

---

# Supplement to “Sample-Efficient Reinforcement Learning for Linearly-Parameterized MDPs with a Generative Model”

---

**Bingyan Wang\***  
Princeton University  
bingyanw@princeton.edu

**Yuling Yan\***  
Princeton University  
yulingy@princeton.edu

**Jianqing Fan**  
Princeton University  
jqfan@princeton.edu

## A Notations

In this section we gather the notations that will be used throughout the appendix.

For any vectors  $\mathbf{u} = [u_i]_{i=1}^n \in \mathbb{R}^n$  and  $\mathbf{v} = [v_i]_{i=1}^n \in \mathbb{R}^n$ , let  $\mathbf{u} \circ \mathbf{v} = [u_i v_i]_{i=1}^n$  denote the Hadamard product of  $\mathbf{u}$  and  $\mathbf{v}$ . We slightly abuse notations to use  $\sqrt{\cdot}$  and  $|\cdot|$  to define entry-wise operation, i.e. for any vector  $\mathbf{v} = [v_i]_{i=1}^n$  denote  $\sqrt{\mathbf{v}} := [\sqrt{v_i}]_{i=1}^n$  and  $|\mathbf{v}| := [|v_i|]_{i=1}^n$ . Furthermore, the binary notations  $\leq$  and  $\geq$  are both defined in entry-wise manner, i.e.  $\mathbf{u} \leq \mathbf{v}$  (resp.  $\mathbf{u} \geq \mathbf{v}$ ) means  $u_i \leq v_i$  (resp.  $u_i \geq v_i$ ) for all  $1 \leq i \leq n$ . For a collection of vectors  $\mathbf{v}_1, \dots, \mathbf{v}_m \in \mathbb{R}^n$  with  $\mathbf{v}_i = [v_{i,j}]_{j=1}^n \in \mathbb{R}^n$ , we define the max operator to be  $\max_{1 \leq i \leq m} \mathbf{v}_i := [\max_{1 \leq i \leq m} v_{i,j}]_{j=1}^n$ .

For any matrix  $\mathbf{M} \in \mathbb{R}^{m \times n}$ ,  $\|\mathbf{M}\|_1$  is defined as the largest row-wise  $\ell_1$  norm of  $\mathbf{M}$ , i.e.  $\|\mathbf{M}\|_1 := \max_i \sum_j |M_{i,j}|$ . In addition, we define  $\mathbf{1}$  to be a vector with all the entries being 1, and  $\mathbf{I}$  be the identity matrix. To express the probability transition function  $P$  in matrix form, we define the matrix  $\mathbf{P} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$  to be a matrix whose  $(s, a)$ -th row  $\mathbf{P}_{s,a}$  corresponds to  $P(\cdot | s, a)$ . In addition, we define  $\mathbf{P}^\pi$  to be the probability transition matrix induced by policy  $\pi$ , i.e.  $\mathbf{P}_{(s,a),(s',a')}^\pi = \mathbf{P}_{s,a}(s') \mathbb{1}_{\pi(s')=a'}$  for all state-action pairs  $(s, a)$  and  $(s', a')$ . We define  $\pi_t$  to be the policy induced by  $Q_t$ , i.e.  $Q_t(s, \pi_t(s)) = \max_a Q_t(s, a)$  for all  $s \in \mathcal{S}$ . Furthermore, we denote the reward function  $r$  by vector  $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ , i.e. the  $(s, a)$ -th element of  $\mathbf{r}$  equals  $r(s, a)$ . In the same manner, we define  $\mathbf{V}^\pi \in \mathbb{R}^{|\mathcal{S}|}$ ,  $\mathbf{V}^* \in \mathbb{R}^{|\mathcal{S}|}$ ,  $\mathbf{V}_t \in \mathbb{R}^{|\mathcal{S}|}$ ,  $\mathbf{Q}^\pi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ ,  $\mathbf{Q}^* \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$  and  $Q_t \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$  to represent  $V^\pi$ ,  $V^*$ ,  $V_t$ ,  $Q^\pi$ ,  $Q^*$  and  $Q_t$  respectively. By using these notations, we can rewrite the Bellman equation as

$$\mathbf{Q}^\pi = \mathbf{r} + \gamma \mathbf{P} \mathbf{V}^\pi = \mathbf{r} + \gamma \mathbf{P}^\pi \mathbf{Q}^\pi. \quad (11)$$

