

Supplementary Material for “Breaking the Sample Size Barrier in Model-Based Reinforcement Learning”

This supplementary material includes details that support the main text “Breaking the Sample Size Barrier in Model-Based Reinforcement Learning” accepted to Operations Research. In particular, the proofs of all theorems, corollaries and technical lemmas are included.

EC.1. Preliminary facts

We begin by recording a few elementary facts about \mathbf{P}^π and \mathbf{P}_π (see definitions in (24)). These are standard results and we omit the proofs for brevity.

LEMMA EC.1. *For any policy π , any probability transition matrix $\mathbf{P} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}| \times |\mathcal{S}|}$ and any $0 < \gamma < 1$, one has*

- (a) $(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} = \sum_{i=0}^{\infty} (\gamma \mathbf{P}_\pi)^i$;
- (b) All entries of the matrix $(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1}$ are non-negative;
- (c) $\|(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1}\|_1 \leq 1/(1 - \gamma)$;
- (d) $(1 - \gamma)(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{1} = \mathbf{1}$;
- (e) For any non-negative vectors $\mathbf{0} \leq \mathbf{r}_1 \leq \mathbf{r}_2$ of compatible dimension, one has $\mathbf{0} \leq (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{r}_1 \leq (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{r}_2$.

The above results continue to hold if \mathbf{P}_π is replaced by \mathbf{P}^π .

EC.2. Proofs of auxiliary lemmas: infinite-horizon MDPs

EC.2.1. Proofs of Lemma 1 and Lemma 2

Auxiliary notation and preliminaries. Before proceeding, we define several $|\mathcal{S}|$ -dimensional auxiliary vectors $\mathbf{r}^{(i)}, \mathbf{V}^{(i)}, \widehat{\mathbf{V}}^{(i)}$ ($1 \leq i \leq m$) recursively as follows

$$\begin{aligned}
 \mathbf{r}^{(0)} &:= \mathbf{r}_\pi, & \mathbf{V}^{(0)} &:= (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{r}^{(0)}, & \widehat{\mathbf{V}}^{(0)} &:= (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} \mathbf{r}^{(0)}, \\
 \mathbf{r}^{(1)} &:= \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(0)}]}, & \mathbf{V}^{(1)} &:= (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{r}^{(1)}, & \widehat{\mathbf{V}}^{(2)} &:= (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} \mathbf{r}^{(1)}, \\
 &\vdots & &\vdots & &\vdots \\
 \mathbf{r}^{(m)} &:= \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(m-1)}]}, & \mathbf{V}^{(m)} &:= (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{r}^{(m)}, & \widehat{\mathbf{V}}^{(m)} &:= (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} \mathbf{r}^{(m)},
 \end{aligned} \tag{EC.1}$$

where m will be specified momentarily.

A crucial quantity that appears repeatedly in analyzing the above terms is $\|(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})}\|_\infty$, whose importance was already made apparent in the work Azar et al. (2013). A widely used upper bound on this quantity, originally due to (Azar et al. 2013, Lemma 8), is given by

$$\left\| (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})} \right\|_\infty \leq \frac{2 \log 2}{\gamma(1 - \gamma)^{1.5}} \|\mathbf{r}\|_\infty. \tag{EC.2}$$

This bound turns out to be loose for our purpose, and we develop an improved bound as follows, whose proof is deferred to Appendix [EC.2.1.1](#).

LEMMA EC.2. Consider any policy π and any probability transition matrix $\mathbf{P} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$. Let \mathbf{V} be a vector obeying $\mathbf{V} = (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{r}_\pi$ for some $|\mathcal{S}|$ -dimensional vector $\mathbf{r}_\pi \geq \mathbf{0}$. For any $0 < \gamma < 1$, one has

$$\left\| (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})} \right\|_\infty \leq \frac{4}{\gamma \sqrt{1 - \gamma}} \|\mathbf{V}\|_\infty. \quad (\text{EC.3})$$

REMARK EC.1. In comparison to the bound (EC.2) derived in (Azar et al. 2013, Lemma 8), Lemma EC.2 offers an improved upper bound stated directly in terms of the properties of \mathbf{V} rather than those of \mathbf{r} .

As it turns out, Lemma EC.2 allows us to obtain an entrywise bound for $\mathbf{V}^{(l)}$ ($1 \leq l \leq m$). To begin with, the first term $\mathbf{V}^{(1)}$ satisfies

$$\|\mathbf{V}^{(1)}\|_\infty = \left\| (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(0)}]} \right\|_\infty \quad (\text{EC.4})$$

since $\mathbf{r}^{(1)} = \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(0)}]}$. Next, for any $l > 1$ one has

$$\begin{aligned} \|\mathbf{V}^{(l)}\|_\infty &= \|(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{r}^{(l)}\|_\infty = \left\| (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(l-1)}]} \right\|_\infty \\ &\leq \frac{4}{\gamma \sqrt{1 - \gamma}} \|\mathbf{V}^{(l-1)}\|_\infty, \end{aligned}$$

where the second identity results from the definition of $\mathbf{r}^{(l)}$, and the last inequality comes from Lemma EC.2. As a consequence, applying this inequality recursively gives

$$\|\mathbf{V}^{(l)}\|_\infty \leq \left(\frac{4}{\gamma \sqrt{1 - \gamma}} \right)^{l-1} \|\mathbf{V}^{(1)}\|_\infty. \quad (\text{EC.5})$$

Main proof. Equipped with the above facts, we are now in a position to prove the lemmas, for which we start with the more general one — Lemma 2. Consider any $0 \leq l \leq m$. We first observe that

$$\begin{aligned} \widehat{\mathbf{V}}^{(l)} - \mathbf{V}^{(l)} &= (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} \mathbf{r}^{(l)} - \mathbf{V}^{(l)} \\ &= (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} (\mathbf{I} - \gamma \mathbf{P}_\pi) \mathbf{V}^{(l)} - (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi) \mathbf{V}^{(l)} \\ &= \gamma (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^{(l)}, \end{aligned} \quad (\text{EC.6})$$

where the second line follows since, by definition, $(\mathbf{I} - \gamma \mathbf{P}_\pi) \mathbf{V}^{(l)} = \mathbf{r}^{(l)}$. Suppose that there exists some quantity $\beta_1 > 0$ such that the following condition

$$\left| (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^{(l)} \right| \leq \sqrt{\frac{\beta_1}{N}} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(l)}]} + \frac{\|\mathbf{V}^{(l)}\|_\infty \beta_1}{N} \mathbf{1} \quad (\text{EC.7})$$

holds uniformly for all $0 \leq l \leq m$. Then this combined with (EC.6) gives

$$\begin{aligned} \|\widehat{\mathbf{V}}^{(l)} - \mathbf{V}^{(l)}\|_\infty &= \gamma \left\| (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^{(l)} \right\|_\infty \\ &\stackrel{(i)}{\leq} \gamma \left\| (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^{(l)} \right\|_\infty \\ &\stackrel{(ii)}{\leq} \gamma \sqrt{\frac{\beta_1}{N}} \left\| (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(l)}]} \right\|_\infty + \frac{\gamma \|\mathbf{V}^{(l)}\|_\infty \beta_1}{N} \left\| (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} \mathbf{1} \right\|_\infty. \end{aligned}$$

Here, (i) follows since $(\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1}$ is a non-negative matrix, (ii) comes from (EC.7) and the triangle inequality. Recalling the definition of $\mathbf{r}^{(l)}$ and $\widehat{\mathbf{V}}^{(l)}$ and invoking Lemma EC.1(d), we can further bound the above as

$$\begin{aligned} \|\widehat{\mathbf{V}}^{(l)} - \mathbf{V}^{(l)}\|_\infty &\leq \gamma \sqrt{\frac{\beta_1}{N}} \|\widehat{\mathbf{V}}^{(l+1)}\|_\infty + \frac{\gamma \|\mathbf{V}^{(l)}\|_\infty \beta_1}{(1-\gamma)N} \\ &\leq \gamma \sqrt{\frac{\beta_1}{N}} \|\widehat{\mathbf{V}}^{(l+1)} - \mathbf{V}^{(l+1)}\|_\infty + \gamma \sqrt{\frac{\beta_1}{N}} \|\mathbf{V}^{(l+1)}\|_\infty + \frac{\gamma \beta_1}{(1-\gamma)N} \|\mathbf{V}^{(l)}\|_\infty, \end{aligned} \quad (\text{EC.8})$$

where the last inequality is a consequence of the triangle inequality.

The above inequality (EC.8) provides a recursive relation that in turn allows for an effective upper bound. Specifically, combining the inequalities (EC.5) and (EC.8) leads to

$$\begin{aligned} \|\widehat{\mathbf{V}}^{(0)} - \mathbf{V}^{(0)}\|_\infty &\leq \gamma \sqrt{\frac{\beta_1}{N}} \|\widehat{\mathbf{V}}^{(1)} - \mathbf{V}^{(1)}\|_\infty + \gamma \sqrt{\frac{\beta_1}{N}} \|\mathbf{V}^{(1)}\|_\infty + \frac{\gamma \beta_1}{(1-\gamma)N} \|\mathbf{V}^{(0)}\|_\infty \\ &=: b_1 \|\widehat{\mathbf{V}}^{(1)} - \mathbf{V}^{(1)}\|_\infty + b_1 \|\mathbf{V}^{(1)}\|_\infty + \frac{\gamma \beta_1}{(1-\gamma)N} \|\mathbf{V}^{(0)}\|_\infty, \end{aligned}$$

and for $l \geq 1$,

$$\begin{aligned} \|\widehat{\mathbf{V}}^{(l)} - \mathbf{V}^{(l)}\|_\infty &\leq \gamma \sqrt{\frac{\beta_1}{N}} \|\widehat{\mathbf{V}}^{(l+1)} - \mathbf{V}^{(l+1)}\|_\infty + \left(4 \sqrt{\frac{\beta_1}{(1-\gamma)N}} + \frac{\gamma \beta_1}{(1-\gamma)N} \right) \left(\frac{4}{\gamma \sqrt{1-\gamma}} \right)^{l-1} \|\mathbf{V}^{(1)}\|_\infty \\ &=: b_1 \|\widehat{\mathbf{V}}^{(l+1)} - \mathbf{V}^{(l+1)}\|_\infty + b_2 b_3^{l-1} \|\mathbf{V}^{(1)}\|_\infty. \end{aligned}$$

Here for notational simplicity, we introduce

$$b_1 := \gamma \sqrt{\frac{\beta_1}{N}}, \quad b_2 := 4 \sqrt{\frac{\beta_1}{(1-\gamma)N}} + \frac{\gamma \beta_1}{(1-\gamma)N}, \quad b_3 := \frac{4}{\gamma \sqrt{1-\gamma}}.$$

