \right| \leq 4\sqrt{\frac{2(1+\sigma) \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}}.$$

1125 E.2.2 Proof of Lemma 16

1126 **Step 1: decomposing the term of interest.** The proof follows the routine of the proof of Lemma 12
 1127 in Appendix C.3.5. To begin with, for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, following the same arguments of (183)
 1128 yields

$$\begin{aligned} \left| \widehat{P}_{s, a}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - P_{s, a}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right| &\leq \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left| (P_{s, a}^0 - \widehat{P}_{s, a}^0) \left[\widehat{V}^{\widehat{\pi}, \sigma} \right]_{\alpha} \right| + \\ &+ \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \sqrt{\sigma} \left| \sqrt{\text{Var}_{\widehat{P}_{s, a}^0} \left(\left[\widehat{V}^{\widehat{\pi}, \sigma} \right]_{\alpha} \right)} - \sqrt{\text{Var}_{P_{s, a}^0} \left(\left[\widehat{V}^{\widehat{\pi}, \sigma} \right]_{\alpha} \right)} \right|. \end{aligned} \quad (187)$$

1129 Invoking the fact in (118) (for proving Lemma 12), the first term in (187) obeys

$$\begin{aligned} \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi}, \sigma}(s), \max_s \widehat{V}^{\widehat{\pi}, \sigma}(s)]} \left| (P_{s, a}^0 - \widehat{P}_{s, a}^0) \left[\widehat{V}^{\widehat{\pi}, \sigma} \right]_{\alpha} \right| &\leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| (P_{s, a}^0 - \widehat{P}_{s, a}^0) \left[\widehat{V}^{\widehat{\pi}, \sigma} \right]_{\alpha} \right| \\ &\leq 4\sqrt{\frac{\log(\frac{3SAN^{3/2}}{(1-\gamma)\delta})}{(1-\gamma)^2 N}} + \frac{2\gamma \varepsilon_{\text{opt}}}{1-\gamma}. \end{aligned} \quad (188)$$

1130 The remainder of the proof will focus on controlling the second term of (187).

1131 **Step 2: controlling the second term of (187).** Towards this, we recall the auxiliary robust MDP
 1132 $\widehat{\mathcal{M}}_{\text{rob}}^{s, u}$ defined in Appendix C.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
 1133 $\widehat{\mathcal{M}}_{\text{rob}}$, we recall the corresponding robust Bellman operator $\widehat{\mathcal{T}}_{s, u}^{\sigma}(\cdot)$ in (107) and the following
 1134 definition in (108)

$$u^* := \widehat{V}^{*, \sigma}(s) - \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(e_s)} \mathcal{P} \widehat{V}^{*, \sigma}. \quad (189)$$

1135 Following the arguments in Appendix C.3.5, it can be verified that there exists a unique fixed point
 1136 $\widehat{Q}_{s, u}^{*, \sigma}$ of the operator $\widehat{\mathcal{T}}_{s, u}^{\sigma}(\cdot)$, which satisfies $0 \leq \widehat{Q}_{s, u}^{*, \sigma} \leq \frac{1}{1-\gamma} \mathbf{1}$. In addition, the corresponding robust
 1137 value function coincides with that of the operator $\widehat{\mathcal{T}}^{\sigma}(\cdot)$, i.e., $\widehat{V}_{s, u}^{*, \sigma} = \widehat{V}^{*, \sigma}$.

1138 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
 1139 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
 1140 independence between $[\widehat{V}_{s, u}^{*, \sigma}]_{\alpha}$ and $\widehat{P}_{s, a}^0$. With this in mind, invoking the fact in (186) and taking the
 1141 union bound over all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and $u \in N_{\varepsilon_2}$ yields that, with probability at least $1 - \delta$,

$$\max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s, a}^0} \left(\left[\widehat{V}_{s, u}^{*, \sigma} \right]_{\alpha} \right)} - \sqrt{\text{Var}_{P_{s, a}^0} \left(\left[\widehat{V}_{s, u}^{*, \sigma} \right]_{\alpha} \right)} \right| \leq 2\sqrt{\frac{2 \log(\frac{24SAN|N_{\varepsilon_2}|}{\delta})}{(1-\gamma)^2 N}} \quad (190)$$

1142 holds for all $(s, a, u) \in \mathcal{S} \times \mathcal{A} \times N_{\varepsilon_2}$.

1143 To continue, we decompose the term of interest in (187) as follows:

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

$$\begin{aligned}
&\leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0} \left(\left[\hat{V}^{\hat{\pi}, \sigma} \right]_{\alpha} \right)} - \sqrt{\text{Var}_{P_{s,a}^0} \left(\left[\hat{V}^{\hat{\pi}, \sigma} \right]_{\alpha} \right)} \right| \\
&\stackrel{(i)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0} \left(\left[\hat{V}^{*, \sigma} \right]_{\alpha} \right)} - \sqrt{\text{Var}_{P_{s,a}^0} \left(\left[\hat{V}^{*, \sigma} \right]_{\alpha} \right)} \right| \\
&\quad + \max_{\alpha \in [0, 1/(1-\gamma)]} \left[\sqrt{\left| \text{Var}_{\hat{P}_{s,a}^0} \left(\left[\hat{V}^{\hat{\pi}, \sigma} \right]_{\alpha} \right) - \text{Var}_{\hat{P}_{s,a}^0} \left(\left[\hat{V}^{*, \sigma} \right]_{\alpha} \right) \right|} \right. \\
&\quad \left. + \sqrt{\left| \text{Var}_{P_{s,a}^0} \left(\left[\hat{V}^{\hat{\pi}, \sigma} \right]_{\alpha} \right) - \text{Var}_{P_{s,a}^0} \left(\left[\hat{V}^{*, \sigma} \right]_{\alpha} \right) \right|} \right] \\
&\stackrel{(ii)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0} \left(\left[\hat{V}^{*, \sigma} \right]_{\alpha} \right)} - \sqrt{\text{Var}_{P_{s,a}^0} \left(\left[\hat{V}^{*, \sigma} \right]_{\alpha} \right)} \right| \\
&\quad + \max_{\alpha \in [0, 1/(1-\gamma)]} 2\sqrt{\frac{2}{(1-\gamma)}} \left\| \left[\hat{V}^{\hat{\pi}, \sigma} \right]_{\alpha} - \left[\hat{V}^{*, \sigma} \right]_{\alpha} \right\|_{\infty} \\
&\leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0} \left(\left[\hat{V}^{*, \sigma} \right]_{\alpha} \right)} - \sqrt{\text{Var}_{P_{s,a}^0} \left(\left[\hat{V}^{*, \sigma} \right]_{\alpha} \right)} \right| + 4\sqrt{\frac{\varepsilon_{\text{opt}}}{(1-\gamma)^2}}, \quad (191)
\end{aligned}$$

1144 where (i) holds by the triangle inequality, (ii) arises from applying Lemma 2, and the last inequality
1145 holds by (48).