Further, for any vector  $\mathbf{V} \in \mathbb{R}^{|\mathcal{S}|}$ , let  $\text{Var}_{\mathbf{P}}(\mathbf{V}) \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$  be

$$\text{Var}_{\mathbf{P}}(\mathbf{V}) := \mathbf{P}(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}\mathbf{V}) \circ (\mathbf{P}\mathbf{V}), \quad (12)$$

and define  $\text{Var}_{\mathbf{P}_{s,a}}(\mathbf{V}) \in \mathbb{R}$  to be

$$\text{Var}_{\mathbf{P}_{s,a}}(\mathbf{V}) := \mathbf{P}_{s,a}(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}_{s,a}\mathbf{V})^2, \quad (13)$$

where  $\mathbf{P}_{s,a}$  is the  $(s, a)$ -th row of  $\mathbf{P}$ .

Next, we reconsider Assumption 1. For any state-action pair  $(s, a)$ , we define vector  $\boldsymbol{\lambda}(s, a) \in \mathbb{R}^K$  (resp.  $\boldsymbol{\phi}(s, a) \in \mathbb{R}^K$ ) with  $\boldsymbol{\lambda}(s, a) = [\lambda_i(s, a)]_{i=1}^K$  (resp.  $\boldsymbol{\phi}(s, a) = [\phi_i(s, a)]_{i=1}^K$ ) and matrix

---

\*Equal contribution.

$\Lambda \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times K}$  (resp.  $\Phi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times K}$ ) whose  $(s, a)$ -th row corresponds to  $\lambda(s, a)^\top$  (resp.  $\phi(s, a)^\top$ ). Define vector  $\psi(s, a) \in \mathbb{R}^K$  with  $\psi(s, a) = [\psi_i(s, a)]_{i=1}^K$  and matrix  $\Psi \in \mathbb{R}^{K \times |\mathcal{S}|}$  whose  $(s, a)$ -th column corresponds to  $\psi(s, a)^\top$ . Further, let  $P_{\mathcal{K}} \in \mathbb{R}^{K \times |\mathcal{S}|}$  (resp.  $\Phi_{\mathcal{K}} \in \mathbb{R}^{K \times K}$ ) to be a submatrix of  $P$  (resp.  $\Phi$ ) formed by concatenating the rows  $\{P_{s,a}, (s, a) \in \mathcal{K}\}$  (resp.  $\{\Phi_{s,a}, (s, a) \in \mathcal{K}\}$ ). By using the previous notations, we can express the relations in Definition 1 and Assumption 1 as  $P_{\mathcal{K}} = \Phi_{\mathcal{K}} \Psi$ ,  $P = \Phi \Psi$  and  $\Phi = \Lambda \Phi_{\mathcal{K}}$ . Note that Assumption 1 suggests  $\Phi_{\mathcal{K}}$  is invertible. Taking these equations collectively yields

$$P = \Phi \Psi = \Phi \Phi_{\mathcal{K}}^{-1} P_{\mathcal{K}} = \Lambda \Phi_{\mathcal{K}} \Phi_{\mathcal{K}}^{-1} P_{\mathcal{K}} = \Lambda P_{\mathcal{K}}, \quad (14)$$

which is reminiscent of the anchor word condition in topic modelling [2]. In addition, for each iteration  $t$ , we denote the collected samples as  $\{s_t(s, a)\}_{(s,a) \in \mathcal{K}}$  and define a matrix  $\widehat{P}_{\mathcal{K}}^{(t)} \in \{0, 1\}^{K \times |\mathcal{S}|}$  to be

$$\widehat{P}_{\mathcal{K}}^{(t)}((s, a), s') := \begin{cases} 1, & \text{if } s' = s_t(s, a) \\ 0, & \text{otherwise} \end{cases} \quad (15)$$

for any  $(s, a) \in \mathcal{K}$  and  $s' \in \mathcal{S}$ . Further, we define  $\widehat{P}_t = \Lambda \widehat{P}_{\mathcal{K}}^{(t)}$ . Then it is obvious to see that  $\widehat{P}_t$  has nonnegative entries and unit  $\ell_1$  norm for each row due to Assumption 1, i.e.  $\|\widehat{P}_t\|_1 = 1$ .