Invoking the above recursive relation, we can arrange terms to reach

$$\|\widehat{\mathbf{V}}^{(0)} - \mathbf{V}^{(0)}\|_\infty \leq \underbrace{b_1^m \|\widehat{\mathbf{V}}^{(m)} - \mathbf{V}^{(m)}\|_\infty}_{=: \alpha_1} + \underbrace{\left(b_1 b_2 \sum_{l=0}^{m-2} (b_1 b_3)^l + b_1 \right)}_{=: \alpha_2} \|\mathbf{V}^{(1)}\|_\infty + \frac{\gamma \beta_1}{(1-\gamma)N} \|\mathbf{V}^{(0)}\|_\infty. \quad (\text{EC.9})$$

Controlling the quantity α_2 . Now it suffices to control the two terms on the right-hand side of the inequality (EC.9) separately, towards which we shall start with the quantity α_2 . Assuming that

$N \geq 64\beta_1/(1-\gamma)$, one can easily verify $b_1b_3 \leq 1/2$. The summation of the geometric sequence thus gives

$$\begin{aligned} \alpha_2 &:= b_1b_2 \sum_{l=0}^{m-2} (b_1b_3)^l + b_1 \leq \frac{b_1b_2}{1-b_1b_3} + b_1 \leq 2b_1b_2 + b_1 \\ &= \gamma\sqrt{\frac{\beta_1}{N}} \left\{ 1 + 8\sqrt{\frac{\beta_1}{(1-\gamma)N}} + \frac{2\gamma\beta_1}{(1-\gamma)N} \right\} \leq 3\gamma\sqrt{\frac{\beta_1}{N}}, \end{aligned} \quad (\text{EC.10})$$

where the last step holds with the assumption $N \geq \frac{64\beta_1}{1-\gamma}$.

Controlling the quantity α_1 . Next, we proceed to the quantity α_1 , which requires the control of $\|\widehat{\mathbf{V}}^{(m)} - \mathbf{V}^{(m)}\|_\infty$. In view of the identity (EC.6), we obtain

$$\begin{aligned} \|\widehat{\mathbf{V}}^{(m)} - \mathbf{V}^{(m)}\|_\infty &= \gamma \left\| (\mathbf{I} - \gamma\widehat{\mathbf{P}}_\pi)^{-1} (\mathbf{P}_\pi - \widehat{\mathbf{P}}_\pi) \mathbf{V}^{(m)} \right\|_\infty \\ &\leq \gamma \left\| (\mathbf{I} - \gamma\widehat{\mathbf{P}}_\pi)^{-1} \left(\sqrt{\frac{\beta_1}{N}} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(m)}]} + \frac{\|\mathbf{V}^{(m)}\|_\infty \beta_1}{N} \mathbf{1} \right) \right\|_\infty, \end{aligned}$$

where the last inequality follows from the Bernstein-type condition (EC.7) and the fact that $(\mathbf{I} - \gamma\widehat{\mathbf{P}}_\pi)^{-1}$ has non-negative entries. By virtue of the simple relation $\sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(m)}]} \leq \|\mathbf{V}^{(m)}\|_\infty$ and the fact that $\|(\mathbf{I} - \gamma\widehat{\mathbf{P}}_\pi)^{-1}\|_1 \leq \frac{1}{1-\gamma}$ (cf. Lemma EC.1(c)), it is further guaranteed that

$$\begin{aligned} \|\widehat{\mathbf{V}}^{(m)} - \mathbf{V}^{(m)}\|_\infty &\leq \gamma \|(\mathbf{I} - \gamma\widehat{\mathbf{P}}_\pi)^{-1}\|_1 \cdot \left\| \sqrt{\frac{\beta_1}{N}} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(m)}]} + \frac{\|\mathbf{V}^{(m)}\|_\infty \beta_1}{N} \mathbf{1} \right\|_\infty \\ &\leq \frac{\gamma}{1-\gamma} \left(\sqrt{\frac{\beta_1}{N}} + \frac{\beta_1}{N} \right) \|\mathbf{V}^{(m)}\|_\infty, \end{aligned}$$

which combined with the bound (EC.5) yields

$$\|\widehat{\mathbf{V}}^{(m)} - \mathbf{V}^{(m)}\|_\infty \leq \frac{\gamma}{1-\gamma} \left(\sqrt{\frac{\beta_1}{N}} + \frac{\beta_1}{N} \right) \left(\frac{4}{\gamma\sqrt{1-\gamma}} \right)^{m-1} \|\mathbf{V}^{(1)}\|_\infty.$$

Putting the above bounds together yields

$$\begin{aligned} \alpha_1 &:= b_1^m \|\widehat{\mathbf{V}}^{(m)} - \mathbf{V}^{(m)}\|_\infty \leq \left(\gamma\sqrt{\frac{\beta_1}{N}} \right)^m \frac{\gamma}{1-\gamma} \left(\sqrt{\frac{\beta_1}{N}} + \frac{\beta_1}{N} \right) \left(\frac{4}{\gamma\sqrt{1-\gamma}} \right)^{m-1} \|\mathbf{V}^{(1)}\|_\infty \\ &= \left(\sqrt{\frac{16\beta_1}{N(1-\gamma)}} \right)^{m-1} \left(\sqrt{\frac{\beta_1}{N}} + 1 \right) \frac{\gamma^2\beta_1}{(1-\gamma)N} \|\mathbf{V}^{(1)}\|_\infty \\ &\leq \left(\frac{1}{e} \right)^{m-1} \frac{1.1\gamma^2\beta_1}{(1-\gamma)N} \|\mathbf{V}^{(1)}\|_\infty, \end{aligned} \quad (\text{EC.11})$$

where the last inequality holds provided that $N > \frac{16e^2}{1-\gamma}\beta_1$.

Putting all this together. Combining the inequalities (EC.9), (EC.10) and (EC.11) gives

$$\|\widehat{\mathbf{V}}^{(0)} - \mathbf{V}^{(0)}\|_\infty \leq \frac{\gamma\beta_1}{(1-\gamma)N} \|\mathbf{V}^{(0)}\|_\infty + \left\{ 3\gamma\sqrt{\frac{\beta_1}{N}} + \left(\frac{1}{e} \right)^{m-1} \frac{1.1\gamma^2\beta_1}{(1-\gamma)N} \right\} \|\mathbf{V}^{(1)}\|_\infty.$$

To finish up, set $m = \log(\frac{e}{1-\gamma})$ and assume that $N > \frac{16e^2}{1-\gamma}\beta_1$. Recognizing that $\mathbf{V}^{(0)} = \mathbf{V}^\pi$ and $\widehat{\mathbf{V}}^{(0)} = \widehat{\mathbf{V}}^\pi$, we arrive at

$$\begin{aligned} \|\widehat{\mathbf{V}}^\pi - \mathbf{V}^\pi\|_\infty &= \|\widehat{\mathbf{V}}^{(0)} - \mathbf{V}^{(0)}\|_\infty \leq \frac{\gamma\beta_1}{(1-\gamma)N} \|\mathbf{V}^{(0)}\|_\infty + \left\{ 3\gamma\sqrt{\frac{\beta_1}{N}} + \left(\frac{1}{e}\right)^{m-1} \frac{1.1\gamma^2\beta_1}{(1-\gamma)N} \right\} \|\mathbf{V}^{(1)}\|_\infty \\ &\leq \frac{\gamma\beta_1}{(1-\gamma)N} \|\mathbf{V}^\pi\|_\infty + 4\gamma\sqrt{\frac{\beta_1}{N}} \left\| (\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^\pi]} \right\|_\infty, \end{aligned} \quad (\text{EC.12})$$

provided that $N \geq \frac{16e^2\beta_1}{1-\gamma}$.

Proof of Lemma 2. Invoking the inequality (EC.2) to bound the second term of (EC.12), we reach

$$\|\widehat{\mathbf{V}}^\pi - \mathbf{V}^\pi\|_\infty \leq \frac{\gamma\beta_1}{(1-\gamma)N} \|\mathbf{V}^\pi\|_\infty + 8\log 2 \sqrt{\frac{\beta_1}{N(1-\gamma)^3}} \|\mathbf{r}\|_\infty \leq 6\sqrt{\frac{\beta_1}{N(1-\gamma)^3}} \|\mathbf{r}\|_\infty, \quad (\text{EC.13})$$

where the last inequality uses the elementary fact that $\|\mathbf{V}^\pi\|_\infty \leq \frac{1}{1-\gamma} \|\mathbf{r}\|_\infty$ and the assumption that $N > \frac{16e^2}{1-\gamma}\beta_1$. We complete the proof of Lemma 2.

Proof of Lemma 1. Finally, to establish Lemma 1, we observe that: for any fixed policy π , the vector $\mathbf{V}^{(l)}$ is independent of $\widehat{\mathbf{P}}_\pi$. The Bernstein inequality (e.g. (Agarwal et al. 2020, Lemma 6)) then reveals that with probability at least $1 - \delta$,

$$\left| (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^{(l)} \right| \leq \sqrt{\frac{2\log\left(\frac{4m|\mathcal{S}|}{\delta}\right)}{N}} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(l)}]} + \frac{\|\mathbf{V}^{(l)}\|_\infty \log\left(\frac{4m|\mathcal{S}|}{\delta}\right)}{N} \mathbf{1} \quad (\text{EC.14})$$

holds uniformly for all $0 \leq l \leq m$. This means that we can take $\beta_1 := 2\log\left(\frac{4m|\mathcal{S}|}{\delta}\right)$ with $m = \log(\frac{e}{1-\gamma})$ for this case. Combining this with the inequality (EC.12), we derive the advertised instance-dependent bound

$$\|\widehat{\mathbf{V}}^\pi - \mathbf{V}^\pi\|_\infty \leq \frac{2\gamma \log\left(\frac{4|\mathcal{S}| \log \frac{e}{1-\gamma}}{\delta}\right)}{(1-\gamma)N} \|\mathbf{V}^\pi\|_\infty + 4\gamma\sqrt{\frac{2\log\left(\frac{4|\mathcal{S}| \log \frac{e}{1-\gamma}}{\delta}\right)}{(1-\gamma)N}} \left\| (\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^\pi]} \right\|_\infty.$$

Further, this taken collectively with (EC.2) and the crude bound $\|\mathbf{V}^\pi\|_\infty \leq \frac{1}{1-\gamma} \|\mathbf{r}\|_\infty$ gives

$$\begin{aligned} \|\widehat{\mathbf{V}}^\pi - \mathbf{V}^\pi\|_\infty &\leq \frac{2\gamma \log\left(\frac{4|\mathcal{S}| \log \frac{e}{1-\gamma}}{\delta}\right)}{(1-\gamma)^2 N} \|\mathbf{r}\|_\infty + 8\log 2 \sqrt{\frac{2\log\left(\frac{4|\mathcal{S}| \log \frac{e}{1-\gamma}}{\delta}\right)}{(1-\gamma)^3 N}} \|\mathbf{r}\|_\infty \\ &\leq 6\sqrt{\frac{2\log\left(\frac{4|\mathcal{S}| \log \frac{e}{1-\gamma}}{\delta}\right)}{(1-\gamma)^3 N}} \|\mathbf{r}\|_\infty, \end{aligned}$$

with the proviso that $N \geq \frac{32e^2}{1-\gamma} \log\left(\frac{4|\mathcal{S}| \log \frac{e}{1-\gamma}}{\delta}\right)$.