1146 Armed with the above facts, invoking the identity $\hat{V}^{*, \sigma} = \hat{V}_{s, u^*}^{*, \sigma}$ leads to that for all $(s, a) \in \mathcal{S} \times \mathcal{A}$,
1147 with probability at least $1 - \delta$,

$$\begin{aligned}
&\max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0} \left(\left[\hat{V}^{*, \sigma} \right]_{\alpha} \right)} - \sqrt{\text{Var}_{P_{s,a}^0} \left(\left[\hat{V}^{*, \sigma} \right]_{\alpha} \right)} \right| \\
&= \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0} \left(\left[\hat{V}_{s, u^*}^{*, \sigma} \right]_{\alpha} \right)} - \sqrt{\text{Var}_{P_{s,a}^0} \left(\left[\hat{V}_{s, u^*}^{*, \sigma} \right]_{\alpha} \right)} \right| \\
&\stackrel{(i)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0} \left(\left[\hat{V}_{s, \bar{u}}^{*, \sigma} \right]_{\alpha} \right)} - \sqrt{\text{Var}_{P_{s,a}^0} \left(\left[\hat{V}_{s, \bar{u}}^{*, \sigma} \right]_{\alpha} \right)} \right| \\
&\quad + \max_{\alpha \in [0, 1/(1-\gamma)]} \left[\sqrt{\left| \text{Var}_{\hat{P}_{s,a}^0} \left(\left[\hat{V}_{s, u^*}^{*, \sigma} \right]_{\alpha} \right) - \text{Var}_{\hat{P}_{s,a}^0} \left(\left[\hat{V}_{s, \bar{u}}^{*, \sigma} \right]_{\alpha} \right) \right|} \right. \\
&\quad \left. + \sqrt{\left| \text{Var}_{P_{s,a}^0} \left(\left[\hat{V}_{s, u^*}^{*, \sigma} \right]_{\alpha} \right) - \text{Var}_{P_{s,a}^0} \left(\left[\hat{V}_{s, \bar{u}}^{*, \sigma} \right]_{\alpha} \right) \right|} \right] \\
&\stackrel{(ii)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0} \left(\left[\hat{V}_{s, \bar{u}}^{*, \sigma} \right]_{\alpha} \right)} - \sqrt{\text{Var}_{P_{s,a}^0} \left(\left[\hat{V}_{s, \bar{u}}^{*, \sigma} \right]_{\alpha} \right)} \right| + 4\sqrt{\frac{\varepsilon_2}{(1-\gamma)}} \\
&\stackrel{(iii)}{\leq} 2\sqrt{\frac{2 \log\left(\frac{24SAN|N_{\varepsilon_2}|}{\delta}\right)}{(1-\gamma)^2 N}} + 4\sqrt{\frac{\varepsilon_2}{(1-\gamma)}} \\
&\leq 6\sqrt{\frac{2 \log\left(\frac{36SAN^2|N_{\varepsilon_2}|}{\delta}\right)}{(1-\gamma)^2 N}}, \quad (192)
\end{aligned}$$

1148 where (i) holds by the triangle inequality, (ii) arises from applying Lemma 2 and the fact
1149 $\left\| \hat{V}_{s, \bar{u}}^{*, \sigma} - \hat{V}_{s, u^*}^{*, \sigma} \right\|_{\infty} \leq \frac{\varepsilon_2}{(1-\gamma)}$ (see (114)), (iii) follows from (190), and the last inequality holds
1150 by letting $\varepsilon_2 = \frac{2 \log\left(\frac{24SAN|N_{\varepsilon_2}|}{\delta}\right)}{(1-\gamma)N}$, which leads to $|N_{\varepsilon_2}| \leq \frac{3}{\varepsilon_2(1-\gamma)} \leq \frac{3N}{2}$.

1151 In summary, inserting (192) back to (191) and (191) leads to with probability at least $1 - \delta$,

$$\max_{\alpha \in [\min_s \hat{V}^{\hat{\pi}, \sigma}(s), \max_s \hat{V}^{\hat{\pi}, \sigma}(s)]} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0} \left(\left[\hat{V}^{\hat{\pi}, \sigma} \right]_{\alpha} \right)} - \sqrt{\text{Var}_{P_{s,a}^0} \left(\left[\hat{V}^{\hat{\pi}, \sigma} \right]_{\alpha} \right)} \right|$$

$$\leq 6\sqrt{\frac{2\sigma \log\left(\frac{36SAN^2|N_{\varepsilon_2}|}{\delta}\right)}{(1-\gamma)^2N}} + 4\sqrt{\frac{\sigma\varepsilon_{\text{opt}}}{(1-\gamma)^2}} \quad (193)$$

1152 holds for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.