## B Analysis of model-based RL (Proof of Theorem 1)

In this section, we will provide complete proof for Theorem 1. As a matter of fact, our proof strategy here justifies a more general version of Theorem 1 that accounts for model misspecification, as stated below.

**Theorem 3.** *Suppose that  $\delta > 0$  and  $\varepsilon \in (0, (1 - \gamma)^{-1/2}]$ . Assume that there exists a probability transition model  $\widetilde{P}$  obeying Definition 1 and Assumption 1 with feature vectors  $\{\phi(s, a)\}_{(s,a) \in \mathcal{S} \times \mathcal{A}} \subset \mathbb{R}^K$  and anchor state-action pairs  $\mathcal{K}$  such that*

$$\|\widetilde{P} - P\|_1 \leq \xi$$

for some  $\xi \geq 0$ . Let  $\widehat{\pi}$  be the policy returned by Algorithm 1. Assume that

$$N \geq \frac{C \log(K / ((1 - \gamma) \delta))}{(1 - \gamma)^3 \varepsilon^2} \quad (16)$$

for some sufficiently large constant  $C > 0$ . Then with probability exceeding  $1 - \delta$ ,

$$Q^*(s, a) - Q^{\widehat{\pi}}(s, a) \leq \varepsilon + \frac{4\varepsilon_{\text{opt}}}{1 - \gamma} + \frac{22\xi}{(1 - \gamma)^2}, \quad (17)$$

for every state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .

Theorem 3 subsumes Theorem 1 as a special case with  $\xi = 0$ . The remainder of this section is devoted to proving Theorem 3.

### B.1 Proof of Theorem 3

The error  $Q^{\widehat{\pi}} - Q^*$  can be decomposed as

$$\begin{aligned} Q^{\widehat{\pi}} - Q^* &= Q^{\widehat{\pi}} - \widehat{Q}^{\widehat{\pi}} + \widehat{Q}^{\widehat{\pi}} - \widehat{Q}^* + \widehat{Q}^* - Q^* \\ &\geq Q^{\widehat{\pi}} - \widehat{Q}^{\widehat{\pi}} + \widehat{Q}^{\widehat{\pi}} - \widehat{Q}^* + \widehat{Q}^{\pi^*} - Q^* \\ &\geq - \left( \|Q^{\widehat{\pi}} - \widehat{Q}^{\widehat{\pi}}\|_{\infty} + \|\widehat{Q}^{\widehat{\pi}} - \widehat{Q}^*\|_{\infty} + \|\widehat{Q}^{\pi^*} - Q^*\|_{\infty} \right) \mathbf{1}. \end{aligned} \quad (18)$$

For policy  $\widehat{\pi}$  satisfying the condition in Theorem 1, we have  $\|\widehat{Q}^{\widehat{\pi}} - \widehat{Q}^*\|_{\infty} \leq \varepsilon_{\text{opt}}$ . It boils down to control  $\|Q^{\widehat{\pi}} - \widehat{Q}^{\widehat{\pi}}\|_{\infty}$  and  $\|\widehat{Q}^{\pi^*} - Q^*\|_{\infty}$ .

To begin with, we can use (11) to further decompose  $\|Q^{\hat{\pi}} - \hat{Q}^{\hat{\pi}}\|_{\infty}$  as