EC.2.1.1. Proof of Lemma EC.2 To begin with, we make the observation that

$$\begin{aligned}\text{Var}_{\mathbf{P}_\pi}(\mathbf{V}) &= \mathbf{P}_\pi(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}_\pi \mathbf{V}) \circ (\mathbf{P}_\pi \mathbf{V}) \\ &= \mathbf{P}_\pi(\mathbf{V} \circ \mathbf{V}) - \frac{1}{\gamma^2}(\mathbf{V} - \mathbf{r}_\pi) \circ (\mathbf{V} - \mathbf{r}_\pi)\end{aligned}\tag{EC.15}$$

$$\begin{aligned}&= \mathbf{P}_\pi(\mathbf{V} \circ \mathbf{V}) - \frac{1}{\gamma^2} \mathbf{V} \circ \mathbf{V} + \frac{2}{\gamma^2} \mathbf{V} \circ \mathbf{r}_\pi - \frac{1}{\gamma^2} \mathbf{r}_\pi \circ \mathbf{r}_\pi \\ &\leq \frac{1}{\gamma^2}(\gamma^2 \mathbf{P}_\pi - \mathbf{I})(\mathbf{V} \circ \mathbf{V}) + \frac{2}{\gamma^2} \mathbf{V} \circ \mathbf{r}_\pi,\end{aligned}\tag{EC.16}$$

where the identity (EC.15) makes use of the relation $\mathbf{V} = \mathbf{r}_\pi + \gamma \mathbf{P}_\pi \mathbf{V}$. In addition, one can deduce that

$$\begin{aligned}\|(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})}\|_\infty &= \frac{1}{1-\gamma} \|(1-\gamma)(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})}\|_\infty \\ &\stackrel{(i)}{\leq} \frac{1}{\sqrt{1-\gamma}} \left\| \sqrt{(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \text{Var}_{\mathbf{P}_\pi}(\mathbf{V})} \right\|_\infty \\ &\stackrel{(ii)}{\leq} \frac{1}{\sqrt{1-\gamma}} \sqrt{\|2(\mathbf{I} - \gamma^2 \mathbf{P}_\pi)^{-1} \text{Var}_{\mathbf{P}_\pi}(\mathbf{V})\|_\infty}.\end{aligned}$$

Here, (i) comes from Jensen's inequality (so that $\mathbb{E}[\sqrt{v}] \leq \sqrt{\mathbb{E}[v]}$) recognizing that each row of $(1-\gamma)(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1}$ is a probability distribution, and Lemma EC.1(d), (ii) is an elementary fact established in (Agarwal et al. 2020, Lemma 4). Combining Lemma EC.1(e) and the inequality (EC.16) further yields

$$\begin{aligned}\|(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})}\|_\infty &\leq \frac{1}{\sqrt{1-\gamma}} \sqrt{\left\| 2(\mathbf{I} - \gamma^2 \mathbf{P}_\pi)^{-1} \left(\frac{1}{\gamma^2} (\gamma^2 \mathbf{P}_\pi - \mathbf{I})(\mathbf{V} \circ \mathbf{V}) + \frac{2}{\gamma^2} \mathbf{V} \circ \mathbf{r}_\pi \right) \right\|_\infty} \\ &\leq \frac{1}{\gamma \sqrt{1-\gamma}} \sqrt{2 \|\mathbf{V} \circ \mathbf{V}\|_\infty} + \frac{2}{\gamma \sqrt{1-\gamma}} \sqrt{\|(\mathbf{I} - \gamma^2 \mathbf{P}_\pi)^{-1} (\mathbf{V} \circ \mathbf{r}_\pi)\|_\infty}.\end{aligned}\tag{EC.17}$$

where the last step arises from the triangle inequality and $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$. This leaves us with two terms to deal with.

Regarding the first term of (EC.17), we observe that $\|(\mathbf{V} \circ \mathbf{V})\|_\infty = \|\mathbf{V}\|_\infty^2$. When it comes to the second term of (EC.17), it is seen that

$$\begin{aligned}\|(\mathbf{I} - \gamma^2 \mathbf{P}_\pi)^{-1} (\mathbf{V} \circ \mathbf{r}_\pi)\|_\infty &\leq \|(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} (\mathbf{V} \circ \mathbf{r}_\pi)\|_\infty \\ &\leq \|\mathbf{V}\|_\infty \|(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{r}_\pi\|_\infty \\ &= \|\mathbf{V}\|_\infty^2.\end{aligned}$$

Here, the first inequality holds true since $(\mathbf{I} - \gamma^2 \mathbf{P}_\pi)^{-1} = \sum_{i=0}^{\infty} \gamma^{2i} \mathbf{P}_\pi^i \leq \sum_{i=0}^{\infty} \gamma^i \mathbf{P}_\pi^i = (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1}$, while the second line holds since \mathbf{V} , \mathbf{r} and $(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1}$ are all non-negative. Substitution into (EC.17) thus yields

$$\|(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})}\|_\infty \leq \frac{\sqrt{2}}{\gamma \sqrt{1-\gamma}} \|\mathbf{V}\|_\infty + \frac{2}{\gamma \sqrt{1-\gamma}} \|\mathbf{V}\|_\infty \leq \frac{4}{\gamma \sqrt{1-\gamma}} \|\mathbf{V}\|_\infty$$

as claimed.

EC.2.2. Proof of Lemma 3

To establish Lemma 3, it suffices to check that \mathbf{V}^* and \mathbf{Q}^* satisfy the Bellman optimality equations underlying \mathcal{M}_{s,a,u^*} . Towards this end, we study the absorbing state-action pair (s, a) and other pairs separately. For notational simplicity, we shall let \mathbf{P}^{abs} and $r^{\text{abs}}(\cdot, \cdot)$ denote respectively the probability transition matrix and the reward function associated with \mathcal{M}_{s,a,u^*} .

First, we observe that, by construction,

$$r^{\text{abs}}(s, a) + \gamma(\mathbf{P}^{\text{abs}}\mathbf{V}^*)_{s,a} = u^* + \gamma V^*(s).$$

Recall that \mathbf{V}^* satisfies the Bellman optimality equation w.r.t. the original MDP, namely, $Q^*(s, a) = r(s, a) + \gamma(\mathbf{P}\mathbf{V}^*)_{s,a}$. This together with our choice of u^* gives

$$u^* + \gamma V^*(s) = Q^*(s, a) - \gamma V^*(s) + \gamma V^*(s) = Q^*(s, a).$$

Putting the above identities together, we arrive at

$$Q^*(s, a) = r^{\text{abs}}(s, a) + \gamma(\mathbf{P}^{\text{abs}}\mathbf{V}^*)_{s,a}. \quad (\text{EC.18})$$

Next, consider any state-action pair $(s', a') \neq (s, a)$. Recalling again the properties of \mathcal{M}_{s,a,u^*} , we reach

$$r^{\text{abs}}(s', a') + \gamma(\mathbf{P}^{\text{abs}}\mathbf{V}^*)_{s',a'} = r(s', a') + \gamma(\mathbf{P}\mathbf{V}^*)_{s',a'} = Q^*(s', a'). \quad (\text{EC.19})$$

Here the last identity is due to the Bellman equation w.r.t. the original MDP. Combining (EC.18) and (EC.19) implies that \mathbf{V}^* and \mathbf{Q}^* satisfy Bellman's optimality equations in \mathcal{M}_{s,a,u^*} , thus concluding the proof.

EC.2.3. Proof of Lemma 4

Our first observation is that $\widehat{\mathbf{Q}}_{s,a,u}^*$ satisfies Lipschitz continuity w.r.t. u in the sense that

$$\|\widehat{\mathbf{Q}}_{s,a,u}^* - \widehat{\mathbf{Q}}_{s,a,u'}^*\|_\infty \leq \frac{1}{1-\gamma}|u - u'|. \quad (\text{EC.20})$$

The proof of this relation is identical to that of (Agarwal et al. 2020, Lemma 8); we omit here for brevity. In view of Lemma 3, if we set $\widehat{u}^* := r(s, a) + \gamma(\widehat{\mathbf{P}}\widehat{\mathbf{V}}^*)_{s,a} - \gamma\widehat{V}^*(s)$, then one has

$$\widehat{\mathbf{Q}}_{s,a,\widehat{u}^*}^* = \widehat{\mathbf{Q}}^* \quad \text{and} \quad \widehat{\mathbf{V}}_{s,a,\widehat{u}^*}^* = \widehat{\mathbf{V}}^*.$$

In addition, there exists a point u_0 in the epsilon-net $\mathcal{N}_{(1-\gamma)\omega/4}$ such that $|\widehat{u}^* - u_0| \leq (1-\gamma)\omega/4$, which combined with the Lipschitz continuity property (EC.20) gives

$$\|\widehat{\mathbf{Q}}^* - \widehat{\mathbf{Q}}_{s,a,u_0}^*\|_\infty = \|\widehat{\mathbf{Q}}_{s,a,\widehat{u}^*}^* - \widehat{\mathbf{Q}}_{s,a,u_0}^*\|_\infty \leq \frac{1}{1-\gamma}|\widehat{u}^* - u_0| \leq \frac{\omega}{4}. \quad (\text{EC.21})$$