1153 **Step 4: finishing up.** Inserting (193) and (188) back to (187), we complete the proof: with
1154 probability at least $1 - \delta$,

$$\begin{aligned} & \left\| \widehat{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - P^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} \right\|_{\infty} \\ & \leq 4\sqrt{\frac{\log\left(\frac{3SAN^{3/2}}{(1-\gamma)\delta}\right)}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} + 6\sqrt{\frac{2\sigma \log\left(\frac{36SAN^2|N_{\varepsilon_2}|}{\delta}\right)}{(1-\gamma)^2N}} + 4\sqrt{\frac{\sigma\varepsilon_{\text{opt}}}{(1-\gamma)^2}} \\ & \leq 12\sqrt{\frac{2(1+\sigma) \log\left(\frac{36SAN^2}{\delta}\right)}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} + 4\sqrt{\frac{\sigma\varepsilon_{\text{opt}}}{(1-\gamma)^2}}. \end{aligned} \quad (194)$$

1155 E.2.3 Proof of Lemma 17

1156 For any $0 \leq \alpha_1, \alpha_2 \leq 1/(1-\gamma)$, one has

$$\begin{aligned} & |J_{s,a}(\alpha_1, V) - J_{s,a}(\alpha_2, V)| \\ & = \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_1})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_1})} \right| - \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_2})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_2})} \right| \\ & \stackrel{(i)}{\leq} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_1})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_1})} \right| + \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_2})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_2})} \right| \\ & \leq \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_1})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_1})} \right| + \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_2})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_2})} \right| \\ & \stackrel{(ii)}{\leq} \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_2}) - \text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_1})} + \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_2}) - \text{Var}_{P_{s,a}^0}([V]_{\alpha_1})} \\ & \stackrel{(iii)}{\leq} \sqrt{\left| \widehat{P}_{s,a}^0([V]_{\alpha_1}) \circ ([V]_{\alpha_1}) - ([V]_{\alpha_2}) \circ ([V]_{\alpha_2}) \right|} + \left| \widehat{P}_{s,a}^0([V]_{\alpha_1} + [V]_{\alpha_2}) \cdot \widehat{P}_{s,a}^0([V]_{\alpha_1} - [V]_{\alpha_2}) \right| \\ & \quad + \sqrt{\left| P_{s,a}^0([V]_{\alpha_1}) \circ ([V]_{\alpha_1}) - ([V]_{\alpha_2}) \circ ([V]_{\alpha_2}) \right|} + \left| P_{s,a}^0([V]_{\alpha_1} + [V]_{\alpha_2}) \cdot P_{s,a}^0([V]_{\alpha_1} - [V]_{\alpha_2}) \right| \\ & \leq 2\sqrt{2(\alpha_1 + \alpha_2)|\alpha_1 - \alpha_2|} \leq 4\sqrt{\frac{|\alpha_1 - \alpha_2|}{1-\gamma}}. \end{aligned} \quad (195)$$

1157 where (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
1158 $\sqrt{x} - \sqrt{y} \leq \sqrt{x - y}$ 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
1159 $P \in \Delta(\mathcal{S})$, (iii) holds by the definition of $\text{Var}_P(\cdot)$ defined in (24), and the last inequality arises from
1160 $0 \leq \alpha_1, \alpha_2 \leq 1/(1-\gamma)$.

1161 F Proof of the lower bound with χ^2 divergence: Theorem 4

1162 To prove Theorem 4, we shall first construct some hard instances and then characterize the sample
1163 complexity requirements over these instances. The structure of the hard instances are the same as the
1164 ones used in the proof of Theorem 2.

1165 F.1 Construction of the hard problem instances

1166 First, note that we shall use the same MDPs defined in Appendix D.1 as follows

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

1167 In particular, we shall keep the structure of the transition kernel in (123), reward function in (125)
1168 and initial state distribution in (126), while p and Δ shall be specified differently later.

1169 **Uncertainty set of the transition kernels.** Recalling the uncertainty set associated with χ^2 diver-
 1170 gence in (172), for any uncertainty level σ , the uncertainty set throughout this section is defined as
 1171 $\mathcal{U}^\sigma(P^\phi)$:

$$\begin{aligned} \mathcal{U}^\sigma(P^\phi) &:= \otimes \mathcal{U}_{\chi^2}^\sigma(P_{s,a}^\phi), \\ \mathcal{U}_{\chi^2}^\sigma(P_{s,a}^\phi) &:= \left\{ P_{s,a} \in \Delta(\mathcal{S}) : \sum_{s' \in \mathcal{S}} \frac{(P(s' | s, a) - P^\phi(s' | s, a))^2}{P^\phi(s' | s, a)} \leq \sigma \right\}. \end{aligned} \quad (196)$$

1172 Clearly, $\mathcal{U}^\sigma(P_{s,a}^\phi) = P_{s,a}^\phi$ whenever the state transition is deterministic for χ^2 divergence. Here,
 1173 q and Δ (whose choice will be specified later in more detail) which determine the instances are
 1174 specified as

$$0 \leq q = \begin{cases} 1 - \gamma & \text{if } \sigma \in (0, \frac{1-\gamma}{4}) \\ \frac{\sigma}{1+\sigma} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty) \end{cases}, \quad p = q + \Delta, \quad (197)$$

1175 and

$$0 < \Delta \leq \begin{cases} \frac{1}{4}(1 - \gamma) & \text{if } \sigma \in (0, \frac{1-\gamma}{4}) \\ \min \left\{ \frac{1}{4}(1 - \gamma), \frac{1}{2(1+\sigma)} \right\} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty) \end{cases}. \quad (198)$$

This directly ensures that

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

1176 since $\gamma \in [\frac{3}{4}, 1)$.

1177 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
 1178 next state s' associated with any perturbed transition kernel $P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^\phi)$ as

$$\underline{P}^\phi(s' | s, a) := \inf_{P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^\phi)} P(s' | s, a). \quad (199)$$

1179 In addition, we denote the transition from state 0 to state 1 as follows, which plays an important role
 1180 in the analysis,

$$\underline{p} := \underline{P}^\phi(1 | 0, \phi), \quad \underline{q} := \underline{P}^\phi(1 | 0, 1 - \phi). \quad (200)$$

1181 Before continuing, we introduce some facts about \underline{p} and \underline{q} which are summarized as the following
 1182 lemma; the proof is postponed to Appendix F.3.1.