Additionally, for any $s' \in \mathcal{S}$ and any $a_1, a_2 \in \mathcal{A}$ with $a_1 \neq a_2$, we have the following decomposition

$$\begin{aligned}
& \widehat{Q}_{s,a,u_0}^*(s', a_1) - \widehat{Q}_{s,a,u_0}^*(s', a_2) \\
&= \widehat{Q}^*(s', a_1) - \widehat{Q}^*(s', a_2) + \widehat{Q}_{s,a,u_0}^*(s', a_1) - \widehat{Q}^*(s', a_1) - \left(\widehat{Q}_{s,a,u_0}^*(s', a_2) - \widehat{Q}^*(s', a_2) \right) \\
&\geq \widehat{Q}^*(s', a_1) - \widehat{Q}^*(s', a_2) - 2 \left\| \widehat{Q}^* - \widehat{Q}_{s,a,u_0}^* \right\|_\infty \\
&\geq \widehat{Q}^*(s', a_1) - \widehat{Q}^*(s', a_2) - \frac{\omega}{2}, \tag{EC.22}
\end{aligned}$$

where the last inequality invokes the inequality (EC.21). Moreover, our separation condition defined in (41) requires that: for any $s' \in \mathcal{S}$ and any $a_2 \neq \widehat{\pi}^*(s')$, one has $\widehat{Q}^*(s', \widehat{\pi}^*(s')) - \widehat{Q}^*(s', a_2) \geq \omega$, which together with (EC.22) reveals that

$$\widehat{Q}_{s,a,u_0}^*(s', \widehat{\pi}^*(s')) - \widehat{Q}_{s,a,u_0}^*(s', a_2) \geq \widehat{Q}^*(s', \widehat{\pi}^*(s')) - \widehat{Q}^*(s', a_2) - \frac{\omega}{2} \geq \frac{\omega}{2}. \tag{EC.23}$$

Given that $\widehat{\pi}_{s,a,u_0}^*(s') := \arg \max_{a'} \widehat{Q}_{s,a,u_0}^*(s', a')$, it is seen from (EC.23) that

$$\widehat{\pi}_{s,a,u_0}^*(s') = \widehat{\pi}^*(s'),$$

which holds true for all $s' \in \mathcal{S}$. This concludes the proof.

EC.2.4. Proof of Lemma 5

We start by bounding $\left\| \widehat{\mathbf{V}}^{\widehat{\pi}^*} - \mathbf{V}^{\widehat{\pi}^*} \right\|_\infty$. Recall the definition of the series $\{\mathbf{V}^{(l)}\}$ in (EC.1). Throughout this proof, we shall write $\mathbf{V}_\pi^{(l)}$ instead in order to make apparent the dependency on the policy π .

For each state-action pair (s, a) , let us construct the epsilon-net $\mathcal{N}_{(1-\gamma)\omega/4}$ as in the expression (42). For every $u \in \mathcal{N}_{(1-\gamma)\omega/4}$, recall that $\widehat{\pi}_{s,a,u}^*$ is defined as the optimal policy with respect to the (s, a) -absorbing MDP $\widehat{\mathcal{M}}_{s,a,u}$. By construction, the set of policies $\widehat{\pi}_{s,a,u}^*$ ($u \in \mathcal{N}_{(1-\gamma)\omega/4}$) is independent of $\widehat{\mathbf{P}}_{s,a}$. The Bernstein inequality (e.g. (Agarwal et al. 2020, Lemma 6)) taken together with the union bound thus guarantees that with probability at least $1 - \delta$,

$$\left| (\widehat{\mathbf{P}} - \mathbf{P})_{s,a} \mathbf{V}_{\widehat{\pi}_{s,a,u}^*}^{(l)} \right| \leq \sqrt{\frac{\beta_1}{N} \sqrt{(\text{Var}_{\mathbf{P}}[\mathbf{V}_{\widehat{\pi}_{s,a,u}^*}^{(l)}])_{s,a}}} + \frac{\left\| \mathbf{V}_{\widehat{\pi}_{s,a,u}^*}^{(l)} \right\|_\infty \beta_1}{N} \tag{EC.24}$$

holds uniformly over all $0 \leq l \leq \log \frac{e}{1-\gamma}$, $u \in \mathcal{N}_{(1-\gamma)\omega/4}$, $(s, a) \in \mathcal{S} \times \mathcal{A}$. Here, β_1 is given by

$$\beta_1 := 2 \log \left(\frac{4 \log \left(\frac{e}{1-\gamma} \right) |\mathcal{N}_{(1-\gamma)\omega/4}| |\mathcal{S}| |\mathcal{A}|}{\delta} \right) \leq 2 \log \left(\frac{32}{(1-\gamma)^2 \omega \delta} |\mathcal{S}| |\mathcal{A}| \log \left(\frac{e}{1-\gamma} \right) \right),$$

where we have used the fact $|\mathcal{N}_{(1-\gamma)\omega/4}| \leq \frac{8}{(1-\gamma)^2 \omega}$. In addition, for any $0 < \omega < 1$, Lemma 4 guarantees that for each state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$, there exists a point $u_0 \in \mathcal{N}_{(1-\gamma)\omega/4}$ such that $\widehat{\pi}^* = \widehat{\pi}_{s,a,u_0}^*$. Invoking this important fact, we obtain

$$\left| (\widehat{\mathbf{P}} - \mathbf{P})_{s,a} \mathbf{V}_{\widehat{\pi}^*}^{(l)} \right| = \left| (\widehat{\mathbf{P}} - \mathbf{P})_{s,a} \mathbf{V}_{\widehat{\pi}_{s,a,u_0}^*}^{(l)} \right|$$

$$\begin{aligned}
&\leq \sqrt{\frac{\beta_1}{N}} \sqrt{(\text{Var}_{\mathcal{P}}[\mathbf{V}_{\hat{\pi}^*, a, u_0}^{(l)}])_{s,a}} + \frac{\|\mathbf{V}_{\hat{\pi}^*, a, u_0}^{(l)}\|_{\infty} \beta_1}{N} \\
&= \sqrt{\frac{\beta_1}{N}} \sqrt{(\text{Var}_{\mathcal{P}}[\mathbf{V}_{\hat{\pi}^*}^{(l)}])_{s,a}} + \frac{\|\mathbf{V}_{\hat{\pi}^*}^{(l)}\|_{\infty} \beta_1}{N}.
\end{aligned}$$

The above inequality further allows us to deduce that, with probability $1 - \delta$,

$$\begin{aligned}
\left| (\hat{\mathbf{P}}_{\hat{\pi}^*} - \mathbf{P}_{\hat{\pi}^*}) \mathbf{V}_{\hat{\pi}^*}^{(l)} \right| &= \left| \mathbf{\Pi}^{\hat{\pi}^*} (\hat{\mathbf{P}} - \mathbf{P}) \mathbf{V}_{\hat{\pi}^*}^{(l)} \right| \\
&\leq \sqrt{\frac{\beta_1}{N}} \sqrt{\mathbf{\Pi}^{\hat{\pi}^*} \text{Var}_{\mathcal{P}}[\mathbf{V}_{\hat{\pi}^*}^{(l)}]} + \frac{\|\mathbf{V}_{\hat{\pi}^*}^{(l)}\|_{\infty} \beta_1}{N} \mathbf{1} \\
&= \sqrt{\frac{\beta_1}{N}} \sqrt{\text{Var}_{\mathcal{P}_{\hat{\pi}^*}}[\mathbf{V}_{\hat{\pi}^*}^{(l)}]} + \frac{\|\mathbf{V}_{\hat{\pi}^*}^{(l)}\|_{\infty} \beta_1}{N} \mathbf{1}.
\end{aligned}$$

The above derivation validates the assumption required for Lemma 2. As a result, if $N > \frac{16e^2}{1-\gamma} \beta_1$ and $0 < \omega < 1$, then Lemma 2 leads to the advertised bound

$$\begin{aligned}
\|\hat{\mathbf{V}}^{\hat{\pi}^*} - \mathbf{V}^{\hat{\pi}^*}\|_{\infty} &\leq \frac{6}{1-\gamma} \sqrt{\frac{\beta_1}{N(1-\gamma)}} \leq \frac{6}{1-\gamma} \sqrt{\frac{2 \log \left(\frac{32}{(1-\gamma)^2 \omega \delta} |\mathcal{S}| |\mathcal{A}| \log \left(\frac{e}{1-\gamma} \right) \right)}{N(1-\gamma)}} \\
&\leq 6 \sqrt{\frac{2 \log \left(\frac{32 |\mathcal{S}| |\mathcal{A}|}{(1-\gamma)^3 \omega \delta} \right)}{N(1-\gamma)^3}}. \tag{EC.25}
\end{aligned}$$

Finally, we move on to the term $\mathbf{V}^* - \hat{\mathbf{V}}^{\pi^*}$. Given that π^* is independent of $\hat{\mathbf{P}}$, invoke Lemma 1 to reach

$$\begin{aligned}
\|\hat{\mathbf{V}}^{\pi^*} - \mathbf{V}^*\|_{\infty} &= \|\hat{\mathbf{V}}^{\pi^*} - \mathbf{V}^{\pi^*}\|_{\infty} \leq 6\sqrt{2} \sqrt{\frac{\log \left(\frac{4|\mathcal{S}|}{\delta} \right) + \log \log \left(\frac{e}{1-\gamma} \right)}{N(1-\gamma)^3}} \\
&\leq 6 \sqrt{\frac{2 \log \left(\frac{32 |\mathcal{S}| |\mathcal{A}|}{(1-\gamma)^2 \omega \delta} \right)}{N(1-\gamma)^3}} \tag{EC.26}
\end{aligned}$$

with probability at least $1 - \delta$. This together with (34) and (EC.25) immediately establishes Lemma 5.

EC.2.5. Proof of Lemma 6

The proofs for Q_p^* and \hat{Q}_p^* are exactly the same; for the sake of conciseness, we shall only provide the proof for Q_p^* . Here we aim to prove a more general result than Lemma 6, namely, with probability at least $1 - \delta$,

$$\forall s \in \mathcal{S} \text{ and } a_1, a_2 \in \mathcal{A} \text{ with } a_1 \neq a_2 : \quad \left| Q_p^*(s, a_1) - Q_p^*(s, a_2) \right| > \frac{\xi \delta (1-\gamma)}{4|\mathcal{S}| |\mathcal{A}|^2}.$$

Consider any state s and any actions $a_1 \neq a_2$. In what follows, we allow $r_p(s, a_1) = \tau$ to vary, while *freezing* the values of all other rewards $\{r_p(\tilde{s}, a) \mid (\tilde{s}, a) \neq (s, a_1)\}$. To streamline notation, we define

- $\mathbf{r}_\tau = [r_\tau(s, a)]_{(s, a) \in \mathcal{S} \times \mathcal{A}}$: the reward vector obeying

$$r_\tau(s, a_1) = \tau \quad \text{and} \quad r_\tau(\tilde{s}, a) = r_p(\tilde{s}, a) \quad \text{for all } (\tilde{s}, a) \neq (s, a_1);$$

- Q_τ^* : the optimal Q-function when the reward vector is \mathbf{r}_τ ;
- V_τ^* : the optimal value function when the reward vector is \mathbf{r}_τ ;
- π_τ^* : the optimal policy when the reward vector is \mathbf{r}_τ .

Additionally, we claim for the moment that there exists a phase transition boundary τ_{th} such that

$$\pi_\tau^*(s) \neq a_1, \quad \text{for all } \tau < \tau_{\text{th}}; \quad (\text{EC.27a})$$

$$\pi_\tau^*(s) = a_1, \quad \text{for all } \tau > \tau_{\text{th}}. \quad (\text{EC.27b})$$

The proof of this claim is deferred to the end of this section. To establish Lemma 6, the idea is to control the size of the set

$$\mathcal{I}_{0, \omega} := \{ \tau \mid |Q_\tau^*(s, a_1) - Q_\tau^*(s, a_2)| < \omega \} \quad (\text{EC.28})$$

for some $\omega > 0$ to be specified shortly. As motivated by (EC.27), we further break down this set into two parts $\mathcal{I}_{0, \omega} = \mathcal{I}_{1, \omega} \cup \mathcal{I}_{2, \omega}$, where

$$\mathcal{I}_{1, \omega} := \{ \tau \mid \tau < \tau_{\text{th}}, |Q_\tau^*(s, a_1) - Q_\tau^*(s, a_2)| < \omega \}, \quad (\text{EC.29a})$$

$$\mathcal{I}_{2, \omega} := \{ \tau \mid \tau \geq \tau_{\text{th}}, |Q_\tau^*(s, a_1) - Q_\tau^*(s, a_2)| < \omega \}. \quad (\text{EC.29b})$$

In what follows, we first control the size of each set separately, and then demonstrate that the probability of these events happening is very small.

Step 1. We begin with $\mathcal{I}_{1, \omega}$ associated with the range $\tau < \tau_{\text{th}}$. In this case, the value function V_τ^* does not vary with τ , since the reward $r_\tau(s, a_1) = \tau$ is never active when calculating V_τ^* (by virtue of (EC.27a)). Thus, the Bellman equation allows us to write

$$Q_\tau^*(s, a_1) = \tau + B_1 \quad \text{and} \quad Q_\tau^*(s, a_2) = B_2$$

for some quantities B_1 and B_2 , where neither B_1 nor B_2 depends on the value of τ . Armed with this observation, we can easily show that: for any $\omega > 0$, the interval $\mathcal{I}_{1, \omega}$ (cf. (EC.29a)) obeys

$$\mathcal{I}_{1, \omega} \subseteq \{ \tau \mid |\tau + B_1 - B_2| < \omega \},$$

and hence has length (or Lebesgue measure) at most 2ω .

Step 2. We then move on to $\mathcal{I}_{2, \omega}$ associated with the range $\tau > \tau_{\text{th}}$ in which case $\pi_\tau^*(s) = a_1$. Towards this, we first make some useful observations.

- To begin with, given the relation $\mathbf{Q}_\tau^* = \mathbf{r}_\tau + \gamma \mathbf{P} \mathbf{V}_\tau^*$, it is easily seen that for any $\tau_2 > \tau_1 > \tau_{\text{th}}$,

$$\mathbf{0} \leq \mathbf{Q}_{\tau_2}^* - \mathbf{Q}_{\tau_1}^* \leq \mathbf{r}_{\tau_2} - \mathbf{r}_{\tau_1} + \gamma \|\mathbf{V}_{\tau_2}^* - \mathbf{V}_{\tau_1}^*\|_\infty. \quad (\text{EC.30})$$

In addition, for any state-action pair $(\tilde{s}, a) \neq (s, a_1)$, by construction we have $r_{\tau_2}(\tilde{s}, a) - r_{\tau_1}(\tilde{s}, a) = 0$, which together with (EC.30) indicates that

$$\forall (\tilde{s}, a) \neq (s, a_1): \quad 0 \leq Q_{\tau_2}^*(\tilde{s}, a) - Q_{\tau_1}^*(\tilde{s}, a) \leq \gamma \|\mathbf{V}_{\tau_2}^* - \mathbf{V}_{\tau_1}^*\|_\infty. \quad (\text{EC.31})$$

- Next, observe that for any $\tau_2 > \tau_1 > \tau_{\text{th}}$,

$$\forall s' \in \mathcal{S}: \quad 0 \leq V_{\tau_2}^*(s') - V_{\tau_1}^*(s') = \max_a Q_{\tau_2}^*(s', a) - \max_a Q_{\tau_1}^*(s', a) \leq \|\mathbf{Q}_{\tau_2}^* - \mathbf{Q}_{\tau_1}^*\|_\infty \quad (\text{EC.32})$$

and hence $\|\mathbf{Q}_{\tau_2}^* - \mathbf{Q}_{\tau_1}^*\|_\infty \geq \|\mathbf{V}_{\tau_2}^* - \mathbf{V}_{\tau_1}^*\|_\infty$. This combined with (EC.31) and the fact $\gamma < 1$ implies that

$$Q_{\tau_2}^*(s, a_1) - Q_{\tau_1}^*(s, a_1) \geq \|\mathbf{V}_{\tau_2}^* - \mathbf{V}_{\tau_1}^*\|_\infty, \quad (\text{EC.33})$$

which together with the facts $V_{\tau_1}^*(s) = Q_{\tau_1}^*(s, a_1)$ and $V_{\tau_2}^*(s) = Q_{\tau_2}^*(s, a_1)$ (by virtue of (EC.27b)) yields

$$\begin{aligned} V_{\tau_2}^*(s) - V_{\tau_1}^*(s) &= Q_{\tau_2}^*(s, a_1) - Q_{\tau_1}^*(s, a_1) \geq \|\mathbf{V}_{\tau_2}^* - \mathbf{V}_{\tau_1}^*\|_\infty \\ \implies Q_{\tau_2}^*(s, a_1) - Q_{\tau_1}^*(s, a_1) &= \|\mathbf{V}_{\tau_2}^* - \mathbf{V}_{\tau_1}^*\|_\infty. \end{aligned}$$

Invoke the Bellman equation to further derive

$$\begin{aligned} Q_{\tau_2}^*(s, a_1) - Q_{\tau_1}^*(s, a_1) &= \|\mathbf{V}_{\tau_2}^* - \mathbf{V}_{\tau_1}^*\|_\infty = \|\mathbf{r}_{\tau_2} + \gamma \mathbf{P} \mathbf{V}_{\tau_2}^* - \mathbf{r}_{\tau_1} - \gamma \mathbf{P} \mathbf{V}_{\tau_1}^*\|_\infty \\ &\geq \|\mathbf{r}_{\tau_2} - \mathbf{r}_{\tau_1}\|_\infty = \tau_2 - \tau_1, \end{aligned} \quad (\text{EC.34})$$

where the last inequality holds since $\mathbf{r}_{\tau_2} + \gamma \mathbf{P} \mathbf{V}_{\tau_2}^* - \mathbf{r}_{\tau_1} - \gamma \mathbf{P} \mathbf{V}_{\tau_1}^* \geq \mathbf{r}_{\tau_2} - \mathbf{r}_{\tau_1} \geq \mathbf{0}$ (due to the monotonicity properties $\mathbf{r}_{\tau_2} \geq \mathbf{r}_{\tau_1}$ and $\mathbf{V}_{\tau_2}^* \geq \mathbf{V}_{\tau_1}^*$), and the last identity follows from the definition of \mathbf{r}_τ .

With the above two properties (EC.31) and (EC.34) in mind, we are ready to locate $\mathcal{I}_{2,\omega}$ by showing that

$$\mathcal{I}_{2,\omega} \subseteq \left[\tau_{\text{th}}, \tau_{\text{th}} + \frac{\omega}{1-\gamma} \right]. \quad (\text{EC.35})$$

Given that $Q_{\tau_{\text{th}}}^*(s, a_1) \geq Q_{\tau_{\text{th}}}^*(s, a_2)$ (in view of (EC.27)), we have for any $\tau \geq \tau_{\text{th}}$ and any $a_2 \neq a_1$ that

$$Q_\tau^*(s, a_1) - Q_\tau^*(s, a_2) \geq (Q_\tau^*(s, a_1) - Q_{\tau_{\text{th}}}^*(s, a_1)) - (Q_\tau^*(s, a_2) - Q_{\tau_{\text{th}}}^*(s, a_2))$$

$$\geq (1 - \gamma)(Q_\tau^*(s, a_1) - Q_{\tau_{\text{th}}}^*(s, a_1)). \quad (\text{EC.36})$$

Here, the last inequality holds since

$$Q_\tau^*(s, a_2) - Q_{\tau_{\text{th}}}^*(s, a_2) \stackrel{\text{(i)}}{\leq} \gamma \|\mathbf{V}_\tau^* - \mathbf{V}_{\tau_{\text{th}}}^*\|_\infty \stackrel{\text{(ii)}}{\leq} \gamma(Q_\tau^*(s, a_1) - Q_{\tau_{\text{th}}}^*(s, a_1)),$$

where (i) follows from (EC.31) and (ii) is due to (EC.33). As a result, for any $\tau > \tau_{\text{th}} + \frac{\omega}{1-\gamma}$, one can invoke (EC.36) and (EC.34) to see that

$$\forall a_2 \neq a_1: \quad Q_\tau^*(s, a_1) - Q_\tau^*(s, a_2) \geq (1 - \gamma)(Q_\tau^*(s, a_1) - Q_{\tau_{\text{th}}}^*(s, a_1)) \geq (1 - \gamma)(\tau - \tau_{\text{th}}) > \omega,$$

which necessarily implies that such a τ does not lie within the interval $\mathcal{I}_{2,\omega}$ as defined in (EC.29b).

This establishes the claimed relation (EC.35).