1183 **Lemma 18.** Consider any $\sigma \in (0, \infty)$ and any p, q, Δ obeying (197) and (198), the following
 1184 properties hold

$$\begin{cases} \frac{1-\gamma}{2} < \underline{q} < 1 - \gamma, & \underline{q} + \frac{3}{4}\Delta \leq \underline{p} \leq \underline{q} + \Delta \leq \frac{5(1-\gamma)}{4} & \text{if } \sigma \in (0, \frac{1-\gamma}{4}), \\ \underline{q} = 0, & \frac{\sigma+1}{2}\Delta \leq \underline{p} \leq (3 + \sigma)\Delta & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty). \end{cases} \quad (201)$$

1185 **Value functions and optimal policies.** Armed with above facts, we are positioned to derive the
 1186 corresponding robust value functions, the optimal policies, and its corresponding optimal robust
 1187 value functions. For any RMDP \mathcal{M}_ϕ with the uncertainty set defined in (196), we denote the robust
 1188 optimal policy as π_ϕ^* , the robust value function of any policy π (resp. the optimal policy π_ϕ^*) as
 1189 $V_\phi^{\pi, \sigma}$ (resp. $V_\phi^{*, \sigma}$). The following lemma describes some key properties of the robust (optimal) value
 1190 functions and optimal policies whose proof is postponed to Appendix F.3.2.

1191 **Lemma 19.** For any $\phi = \{0, 1\}$ and any policy π , one has

$$V_\phi^{\pi, \sigma}(0) = \frac{\gamma z_\phi^\pi}{(1 - \gamma) \left(1 - \gamma(1 - z_\phi^\pi) \right)}, \quad (202)$$

1192 where z_ϕ^π is defined as

$$z_\phi^\pi := \underline{p}\pi(\phi | 0) + \underline{q}\pi(1 - \phi | 0). \quad (203)$$

1193 In addition, the optimal value functions and the optimal policies obey

$$V_\phi^{*, \sigma}(0) = \frac{\gamma \underline{p}}{(1 - \gamma) (1 - \gamma(1 - \underline{p}))}, \quad (204a)$$

$$\pi_\phi^*(\phi | s) = 1, \quad \text{for } s \in \mathcal{S}. \quad (204b)$$

1194 **F.2 Establishing the minimax lower bound**

1195 Our goal is to control the performance gap w.r.t. any policy estimator $\hat{\pi}$ based on the generated
1196 dataset and the chosen initial distribution φ in (126), which gives

$$\langle \varphi, V_{\phi}^{*,\sigma} - V_{\phi}^{\hat{\pi},\sigma} \rangle = V_{\phi}^{*,\sigma}(0) - V_{\phi}^{\hat{\pi},\sigma}(0). \quad (205)$$

1197 **Step 1: converting the goal to estimate ϕ .** To achieve the goal, we first introduce the following
1198 fact which shall be verified in Appendix F.3.3: given

$$\varepsilon \leq \begin{cases} \frac{1}{72(1-\gamma)} & \text{if } \sigma \in (0, \frac{1-\gamma}{4}), \\ \frac{1}{256(1+\sigma)(1-\gamma)} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \frac{1}{3(1-\gamma)}), \\ \frac{3}{32} & \text{if } \sigma > \frac{1}{3(1-\gamma)}. \end{cases} \quad (206)$$

1199 choosing

$$\Delta = \begin{cases} 18(1-\gamma)^2\varepsilon & \text{if } \sigma \in (0, \frac{1-\gamma}{4}), \\ 64(1+\sigma)(1-\gamma)^2\varepsilon & \text{if } \sigma \in [\frac{1-\gamma}{4}, \frac{1}{3(1-\gamma)}), \\ \frac{16}{3(1+\sigma)}\varepsilon & \text{if } \sigma > \frac{1}{3(1-\gamma)}. \end{cases} \quad (207)$$

1200 which satisfies the requirement of Δ in (197), it holds that for any policy $\hat{\pi}$,

$$\langle \varphi, V_{\phi}^{*,\sigma} - V_{\phi}^{\hat{\pi},\sigma} \rangle \geq 2\varepsilon(1 - \hat{\pi}(\phi|0)). \quad (208)$$

1201 **Step 2: arriving at the final results.** To continue, following the same definitions and argument in
1202 Appendix D.2, we recall the minimax probability of the error and its property as follows:

$$p_e \geq \frac{1}{4} \exp \left\{ -N \left(\text{KL}(P^0(\cdot|0,0) \| P^1(\cdot|0,0)) + \text{KL}(P^0(\cdot|0,1) \| P^1(\cdot|0,1)) \right) \right\}, \quad (209)$$

1203 then we can complete the proof by showing $p_e \geq \frac{1}{8}$ given the bound for the sample size N . In the
1204 following, we shall control the KL divergence terms in (209) in three different cases.

1205 • Case 1: $\sigma \in (0, \frac{1-\gamma}{4})$. In this case, applying $\gamma \in [\frac{3}{4}, 1)$ yields

$$\begin{aligned} 1-q &> 1-p = 1-q - \Delta > \gamma - \frac{1-\gamma}{4} > \frac{3}{4} - \frac{1}{16} > \frac{1}{2}, \\ p &\geq q = 1-\gamma. \end{aligned} \quad (210)$$

1206 Armed with the above facts, applying Lemma 1 (cf. (23)) yields

$$\begin{aligned} \text{KL}(P^0(\cdot|0,1) \| P^1(\cdot|0,1)) &= \text{KL}(p \| q) \leq \frac{(p-q)^2}{(1-p)p} \stackrel{(i)}{=} \frac{\Delta^2}{p(1-p)} \\ &\stackrel{(ii)}{=} \frac{324(1-\gamma)^4\varepsilon^2}{p(1-p)} \\ &\stackrel{(iii)}{\leq} 648(1-\gamma)^3\varepsilon^2, \end{aligned} \quad (211)$$

1207 where (i) follows from the definition in (197), (ii) holds by plugging in the expression of Δ in
1208 (207), and (iii) arises from (210). The same bound can be established for $\text{KL}(P_1^0(\cdot|0,0) \|$
1209 $P_1^1(\cdot|0,0))$. Substituting (211) back into (209) demonstrates that: if the sample size is
1210 chosen as

$$N \leq \frac{\log 2}{1296(1-\gamma)^3\varepsilon^2}, \quad (212)$$

1211 then one necessarily has

$$p_e \geq \frac{1}{4} \exp \left\{ -N \cdot 1296(1-\gamma)^3\varepsilon^2 \right\} \geq \frac{1}{8}. \quad (213)$$