Step 3. Putting the results in the above two steps together, we see the set $\mathcal{I}_{0,\omega}$ (cf. (EC.28)) has total length (or Lebesgue measure) at most $\frac{3\omega}{1-\gamma}$. Given that $r_p(s, a) = r(s, a) + \zeta(s, a)$ with $\zeta(s, a) \sim \text{Unif}(0, \xi)$, one has

$$\mathbb{P} \left\{ |Q_p^*(s, a_1) - Q_p^*(s, a_2)| < \omega \right\} \leq \mathbb{P} \left\{ r(s, a) + \zeta(s, a) \in \mathcal{I}_{0,\omega} \right\} \leq \frac{3\omega}{\xi(1-\gamma)}.$$

By setting $\omega = \frac{\delta(1-\gamma)\xi}{3|\mathcal{S}||\mathcal{A}|^2}$, we arrive at

$$\mathbb{P} \left\{ |Q_p^*(s, a_1) - Q_p^*(s, a_2)| < \frac{\delta(1-\gamma)\xi}{3|\mathcal{S}||\mathcal{A}|^2} \right\} \leq \frac{\delta}{|\mathcal{S}||\mathcal{A}|^2}.$$

Finally, taking the union bound over all s, a_1, a_2 , we conclude that

$$\mathbb{P} \left\{ \exists s, a_1 \neq a_2: |Q_p^*(s, a_1) - Q_p^*(s, a_2)| < \frac{\delta(1-\gamma)\xi}{3|\mathcal{S}||\mathcal{A}|^2} \right\} \leq \delta,$$

thus establishing Lemma 6 as long as the claim (EC.27) is valid.

Proof of the claim (EC.27). To establish the claim, it suffices to take

$$\tau_{\text{th}} = \sup \{u \mid \pi_\tau^*(s) \neq a_1 \text{ for all } \tau < u\}. \quad (\text{EC.37})$$

It thus suffices to verify (EC.27b) for our choice (EC.37). Towards this, suppose instead that there exist some $\tau_3 < \tau_{\text{th}} \leq \tau_2 < \tau_1$ such that

$$\pi_{\tau_3}^*(s) \neq a_1, \quad \pi_{\tau_2}^*(s) = a_1, \quad \text{and} \quad \pi_{\tau_1}^*(s) \neq a_1.$$

It is straightforward to see that $V_{\tau_1}^* = V_{\tau_3}^*$, since in both cases, the reward $r_\tau(s, a_1)$ does not enter the calculation of the optimal value function (while the rewards in other state-action pairs are identical in both cases). In view of the monotonicity of the value function w.r.t. the reward vector, we have

$$V_{\tau_1}^* = V_{\tau_2}^* = V_{\tau_3}^*.$$

However, this contradicts our assumption that a_1 is the optimal action for state s at τ_2 but not at τ_3 , since enlarging τ_2 to τ_3 otherwise will enlarge the optimal value function $V_{\tau_3}^*$. We have thus established (EC.27). \square

EC.2.6. Proof of Lemma 7

To begin, we find it helpful to introduce a modified reward function $\tilde{\mathbf{r}} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ as follows

$$\tilde{r}(s, a) := \begin{cases} r(s, a) + \hat{V}^*(s) - \hat{Q}^*(s, a), & \text{if } a = \hat{\pi}_c(s), \\ r(s, a), & \text{otherwise.} \end{cases} \quad (\text{EC.38})$$

Armed with this new reward function, we subsequently define a vector $\tilde{\mathbf{Q}} = [\tilde{Q}(s, a)]$ as follows

$$\tilde{\mathbf{Q}} := \tilde{\mathbf{r}} + \gamma \hat{\mathbf{P}} \hat{\mathbf{V}}^*. \quad (\text{EC.39})$$

In view of the Bellman optimality equation $\hat{\mathbf{Q}}^* = \mathbf{r} + \gamma \hat{\mathbf{P}} \hat{\mathbf{V}}^*$, we see that $\tilde{\mathbf{Q}}$ satisfies $\tilde{\mathbf{Q}} = \tilde{\mathbf{r}} + \hat{\mathbf{Q}}^* - \mathbf{r}$, which combined with the construction (EC.38) gives

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \tilde{Q}(s, a) = \begin{cases} \hat{V}^*(s), & \text{if } a = \hat{\pi}_c(s), \\ \hat{Q}^*(s, a), & \text{else.} \end{cases}$$

As a consequence, it is easily seen that

$$\forall s \in \mathcal{S}: \quad \max_a \tilde{Q}(s, a) = \hat{V}^*(s) \quad \text{and} \quad \tilde{Q}(s, \hat{\pi}_c(s)) = \hat{V}^*(s)$$

This taken collectively with (EC.39) demonstrates that $\tilde{\mathbf{Q}}$ and $\hat{\mathbf{V}}^*$ are respectively the optimal Q-function and optimal value function of the MDP $\tilde{\mathcal{M}} = (\mathcal{S}, \mathcal{A}, \tilde{\mathbf{r}}, \hat{\mathbf{P}}, \gamma)$, since they satisfy the Bellman optimality condition w.r.t. $\tilde{\mathcal{M}}$. In addition, if we let $\tilde{\mathbf{V}}^{\hat{\pi}_c}$ represent the value function of the policy $\hat{\pi}_c$ in $\tilde{\mathcal{M}}$, then the preceding relation clearly implies that $\tilde{\mathbf{V}}^{\hat{\pi}_c} = \hat{\mathbf{V}}^*$.

Using the above properties, one can deduce that

$$\begin{aligned} \hat{\mathbf{V}}^* - \hat{\mathbf{V}}^{\hat{\pi}_c} &= \tilde{\mathbf{V}}^{\hat{\pi}_c} - \hat{\mathbf{V}}^{\hat{\pi}_c} \stackrel{(i)}{=} (\mathbf{I} - \gamma \hat{\mathbf{P}}_{\hat{\pi}_c})^{-1} \tilde{\mathbf{r}} - (\mathbf{I} - \gamma \hat{\mathbf{P}}_{\hat{\pi}_c})^{-1} \mathbf{r} = (\mathbf{I} - \gamma \hat{\mathbf{P}}_{\hat{\pi}_c})^{-1} (\tilde{\mathbf{r}} - \mathbf{r}) \\ &\leq \|(\mathbf{I} - \gamma \hat{\mathbf{P}}_{\hat{\pi}_c})^{-1}\|_1 \|\tilde{\mathbf{r}} - \mathbf{r}\|_\infty \mathbf{1} \stackrel{(ii)}{=} \frac{1}{1 - \gamma} \|\tilde{\mathbf{r}} - \mathbf{r}\|_\infty \mathbf{1} \leq \frac{\varsigma}{1 - \gamma} \mathbf{1} \stackrel{(iii)}{\leq} \frac{\xi}{1 - \gamma} \mathbf{1}. \end{aligned}$$

Here, (i) is due to the Bellman equation, (ii) relies on the fact $\|(\mathbf{I} - \gamma \hat{\mathbf{P}}_{\hat{\pi}_c})^{-1}\|_1 = \frac{1}{1 - \gamma}$, (iii) arises since $\hat{V}^*(s) - \hat{Q}^*(s, \hat{\pi}_c(s)) \leq \varsigma$ by construction of $\hat{\pi}_c$, whereas (iv) is valid since $\varsigma \in [0, \xi]$. The lemma then follows by recognizing that $\hat{\mathbf{V}}^{\pi^*} \leq \hat{\mathbf{V}}^*$ due to the optimality of $\hat{\mathbf{V}}^*$.

EC.2.7. Proof of Lemma 8

We first make the key observation that, with probability at least $1 - \delta$, the following event holds true:

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \quad |\hat{V}(s) - \hat{Q}(s, a) - \varsigma| > \frac{\xi \delta}{2|\mathcal{S}||\mathcal{A}|}. \quad (\text{EC.40})$$

To justify this claim, note that its complementary event satisfies

$$\mathbb{P}\left(\exists (s, a) \in \mathcal{S} \times \mathcal{A}, |\hat{V}(s) - \hat{Q}(s, a) - \varsigma| \leq \frac{\xi \delta}{2|\mathcal{S}||\mathcal{A}|}\right) \leq \sum_{s, a} \mathbb{P}\left(\varsigma \in \left[\hat{V}(s) - \hat{Q}(s, a) \pm \frac{\xi \delta}{2|\mathcal{S}||\mathcal{A}|}\right]\right)$$

$$\leq |\mathcal{S}||\mathcal{A}| \cdot \frac{\xi\delta}{\xi|\mathcal{S}||\mathcal{A}|} = \delta,$$

where the first inequality applies the union bound, and the last line follows since $\varsigma \sim \text{Unif}(0, \xi)$. Here, we abbreviate $[a \pm b] := [a - b, a + b]$. Combining (EC.40) with the assumption (55), we further reach

$$|V(s) - Q(s, a) - \varsigma| \geq |\widehat{V}(s) - \widehat{Q}(s, a) - \varsigma| - 2 \max_{s,a} |Q(s, a) - \widehat{Q}(s, a)| > \frac{\xi\delta}{4|\mathcal{S}||\mathcal{A}|}$$

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

$$V(s) - Q(s, a) < \varsigma - \frac{\xi\delta}{4|\mathcal{S}||\mathcal{A}|} \quad \text{or} \quad V(s) - Q(s, a) > \varsigma + \frac{\xi\delta}{4|\mathcal{S}||\mathcal{A}|} \quad (\text{EC.41})$$

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

We are now prepared to justify the claim of this lemma. To begin with, consider any action $\widehat{a} \in \{a \in \mathcal{A} : Q(s, a) > V(s) - \varsigma\}$. Comparing this with the condition (EC.41), we can easily see that the only possibility is

$$V(s) - Q(s, \widehat{a}) < \varsigma - \frac{\xi\delta}{4|\mathcal{S}||\mathcal{A}|}.$$

Therefore, by invoking a basic decomposition, we ensure that

$$\begin{aligned} \widehat{V}(s) - \widehat{Q}(s, \widehat{a}) &= \widehat{V}(s) - V(s) + V(s) - Q(s, \widehat{a}) + Q(s, \widehat{a}) - \widehat{Q}(s, \widehat{a}) \\ &\leq V(s) - Q(s, \widehat{a}) + |V(s) - \widehat{V}(s)| + |\widehat{Q}(s, \widehat{a}) - Q(s, \widehat{a})| \\ &< \varsigma - \frac{\xi\delta}{4|\mathcal{S}||\mathcal{A}|} + 2 \max_{s,a} |Q(s, a) - \widehat{Q}(s, a)| \\ &\leq \varsigma - \frac{\xi\delta}{4|\mathcal{S}||\mathcal{A}|} + 2 \cdot \frac{\xi\delta}{8|\mathcal{S}||\mathcal{A}|} = \varsigma. \end{aligned}$$

This essentially implies that

$$\{a \in \mathcal{A} : Q(s, a) > V(s) - \varsigma\} \subseteq \{a \in \mathcal{A} : \widehat{Q}(s, a) > \widehat{V}(s) - \varsigma\}. \quad (\text{EC.42})$$

Applying exactly the same argument for any action $\widehat{a} \in \{a \in \mathcal{A} : Q(s, a) \leq V(s) - \varsigma\}$, we can also derive

$$\{a \in \mathcal{A} : \widehat{Q}(s, a) > \widehat{V}(s) - \varsigma\} \subseteq \{a \in \mathcal{A} : Q(s, a) > V(s) - \varsigma\}. \quad (\text{EC.43})$$

These two set inequalities taken collectively establish the lemma.

EC.3. Proofs of auxiliary lemmas: finite-horizon MDPs

To begin, we find it helpful to control the entrywise magnitudes of $\{\mathbf{V}_h^{(l)}\}$. This is accomplished via the following lemma, with the proof postponed to Section EC.3.2.

LEMMA EC.3. *The vectors $\{\mathbf{V}_h^{(l)}\}$ defined in (62) obey*

$$\|\mathbf{V}_h^{(l)}\|_\infty \leq (\sqrt{3H})^l H \quad (\text{EC.44})$$

for all $1 \leq h \leq H$ and all $l \geq 0$.

EC.3.1. Proof of Lemma 9

Consider any l obeying $1 \leq l \leq m := \log_2 H$. By construction (cf. (62)), we see that

$$\begin{aligned}\widehat{\mathbf{V}}_h^{(l)} - \mathbf{V}_h^{(l)} &= \widehat{\mathbf{P}}_{h,\pi} \widehat{\mathbf{V}}_{h+1}^{(l)} - \mathbf{P}_{h,\pi} \mathbf{V}_{h+1}^{(l)} \\ &= \widehat{\mathbf{P}}_{h,\pi} (\widehat{\mathbf{V}}_{h+1}^{(l)} - \mathbf{V}_{h+1}^{(l)}) + (\widehat{\mathbf{P}}_{h,\pi} - \mathbf{P}_{h,\pi}) \mathbf{V}_{h+1}^{(l)},\end{aligned}$$

which connects $\mathbf{V}_h^{(l)}$ and $\widehat{\mathbf{V}}_h^{(l)}$ with $\mathbf{V}_{h+1}^{(l)}$ and $\widehat{\mathbf{V}}_{h+1}^{(l)}$. Apply the above relation recursively and make use of the conditions (61) to arrive at

$$\widehat{\mathbf{V}}_h^{(l)} - \mathbf{V}_h^{(l)} = \sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \widehat{\mathbf{P}}_{i,\pi} (\widehat{\mathbf{P}}_{j,\pi} - \mathbf{P}_{j,\pi}) \mathbf{V}_{j+1}^{(l)},$$

where we adopt the convenient notation and let $\prod_{i=h}^{h-1} \widehat{\mathbf{P}}_{i,\pi} = \mathbf{I}$.

According to the triangle inequality, we can further deduce that

$$\begin{aligned}|\widehat{\mathbf{V}}_h^{(l)} - \mathbf{V}_h^{(l)}| &\leq \sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \widehat{\mathbf{P}}_{i,\pi} \left| (\widehat{\mathbf{P}}_{j,\pi} - \mathbf{P}_{j,\pi}) \mathbf{V}_{j+1}^{(l)} \right| \\ &\leq \sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \widehat{\mathbf{P}}_{i,\pi} \left(\sqrt{\frac{\beta_1}{N}} \sqrt{\text{Var}_{\mathbf{P}_{j,\pi}} [\mathbf{V}_{j+1}^{(l)}]} + \frac{\beta_1 \|\mathbf{V}_{j+1}^{(l)}\|_\infty}{N} \mathbf{1} \right) \\ &= \sqrt{\frac{\beta_1}{N}} \left(\sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \widehat{\mathbf{P}}_{i,\pi} \mathbf{r}_j^{(l+1)} \right) + \sum_{j=h}^{H-1} \frac{\beta_1 \|\mathbf{V}_{j+1}^{(l)}\|_\infty}{N} \mathbf{1} \\ &\leq \sqrt{\frac{\beta_1}{N}} \left(\sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \widehat{\mathbf{P}}_{i,\pi} \mathbf{r}_j^{(l+1)} \right) + \frac{\beta_1 H \max_j \|\mathbf{V}_j^{(l)}\|_\infty}{N} \mathbf{1},\end{aligned}\tag{EC.45}$$

where the second line follows from the assumption (63), and the third line makes use of the definition (62) of $\mathbf{r}_h^{(l)}$ and the elementary identity $\prod_{i=h}^{j-1} \widehat{\mathbf{P}}_{i,\pi} \mathbf{1} = \mathbf{1}$ (since each $\widehat{\mathbf{P}}_{i,\pi}$ is a probability transition matrix). In view of the construction (62), we can also derive recursively that

$$\widehat{\mathbf{V}}_h^{(l)} = \sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \widehat{\mathbf{P}}_{i,\pi} \mathbf{r}_j^{(l)},\tag{EC.46}$$

which combined with (EC.45) yields

$$|\widehat{\mathbf{V}}_h^{(l)} - \mathbf{V}_h^{(l)}| \leq \sqrt{\frac{\beta_1}{N}} \widehat{\mathbf{V}}_h^{(l+1)} + \frac{\beta_1 H \max_j \|\mathbf{V}_j^{(l)}\|_\infty}{N} \mathbf{1}.\tag{EC.47}$$

Further, the above inequality together with the triangle inequality immediately results in the following recursive relation

$$\begin{aligned}\|\widehat{\mathbf{V}}_h^{(l)} - \mathbf{V}_h^{(l)}\|_\infty &\leq \sqrt{\frac{\beta_1}{N}} \|\widehat{\mathbf{V}}_h^{(l+1)}\|_\infty + \frac{\beta_1 H}{N} \max_j \|\mathbf{V}_j^{(l)}\|_\infty \\ &\leq \sqrt{\frac{\beta_1}{N}} \|\widehat{\mathbf{V}}_h^{(l+1)} - \mathbf{V}_h^{(l+1)}\|_\infty + \sqrt{\frac{\beta_1}{N}} \|\mathbf{V}_h^{(l+1)}\|_\infty + \frac{\beta_1 H}{N} \max_j \|\mathbf{V}_j^{(l)}\|_\infty,\end{aligned}$$

thus revealing a useful connection between $\|\widehat{\mathbf{V}}_h^{(l)} - \mathbf{V}_h^{(l)}\|_\infty$ and $\|\widehat{\mathbf{V}}_h^{(l+1)} - \mathbf{V}_h^{(l+1)}\|_\infty$. Applying this relation recursively with a little algebra leads to

$$\begin{aligned} \|\widehat{\mathbf{V}}_h^{(0)} - \mathbf{V}_h^{(0)}\|_\infty &\leq \left(\sqrt{\frac{\beta_1}{N}}\right)^m \|\widehat{\mathbf{V}}_h^{(m)} - \mathbf{V}_h^{(m)}\|_\infty \\ &\quad + \sum_{l=1}^m \left(\sqrt{\frac{\beta_1}{N}}\right)^l \|\mathbf{V}_h^{(l)}\|_\infty + \frac{\beta_1 H}{N} \sum_{l=1}^m \left(\sqrt{\frac{\beta_1}{N}}\right)^{l-1} \max_j \|\mathbf{V}_j^{(l-1)}\|_\infty. \end{aligned} \quad (\text{EC.48})$$

Additionally, it is easily seen from the definition (62) that

$$\mathbf{r}_h^{(l+1)} \leq \mathbf{V}_{h+1}^{(l)} \leq \|\mathbf{V}_{h+1}^{(l)}\|_\infty \mathbf{1},$$

which taken together with (EC.45) and the elementary identity $\prod_{i=h}^{j-1} \widehat{\mathbf{P}}_{i,\pi} \mathbf{1} = \mathbf{1}$ implies that

$$\begin{aligned} \|\widehat{\mathbf{V}}_h^{(m)} - \mathbf{V}_h^{(m)}\|_\infty &\leq \sqrt{\frac{\beta_1}{N}} \max_j \|\mathbf{r}_j^{(m+1)}\|_\infty \left(\sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \widehat{\mathbf{P}}_{i,\pi} \mathbf{1} \right) + \frac{\beta_1 H \max_j \|\mathbf{V}_j^{(m)}\|_\infty}{N} \mathbf{1} \\ &\leq \sqrt{\frac{\beta_1}{N}} \max_j \|\mathbf{V}_j^{(m)}\|_\infty \left(\sum_{j=h}^{H-1} \mathbf{1} \right) + \frac{\beta_1 H \max_j \|\mathbf{V}_j^{(m)}\|_\infty}{N} \mathbf{1} \\ &\leq H \left(\sqrt{\frac{\beta_1}{N}} + \frac{\beta_1}{N} \right) \max_j \|\mathbf{V}_j^{(m)}\|_\infty \mathbf{1}. \end{aligned}$$

Substitution into (EC.48) results in

$$\begin{aligned} \|\widehat{\mathbf{V}}_h^{(0)} - \mathbf{V}_h^{(0)}\|_\infty &\leq H \left(\sqrt{\frac{\beta_1}{N}} \right)^m \left(\sqrt{\frac{\beta_1}{N}} + \frac{\beta_1}{N} \right) \max_j \|\mathbf{V}_j^{(m)}\|_\infty \\ &\quad + \sum_{l=1}^m \left(\sqrt{\frac{\beta_1}{N}} \right)^l \|\mathbf{V}_h^{(l)}\|_\infty + \frac{\beta_1 H}{N} \sum_{l=1}^m \left(\sqrt{\frac{\beta_1}{N}} \right)^{l-1} \max_j \|\mathbf{V}_j^{(l-1)}\|_\infty. \end{aligned} \quad (\text{EC.49})$$

To finish up, it remains to control the terms $\{\|\mathbf{V}_h^{(l)}\|_\infty\}$. Towards this, combining Lemma EC.3 with the above inequality (EC.49) yields

$$\begin{aligned} \|\widehat{\mathbf{V}}_h^{(0)} - \mathbf{V}_h^{(0)}\|_\infty &\leq H^2 \left(\sqrt{\frac{3\beta_1 H}{N}} \right)^m \left(\sqrt{\frac{\beta_1}{N}} + \frac{\beta_1}{N} \right) + \sum_{l=1}^m H \left(\sqrt{\frac{3\beta_1 H}{N}} \right)^l + \sum_{l=1}^m \frac{\beta_1 H^2}{N} \left(\sqrt{\frac{3\beta_1 H}{N}} \right)^{l-1} \\ &= H^2 \left(\sqrt{\frac{3\beta_1 H}{N}} \right)^m \left(\sqrt{\frac{\beta_1}{N}} + \frac{\beta_1}{N} \right) + \left\{ H \sqrt{\frac{3\beta_1 H}{N}} + \frac{\beta_1 H^2}{N} \right\} \sum_{l=1}^m \left(\sqrt{\frac{3\beta_1 H}{N}} \right)^{l-1} \\ &\leq H^2 \left(\sqrt{\frac{3\beta_1 H}{N}} \right)^m \left(\sqrt{\frac{\beta_1}{N}} + \frac{\beta_1}{N} \right) + 2 \left\{ H \sqrt{\frac{3\beta_1 H}{N}} + \frac{\beta_1 H^2}{N} \right\}, \end{aligned} \quad (\text{EC.50})$$