1212

- Case 2: $\sigma \in \left[\frac{1-\gamma}{4}, \frac{1}{3(1-\gamma)} \right)$. Applying the facts of Δ in (198), one has

$$\begin{aligned} 1 - q > 1 - p = 1 - q - \Delta &\geq \frac{1}{1 + \sigma} - \frac{1}{2(1 + \sigma)} = \frac{1}{2(1 + \sigma)}, \\ p \geq q &= \frac{\sigma}{1 + \sigma}. \end{aligned} \quad (214)$$

1213

Given (214), applying Lemma 1 (cf. (23)) yields

$$\begin{aligned} \text{KL}(P^0(\cdot | 0, 1) \| P^1(\cdot | 0, 1)) &= \text{KL}(p \| q) \leq \frac{(p - q)^2}{(1 - p)p} \stackrel{(i)}{=} \frac{\Delta^2}{p(1 - p)} \\ &\stackrel{(ii)}{\leq} \frac{4096(1 + \sigma)^2(1 - \gamma)^4 \varepsilon^2}{p(1 - p)} \\ &\stackrel{(iii)}{\leq} \frac{4096(1 + \sigma)^2(1 - \gamma)^4 \varepsilon^2}{\frac{\sigma}{2(1 + \sigma)^2}} \leq \frac{8192(1 - \gamma)^4(1 + \sigma)^4 \varepsilon^2}{\sigma}, \end{aligned} \quad (215)$$

1214

where (i) follows from the definition in (197), (ii) holds by plugging in the expression of Δ in (207), and (iii) arises from (214). The same bound can be established for $\text{KL}(P_1^0(\cdot | 0, 0) \| P_1^1(\cdot | 0, 0))$.

1215

1216

1217

Substituting (215) back into (142) demonstrates that: if the sample size is chosen as

$$N \leq \frac{\sigma \log 2}{16384(1 - \gamma)^4(1 + \sigma)^4 \varepsilon^2}, \quad (216)$$

1218

then one necessarily has

$$p_e \geq \frac{1}{4} \exp \left\{ -N \frac{16384(1 - \gamma)^4(1 + \sigma)^4 \varepsilon^2}{\sigma} \right\} \geq \frac{1}{8}. \quad (217)$$

1219

- Case 3: $\sigma > \frac{1}{3(1-\gamma)} \geq \frac{1}{3}$. Regarding this, one gives

$$\begin{aligned} 1 - q > 1 - p = 1 - q - \Delta &\geq \frac{1}{1 + \sigma} - \frac{1}{4(1 + \sigma)} \geq \frac{1}{2(1 + \sigma)}, \\ p \geq q &\geq \frac{1}{4}. \end{aligned} \quad (218)$$

1220

Given $p \geq q \geq 1/2$ and (218), applying Lemma 1 (cf. (23)) yields

$$\begin{aligned} \text{KL}(P^0(\cdot | 0, 1) \| P^1(\cdot | 0, 1)) &= \text{KL}(p \| q) \leq \frac{(p - q)^2}{(1 - p)p} \stackrel{(i)}{=} \frac{\Delta^2}{p(1 - p)} \\ &\stackrel{(ii)}{\leq} \frac{\frac{64}{(1 + \sigma)^2} \varepsilon^2}{p(1 - p)} \\ &\stackrel{(iii)}{\leq} \frac{492 \varepsilon^2}{\sigma}, \end{aligned} \quad (219)$$

1221

1222

1223

1224

where (i) follows from the definition in (197), (ii) holds by plugging in the expression of Δ in (207), and (iii) arises from (218). The same bound can be established for $\text{KL}(P_1^0(\cdot | 0, 0) \| P_1^1(\cdot | 0, 0))$. Substituting (219) back into (142) demonstrates that: if the sample size is chosen as

$$N \leq \frac{\sigma \log 2}{984 \varepsilon^2}, \quad (220)$$

1225

then one necessarily has

$$p_e \geq \frac{1}{4} \exp \left\{ -N \frac{984 \varepsilon^2}{\sigma} \right\} \geq \frac{1}{8}. \quad (221)$$

1226 **Step 3: putting things together.** Finally, summing up the results in (212), (216), and (220),
 1227 combined with the requirement in (206), one has when

$$\varepsilon \leq c_1 \begin{cases} \frac{1}{1-\gamma} & \text{if } \sigma \in (0, \frac{1-\gamma}{4}) \\ \max \left\{ \frac{1}{(1+\sigma)(1-\gamma)}, 1 \right\} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty) \end{cases}, \quad (222)$$

1228 taking

$$N \leq c_2 \begin{cases} \frac{1}{(1-\gamma)^3 \varepsilon^2} & \text{if } \sigma \in (0, \frac{1-\gamma}{4}) \\ \frac{\sigma}{\min\{1, (1-\gamma)^4 (1+\sigma)^4\} \varepsilon^2} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty) \end{cases} \quad (223)$$

1229 leads to $p_e \geq \frac{1}{8}$, for some universal constants $c_1, c_2 > 0$.

1230 F.3 Proof of the auxiliary facts

1231 We begin with some basic facts about the χ^2 divergence defined in (22) for any two Bernoulli
 1232 distributions $\text{Ber}(w)$ and $\text{Ber}(x)$, denoted as

$$f(w, x) := \chi^2(x \parallel w) = \frac{(w-x)^2}{w} + \frac{(1-w-(1-x))^2}{1-w} = \frac{(w-x)^2}{w(1-w)}. \quad (224)$$

1233 For $x \in [0, w)$, it is easily verified that the partial derivative w.r.t. x obeys $\frac{\partial f(w, x)}{\partial x} = \frac{2(x-w)}{w(1-w)} < 0$,
 1234 implying that

$$\forall x_1 < x_2 \in [0, w), \quad f(w, x_1) > f(w, x_2). \quad (225)$$

1235 In other words, the χ^2 divergence $f(w, x)$ increases as x decreases from w to 0.

1236 Next, we introduce the following function for any fixed $\sigma \in (0, \infty)$ and any $x \in [\frac{\sigma}{1+\sigma}, 1)$:

$$f_\sigma(x) := \inf_{\{y: \chi^2(y \parallel x) \leq \sigma, y \in [0, x]\}} y \stackrel{(i)}{=} \max \left\{ 0, x - \sqrt{\sigma x(1-x)} \right\} = x - \sqrt{\sigma x(1-x)}, \quad (226)$$

1237 where (i) has been verified in Yang et al. (2022, Corollary B.2), and the last equality holds since
 1238 $x \geq \frac{\sigma}{1+\sigma}$. The next lemma summarizes some useful facts about $f_\sigma(\cdot)$, which again has been verified
 1239 in Yang et al. (2022, Lemma B.12 and Corollary B.2).

1240 **Lemma 20.** Consider any $\sigma \in (0, \infty)$. For $x \in [\frac{\sigma}{1+\sigma}, 1)$, $f_\sigma(x)$ is convex and differentiable, which
 1241 obeys

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

1242 F.3.1 Proof of Lemma 18

1243 Let us control \underline{q} and \underline{p} respectively.

1244 **Step 1: controlling \underline{q} .** We shall control \underline{q} in different cases w.r.t. the uncertainty level σ .

1245 • Case 1: $\sigma \in (0, \frac{1-\gamma}{4})$. In this case, recall that $q = 1 - \gamma$ defined in (197), applying (226)
 1246 with $x = q$ leads to

$$1 - \gamma = q > \underline{q} = f_\sigma(q) = 1 - \gamma - \sqrt{\sigma \gamma (1 - \gamma)} \geq 1 - \gamma - \sqrt{\frac{1-\gamma}{4} \gamma (1 - \gamma)} > \frac{1-\gamma}{2}. \quad (227)$$

1247 • Case 2: $\sigma \in [\frac{1-\gamma}{4}, \infty)$. Note that it suffices to treat $P_{0,1-\phi}^\phi$ as a Bernoulli distribution $\text{Ber}(q)$
 1248 over states 1 and 0, since we do not allow transition to other states. Recalling $q = \frac{\sigma}{1+\sigma}$ in
 1249 (197) and noticing the fact that

$$f(q, 0) = \frac{q^2}{q} + \frac{(1-(1-q))^2}{1-q} = \frac{q}{(1-q)} = \sigma, \quad (228)$$

1250 one has the probability $\text{Ber}(0)$ falls into the uncertainty set of $\text{Ber}(q)$ of size σ . As a result,
 1251 recalling the definition (200) leads to

$$\underline{q} = P^\phi(1 | 0, 1 - \phi) = 0, \quad (229)$$

1252 since $\underline{q} \geq 0$.

1253 **Step 2: controlling \underline{p} .** To characterize the value of \underline{p} , we also divide into several cases separately.

1254 • 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
1255 that $f_\sigma(\cdot)$ is convex and the form of its derivative in Lemma 20, one has

$$\begin{aligned} \underline{p} &= f_\sigma(p) \geq f_\sigma(q) + f'_\sigma(q)(p - q) \\ &= \underline{q} + \left(1 + \frac{\sqrt{\sigma}(2q - 1)}{2\sqrt{q(1 - q)}}\right) \Delta \geq \underline{q} + \left(1 - \frac{\sqrt{\frac{1-\gamma}{4}}(1 - 2(1 - \gamma))}{2\sqrt{(1 - \gamma)\gamma}}\right) \Delta \geq \underline{q} + \frac{3\Delta}{4}. \end{aligned} \quad (230)$$

1256 Similarly, applying Lemma 20 leads to

$$\begin{aligned} \underline{p} &= f_\sigma(p) \leq f_\sigma(q) + f'_\sigma(p)(p - q) \\ &= \underline{q} + \left(1 - \frac{\sqrt{\sigma}(1 - 2p)}{2\sqrt{p(1 - p)}}\right) \Delta \leq \underline{q} + \Delta, \end{aligned} \quad (231)$$

1257 where the last inequality holds by $1 - 2p > 0$ due to the fact $p = q + \Delta \leq \frac{5}{4}(1 - \gamma) \leq \frac{5}{16} < \frac{1}{2}$
1258 (cf. (198) and $\gamma \in [\frac{3}{4}, 1)$). To sum up, given $\sigma \in (0, \frac{1-\gamma}{4})$, combined with (227), we arrive at
1259

$$\underline{q} + \frac{3}{4}\Delta \leq \underline{p} \leq \underline{q} + \Delta \leq \frac{5(1 - \gamma)}{4}, \quad (232)$$

1260 where the last inequality holds by $\Delta \leq \frac{1}{4}(1 - \gamma)$ (see (197)).

1261 • Case 2: $\sigma \in [\frac{1-\gamma}{4}, \infty)$. We recall that $p = q + \Delta > q = \frac{\sigma}{1+\sigma}$ in (197). To derive the lower
1262 bound for \underline{p} in (200), similar to (230), one has

$$\begin{aligned} \underline{p} &= f_\sigma(p) \geq f_\sigma(q) + f'_\sigma(q)(p - q) \\ &= \underline{q} + \left(1 + \frac{\sqrt{\sigma}(2q - 1)}{2\sqrt{q(1 - q)}}\right) \Delta \\ &\stackrel{(i)}{=} 0 + \left(1 + \frac{\sqrt{\sigma}\frac{\sigma-1}{1+\sigma}}{2\sqrt{\frac{\sigma}{1+\sigma}\frac{1}{1+\sigma}}}\right) \Delta = \left(1 + \frac{\sigma - 1}{2}\right) \Delta = \left(\frac{\sigma + 1}{2}\right) \Delta, \end{aligned} \quad (233)$$

1263 where (i) follows from $q = \frac{\sigma}{1+\sigma}$ and $\underline{q} = 0$ (see (229)). For the other direction, similar to
1264 (231), we have

$$\begin{aligned} \underline{p} &= f_\sigma(p) \leq f_\sigma(q) + f'_\sigma(p)(p - q) = \underline{q} + \left(1 + \frac{\sqrt{\sigma}(2p - 1)}{2\sqrt{p(1 - p)}}\right) \Delta \\ &\stackrel{(i)}{=} \left(1 + \frac{\sqrt{\sigma}(2p - 1)}{2\sqrt{p(1 - p)}}\right) \Delta \stackrel{(ii)}{=} \left(1 + \frac{\sqrt{\sigma}\left(\frac{\sigma-1}{1+\sigma} + 2\Delta\right)}{2\sqrt{\left(\frac{\sigma}{1+\sigma} + \Delta\right)\left(\frac{1}{1+\sigma} - \Delta\right)}}\right) \Delta \\ &\stackrel{(iii)}{\leq} \left(1 + \frac{\sqrt{\sigma}(1 + 2\Delta)}{2\sqrt{\frac{\sigma}{1+\sigma} \cdot \frac{1}{2(1+\sigma)}}}\right) \Delta \stackrel{(iv)}{\leq} \left(1 + (1 + \sigma)\left(1 + \frac{1}{1 + \sigma}\right)\right) \Delta = (3 + \sigma)\Delta, \end{aligned} \quad (234)$$