Here, the last inequality holds true since

$$\sum_{l=1}^m \left(\sqrt{\frac{3\beta_1 H}{N}} \right)^{l-1} \leq \frac{1}{1 - \sqrt{\frac{3\beta_1 H}{N}}} \leq 2,$$

provided that $\sqrt{\frac{3\beta_1 H}{N}} \leq 1/2$. Invoking the assumption $N \geq 12H\beta_1$ again and taking $m = \log_2 H$, we have $(\sqrt{\frac{3\beta_1 H}{N}})^m \leq 1/H$. This combined with (EC.51) immediately leads to

$$\|\widehat{\mathbf{V}}_h^{(0)} - \mathbf{V}_h^{(0)}\|_\infty \leq 3 \left\{ H \sqrt{\frac{3\beta_1 H}{N}} + \frac{\beta_1 H^2}{N} \right\} \leq 6H \sqrt{\frac{3\beta_1 H}{N}}, \quad (\text{EC.51})$$

with the proviso that $N \geq 3\beta_1 H$. This concludes the proof.

EC.3.2. Proof of Lemma EC.3

To begin with, it is seen from the notation (28) that for all $1 \leq j \leq H$,

$$\begin{aligned} \text{Var}_{\mathcal{P}_{j,\pi}}(\mathbf{V}_{j+1}^{(l)}) &= \mathcal{P}_{j,\pi}(\mathbf{V}_{j+1}^{(l)} \circ \mathbf{V}_{j+1}^{(l)}) - (\mathcal{P}_{j,\pi} \mathbf{V}_{j+1}^{(l)}) \circ (\mathcal{P}_{j,\pi} \mathbf{V}_{j+1}^{(l)}) \\ &= \mathcal{P}_{j,\pi}(\mathbf{V}_{j+1}^{(l)} \circ \mathbf{V}_{j+1}^{(l)}) - (\mathbf{V}_j^{(l)} - \mathbf{r}_j^{(l)}) \circ (\mathbf{V}_j^{(l)} - \mathbf{r}_j^{(l)}) \\ &= \mathcal{P}_{j,\pi}(\mathbf{V}_{j+1}^{(l)} \circ \mathbf{V}_{j+1}^{(l)}) - \mathbf{V}_j^{(l)} \circ \mathbf{V}_j^{(l)} + 2\mathbf{V}_j^{(l)} \circ \mathbf{r}_j^{(l)} - \mathbf{r}_j^{(l)} \circ \mathbf{r}_j^{(l)} \\ &\leq \mathcal{P}_{j,\pi}(\mathbf{V}_{j+1}^{(l)} \circ \mathbf{V}_{j+1}^{(l)}) - \mathbf{V}_j^{(l)} \circ \mathbf{V}_j^{(l)} + 2\mathbf{V}_j^{(l)} \circ \mathbf{r}_j^{(l)}, \end{aligned} \quad (\text{EC.52})$$

where the second identity makes use of the fact that $\mathbf{V}_j^{(l)} = \mathbf{r}_j^{(l)} + \mathcal{P}_{j,\pi} \mathbf{V}_{j+1}^{(l)}$ (cf. (62)). Moreover, from the construction (62) we can easily derive

$$\mathbf{V}_h^{(l)} = \sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \mathcal{P}_{i,\pi} \mathbf{r}_j^{(l)} = \sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \mathcal{P}_{i,\pi} \sqrt{\text{Var}_{\mathcal{P}_{j,\pi}}[\mathbf{V}_{j+1}^{(l-1)}]}. \quad (\text{EC.53})$$

The above two results taken collectively give

$$\begin{aligned} |\mathbf{V}_h^{(l+1)}| &= \sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \mathcal{P}_{i,\pi} \sqrt{\text{Var}_{\mathcal{P}_{j,\pi}}[\mathbf{V}_{j+1}^{(l)}]} \stackrel{(i)}{\leq} \sum_{j=h}^{H-1} \sqrt{\prod_{i=h}^{j-1} \mathcal{P}_{i,\pi} \text{Var}_{\mathcal{P}_{j,\pi}}[\mathbf{V}_{j+1}^{(l)}]} \\ &\stackrel{(ii)}{\leq} \sqrt{H \sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \mathcal{P}_{i,\pi} \text{Var}_{\mathcal{P}_{j,\pi}}[\mathbf{V}_{j+1}^{(l)}]} \\ &\stackrel{(iii)}{\leq} \sqrt{H \sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \mathcal{P}_{i,\pi} \left[\mathcal{P}_{j,\pi}(\mathbf{V}_{j+1}^{(l)} \circ \mathbf{V}_{j+1}^{(l)}) - \mathbf{V}_j^{(l)} \circ \mathbf{V}_j^{(l)} + 2\mathbf{V}_j^{(l)} \circ \mathbf{r}_j^{(l)} \right]} \\ &= \sqrt{H \sum_{j=h}^{H-1} \left(\prod_{i=h}^j \mathcal{P}_{i,\pi}(\mathbf{V}_{j+1}^{(l)} \circ \mathbf{V}_{j+1}^{(l)}) - \prod_{i=h}^{j-1} \mathcal{P}_{i,\pi}(\mathbf{V}_j^{(l)} \circ \mathbf{V}_j^{(l)}) \right) + 2H \sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \mathcal{P}_{i,\pi}(\mathbf{V}_j^{(l)} \circ \mathbf{r}_j^{(l)})}. \end{aligned}$$

Here, (i) arises from Jensen's inequality; (ii) holds true due to the Cauchy-Schwarz inequality; (iii) follows from (EC.52). By telescoping summation, one further arrives at

$$|\mathbf{V}_h^{(l+1)}| \leq \sqrt{H \left[\prod_{i=h}^{H-1} \mathcal{P}_{i,\pi}(\mathbf{V}_H^{(l)} \circ \mathbf{V}_H^{(l)}) - \mathbf{V}_h^{(l)} \circ \mathbf{V}_h^{(l)} \right] + 2H \max_j \|\mathbf{V}_j^{(l)}\|_\infty \sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \mathcal{P}_{i,\pi} \mathbf{r}_j^{(l)}}$$

$$\begin{aligned}
&\stackrel{\text{(iv)}}{=} \sqrt{H \left[\prod_{i=h}^{H-1} \mathbf{P}_{i,\pi}(\mathbf{V}_H^{(l)} \circ \mathbf{V}_H^{(l)}) - \mathbf{V}_h^{(l)} \circ \mathbf{V}_h^{(l)} \right] + 2H \max_j \|\mathbf{V}_j^{(l)}\|_\infty \mathbf{V}_h^{(l)}} \\
&\leq \sqrt{H \prod_{i=h}^{H-1} \mathbf{P}_{i,\pi}(\mathbf{V}_H^{(l)} \circ \mathbf{V}_H^{(l)}) + 2H \left(\max_j \|\mathbf{V}_j^{(l)}\|_\infty \right)^2 \mathbf{1}} \\
&\stackrel{\text{(v)}}{\leq} \sqrt{3H} \max_j \|\mathbf{V}_j^{(l)}\|_\infty \mathbf{1}.
\end{aligned}$$

Here (iv) invokes the relation $\mathbf{V}_h^{(l)} = \sum_{j=h}^{H-1} \prod_{i=h}^{j-1} \mathbf{P}_{i,\pi} \mathbf{r}_j^{(l)}$ (see (EC.53)); and (v) holds true since

$$\left\| \prod_{i=h}^{H-1} \mathbf{P}_{i,\pi}(\mathbf{V}_H^{(l)} \circ \mathbf{V}_H^{(l)}) \right\|_\infty \leq \left\| \prod_{i=h}^{H-1} \mathbf{P}_{i,\pi} \right\|_1 \left\| (\mathbf{V}_H^{(l)} \circ \mathbf{V}_H^{(l)}) \right\|_\infty \leq \max_j \|\mathbf{V}_j^{(l)}\|_\infty^2.$$

As a consequence, the above inequality allows one to deduce that

$$\max_j \|\mathbf{V}_j^{(l+1)}\|_\infty \leq \sqrt{3H} \max_j \|\mathbf{V}_j^{(l)}\|_\infty,$$

and therefore,

$$\max_j \|\mathbf{V}_j^{(l)}\|_\infty \leq (\sqrt{3H})^l \max_j \|\mathbf{V}_j^{(0)}\|_\infty \leq (\sqrt{3H})^l H,$$

where the last inequality arises from the trivial upper bound $\max_j \|\mathbf{V}_j^{(0)}\|_\infty \leq H$.

Endnotes

1. Here and throughout, the “model” refers to the transition kernel and the rewards of the MDP taken collectively.
2. Note that perturbation is only invoked when running the planning algorithms and does not require collecting new samples.

Acknowledgments

Y. Wei is supported in part by the the NSF grants CCF-2106778, DMS-2147546/2015447 and CAREER award DMS-2143215. Y. Chi is supported in part by the ONR grants N00014-18-1-2142 and N00014-19-1-2404, the NSF grants CCF-1806154, CCF-2007911 and CCF-2106778. Y. Chen is supported in part by the Alfred P. Sloan Research Fellowship, the Google Research Scholar Award, the AFOSR grants FA9550-19-1-0030 and FA9550-22-1-0198, the ONR grant N00014-22-1-2354, and the NSF grants CCF-2221009, CCF-1907661, DMS-2014279, IIS-2218713 and IIS-2218773. We thank Qiwen Cui for pointing out an issue in Section EC.2.5 in an early version of this paper, and thank Shicong Cen, Chen Cheng and Cong Ma for numerous discussions. Part of this work was done while G. Li, Y. Wei and Y. Chen were visiting the Simons Institute for the Theory of Computing.

References

- Agarwal A, Kakade S, Yang LF (2020) Model-based reinforcement learning with a generative model is minimax optimal. *Conference on Learning Theory (COLT)* .
- Azar MG, Munos R, Kappen HJ (2013) Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning* 91(3):325–349.