1265 where (i) holds by $\underline{q} = 0$ (see (229)), (ii) follows from plugging in $p = q + \Delta = \frac{\sigma}{1+\sigma} + \Delta$,
1266 and (iii) and (iv) arises from $\Delta = \min\left\{\frac{1}{4}(1 - \gamma), \frac{1}{2(1+\sigma)}\right\} \leq 1$ in (198). Combining (233)
1267 and (234) yields

$$\frac{\sigma + 1}{2}\Delta \leq \underline{p} \leq (3 + \sigma)\Delta. \quad (235)$$

1268 **Step 3: combining all the results.** Finally, summing up the results for both \underline{q} (in (227) and (229))
 1269 and \underline{p} (in (232) and (235)), we arrive at the advertised bound.

1270 F.3.2 Proof of Lemma 19

1271 **The robust value function for any policy π .** For any \mathcal{M}_ϕ with $\phi \in \{0, 1\}$, we first characterize
 1272 the robust value function of any policy π over different states.

1273 Towards this, it is easily observed that for any policy π , the robust value functions at state $s = 1$ or
 1274 any $s \in \{2, 3, \dots, S-1\}$ obey

$$V_\phi^{\pi, \sigma}(1) \stackrel{(i)}{=} 1 + \gamma V_\phi^{\pi, \sigma}(1) = \frac{1}{1 - \gamma} \quad (236a)$$

1275 and

$$\forall s \in \{2, 3, \dots, S\} : \quad V_\phi^{\pi, \sigma}(s) \stackrel{(ii)}{=} 0 + \gamma V_\phi^{\pi, \sigma}(1) = \frac{\gamma}{1 - \gamma}, \quad (236b)$$

1276 where (i) and (ii) is according to the facts that the transitions defined over states $s \geq 1$ in (123) give
 1277 only one possible next state 1, leading to a non-random transition in the uncertainty set associated with
 1278 χ^2 divergence, and $r(1, a) = 1$ for all $a \in \mathcal{A}$ and $r(s, a) = 0$ holds all $(s, a) \in \{2, 3, \dots, S-1\} \times \mathcal{A}$.
 1279 To continue, the robust value function at state 0 with policy π satisfies

$$\begin{aligned} V_\phi^{\pi, \sigma}(0) &= \mathbb{E}_{a \sim \pi(\cdot | 0)} \left[r(0, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,a}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} \right] \\ &= 0 + \gamma \pi(\phi | 0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,\phi}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} + \gamma \pi(1 - \phi | 0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,1-\phi}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} \end{aligned} \quad (237)$$

$$\stackrel{(i)}{\leq} \frac{\gamma}{1 - \gamma}, \quad (238)$$

1280 where (i) holds by that $\|V_\phi^{\pi, \sigma}\|_\infty \leq \frac{1}{1 - \gamma}$. Summing up the results in (236b) and (238) leads to

$$\forall s \in \{2, 3, \dots, S\}, \quad V_\phi^{\pi, \sigma}(1) > V_\phi^{\pi, \sigma}(s) \geq V_\phi^{\pi, \sigma}(0). \quad (239)$$

1281 With the transition kernel in (123) over state 0 and the fact in (239), (237) can be rewritten as

$$\begin{aligned} V_\phi^{\pi, \sigma}(0) &= \gamma \pi(\phi | 0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,\phi}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} + \gamma \pi(1 - \phi | 0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0,1-\phi}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} \\ &\stackrel{(i)}{=} \gamma \pi(\phi | 0) \left[\underline{p} V_\phi^{\pi, \sigma}(1) + (1 - \underline{p}) V_\phi^{\pi, \sigma}(0) \right] + \gamma \pi(1 - \phi | 0) \left[\underline{q} V_\phi^{\pi, \sigma}(1) + (1 - \underline{q}) V_\phi^{\pi, \sigma}(0) \right] \\ &\stackrel{(ii)}{=} \gamma z_\phi^\pi V_\phi^{\pi, \sigma}(1) + \gamma (1 - z_\phi^\pi) V_\phi^{\pi, \sigma}(0) \\ &= \frac{\gamma z_\phi^\pi}{(1 - \gamma) \left(1 - \gamma(1 - z_\phi^\pi) \right)}, \end{aligned} \quad (240)$$

1282 where (i) holds by the definition of \underline{p} and \underline{q} in (200), (ii) follows from the definition of z_ϕ^π in (203),
 1283 and the last line holds by applying (236a) and solving the resulting linear equation for $V_\phi^{\pi, \sigma}(0)$.

1284 **Optimal policy and its optimal value function.** To continue, observing that $V_\phi^{\pi, \sigma}(0) =: f(z_\phi^\pi)$ is
 1285 increasing in z_ϕ^π since the derivative of $f(z_\phi^\pi)$ w.r.t. z_ϕ^π obeys

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

1286 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
 1287 $\underline{p} \geq \underline{q}$ in (200)), the optimal robust policy in state 0 thus obeys

$$\pi_\phi^*(\phi | 0) = 1. \quad (241)$$

1288 Considering that the action does not influence the state transition for all states $s > 0$, without loss of
 1289 generality, we choose the optimal robust policy to obey

$$\forall s > 0 : \quad \pi_\phi^*(\phi | s) = 1. \quad (242)$$

1290 Taking $\pi = \pi_\phi^*$ and $z_\phi^{\pi_\phi^*} = \underline{p}$ in (240), we complete the proof by showing the corresponding optimal
 1291 robust value function at state 0 as follows:

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

1292 F.3.3 Proof of the claim (208)

1293 Plugging in the definition of φ , we arrive at that for any policy π ,

$$\begin{aligned} \langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\pi,\sigma} \rangle &= V_\phi^{*,\sigma}(0) - V_\phi^{\pi,\sigma}(0) \\ &\stackrel{(i)}{=} \frac{\gamma \underline{p}}{(1-\gamma) \left(1 - \gamma (1 - \underline{p})\right)} - \frac{\gamma z_\phi^\pi}{(1-\gamma) \left(1 - \gamma (1 - z_\phi^\pi)\right)} \\ &= \frac{\gamma (\underline{p} - z_\phi^\pi)}{(1-\gamma(1-\underline{p})) \left(1 - \gamma(1 - z_\phi^\pi)\right)} \stackrel{(ii)}{\geq} \frac{\gamma (\underline{p} - z_\phi^\pi)}{(1-\gamma(1-\underline{p}))^2} \\ &\stackrel{(iii)}{=} \frac{\gamma(\underline{p} - \underline{q})(1 - \pi(\phi | 0))}{(1-\gamma(1-\underline{p}))^2}, \end{aligned} \quad (243)$$

1294 where (i) holds by applying Lemma 19, (ii) arises from $z_\phi^\pi \leq \underline{p}$ (see the definition of z_ϕ^π in (203) and
 1295 the fact $\underline{p} \geq \underline{q} + \frac{3\Delta}{4}$ in (200)), and (iii) follows from the definition of z_ϕ^π in (203).

1296 To further control (243), we consider it in two cases separately:

1297 • Case 1: $\sigma \in (0, \frac{1-\gamma}{4})$. In this case, applying Lemma 18 to (243) yields

$$\begin{aligned} \langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\pi,\sigma} \rangle &\geq \frac{\gamma(\underline{p} - \underline{q})(1 - \pi(\phi | 0))}{(1-\gamma(1-\underline{p}))^2} \geq \frac{\gamma \frac{3\Delta}{4} (1 - \pi(\phi | 0))}{\left(1 - \gamma \left(1 - \frac{5(1-\gamma)}{4}\right)\right)^2} \\ &\geq \frac{\Delta(1 - \pi(\phi | 0))}{9(1-\gamma)^2} = 2\varepsilon(1 - \pi(\phi | 0)), \end{aligned} \quad (244)$$

1298 where the penultimate inequality follows from $\gamma \geq 3/4$, and the last inequality holds by
 1299 taking the specification of Δ in (207) as follows:

$$\Delta = 18(1-\gamma)^2\varepsilon. \quad (245)$$

1300 It is easily verified that taking $\varepsilon \leq \frac{1}{72(1-\gamma)}$ as in (206) directly leads to meeting the
 1301 requirement in (198), i.e., $\Delta \leq \frac{1}{4}(1-\gamma)$.

1302 • Case 2: $\sigma \in [\frac{1-\gamma}{4}, \infty)$. Similarly, applying Lemma 18 to (243) gives

$$\langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\pi,\sigma} \rangle \geq \frac{\gamma(\underline{p} - \underline{q})(1 - \pi(\phi | 0))}{(1-\gamma(1-\underline{p}))^2} \geq \frac{\gamma \frac{\sigma+1}{2} \Delta (1 - \pi(\phi | 0))}{\min \left\{ 1, (1-\gamma(1-(3+\sigma)\Delta))^2 \right\}} \quad (246)$$

1303 Before continuing, it can be verified that

$$\begin{aligned} 1 - \gamma(1 - (3 + \sigma)\Delta) &= 1 - \gamma + \gamma(3 + \sigma)\Delta \\ &\stackrel{(i)}{\leq} 1 - \gamma + (3 + \sigma) \min \left\{ \frac{1}{4}(1 - \gamma), \frac{1}{2(\sigma + 1)} \right\} \end{aligned}$$

$$\leq \min \left\{ 2(1 + \sigma)(1 - \gamma), \frac{3}{2} \right\}, \quad (247)$$

1304 where (i) is obtained by $\Delta \leq \min \left\{ \frac{1}{4}(1 - \gamma), \frac{1}{2(1 + \sigma)} \right\}$ (see (197)). Applying the above fact
1305 to (246) gives

$$\begin{aligned} \langle \varphi, V_{\phi}^{*,\sigma} - V_{\phi}^{\pi,\sigma} \rangle &\geq \frac{\gamma^{\frac{\sigma+1}{2}} \Delta (1 - \pi(\phi|0))}{\min \left\{ 1, (1 - \gamma(1 - (3 + \sigma)\Delta))^2 \right\}} \stackrel{(i)}{\geq} \frac{3(\sigma + 1)\Delta(1 - \pi(\phi|0))}{8 \min \{4(1 + \sigma)^2(1 - \gamma)^2, 1\}} \\ &\geq \frac{\Delta(1 - \pi(\phi|0))}{\min \left\{ 32(1 + \sigma)(1 - \gamma)^2, \frac{8}{3(1 + \sigma)} \right\}} = 2\varepsilon(1 - \pi(\phi|0)), \end{aligned} \quad (248)$$

1306 where (i) holds by $\gamma \geq \frac{3}{4}$ and (246), and the last equality holds by the specification in (207):

$$\Delta = \begin{cases} 64(1 + \sigma)(1 - \gamma)^2 \varepsilon & \text{if } \sigma \in \left[\frac{1 - \gamma}{4}, \frac{1}{3(1 - \gamma)} \right), \\ \frac{16}{3(1 + \sigma)} \varepsilon & \text{if } \sigma > \frac{1}{3(1 - \gamma)}. \end{cases} \quad (249)$$

1307 As a result, it is easily verified that the requirement in (198)

$$\Delta \leq \min \left\{ \frac{1}{4}(1 - \gamma), \frac{1}{2(1 + \sigma)} \right\} \quad (250)$$

1308 is met if we let

$$\varepsilon \leq \begin{cases} \frac{1}{256(1 + \sigma)(1 - \gamma)} & \text{if } \sigma \in \left[\frac{1 - \gamma}{4}, \frac{1}{3(1 - \gamma)} \right), \\ \frac{3}{32} & \text{if } \sigma > \frac{1}{3(1 - \gamma)}, \end{cases} \quad (251)$$

1309 as in (206).

1310 The proof is then completed by summing up the results in the above two cases.

1311