$$\begin{aligned}
\|Q^{\hat{\pi}} - \hat{Q}^{\hat{\pi}}\|_{\infty} &= \left\| (I - \gamma P^{\hat{\pi}})^{-1} r - (I - \gamma \hat{P}^{\hat{\pi}})^{-1} r \right\|_{\infty} \\
&= \left\| (I - \gamma P^{\hat{\pi}})^{-1} \left[ (I - \gamma \hat{P}^{\hat{\pi}}) - (I - \gamma P^{\hat{\pi}}) \right] \hat{Q}^{\hat{\pi}} \right\|_{\infty} \\
&= \left\| \gamma (I - \gamma P^{\hat{\pi}})^{-1} (P - \hat{P}) \hat{V}^{\hat{\pi}} \right\|_{\infty} \\
&\leq \left\| \gamma (I - \gamma P^{\hat{\pi}})^{-1} (P - \hat{P}) \hat{V}^{\star} \right\|_{\infty} + \left\| \gamma (I - \gamma P^{\hat{\pi}})^{-1} (P - \hat{P}) (\hat{V}^{\hat{\pi}} - \hat{V}^{\star}) \right\|_{\infty} \\
&\leq \left\| \gamma (I - \gamma P^{\hat{\pi}})^{-1} (P - \hat{P}) \hat{V}^{\star} \right\|_{\infty} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}. \tag{19}
\end{aligned}$$

Here the last inequality is due to

$$\begin{aligned}
&\left\| \gamma (I - \gamma P^{\hat{\pi}})^{-1} (P - \hat{P}) (\hat{V}^{\hat{\pi}} - \hat{V}^{\star}) \right\|_{\infty} \\
&\leq \gamma \left\| (I - \gamma P^{\hat{\pi}})^{-1} \right\|_1 \left\| (P - \hat{P}) (\hat{V}^{\hat{\pi}} - \hat{V}^{\star}) \right\|_{\infty} \\
&\leq \gamma \left\| (I - \gamma P^{\hat{\pi}})^{-1} \right\|_1 (\|P\|_1 + \|\hat{P}\|_1) \|\hat{V}^{\hat{\pi}} - \hat{V}^{\star}\|_{\infty} \\
&\leq \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma},
\end{aligned}$$

where we use the fact that  $\|(I - \gamma P^{\hat{\pi}})^{-1}\|_1 \leq 1/(1-\gamma)$  and  $\|P\|_1 = \|\hat{P}\|_1 = 1$ .

Similarly, for the term  $\|\hat{Q}^{\pi^{\star}} - Q^{\star}\|_{\infty}$  in (18), we have

$$\begin{aligned}
\|\hat{Q}^{\pi^{\star}} - Q^{\star}\|_{\infty} &= \left\| \gamma (I - \gamma P^{\pi^{\star}})^{-1} (P - \hat{P}) \hat{V}^{\pi^{\star}} \right\|_{\infty} \\
&\leq \left\| \gamma (I - \gamma P^{\pi^{\star}})^{-1} (P - \hat{P}) \hat{V}^{\pi^{\star}} \right\|_{\infty}. \tag{20}
\end{aligned}$$

As can be seen from (19) and (20), it boils down to bound  $\|(P - \hat{P})\hat{V}^{\star}\|$  and  $\|(P - \hat{P})\hat{V}^{\pi^{\star}}\|$ . We have the following lemma.

**Lemma 1.** *With probability exceeding  $1 - \delta$ , one has*

$$\begin{aligned}
\left| (P - \hat{P})_{s,a} \hat{V}^{\star} \right| &\leq \frac{10\xi}{1-\gamma} + 4\sqrt{\frac{2\log(4K/\delta)}{N}} + \frac{4\log(8K/((1-\gamma)\delta))}{(1-\gamma)N} \\
&\quad + \sqrt{\frac{4\log(8K/((1-\gamma)\delta))}{N}} \sqrt{\text{Var}_{P_{s,a}}(\hat{V}^{\star})}, \tag{21}
\end{aligned}$$

$$\begin{aligned}
\left| (P - \hat{P})_{s,a} \hat{V}^{\pi^{\star}} \right| &\leq \frac{10\xi}{1-\gamma} + 4\sqrt{\frac{2\log(4K/\delta)}{N}} + \frac{4\log(8K/((1-\gamma)\delta))}{(1-\gamma)N} \\
&\quad + \sqrt{\frac{4\log(8K/((1-\gamma)\delta))}{N}} \sqrt{\text{Var}_{P_{s,a}}(\hat{V}^{\pi^{\star}})}. \tag{22}
\end{aligned}$$

*Proof.* See Appendix B.2. □

Applying (21) to (19) reveals that

$$\begin{aligned} \left\| \mathbf{Q}^{\hat{\pi}} - \hat{\mathbf{Q}}^{\hat{\pi}} \right\|_{\infty} &\leq \sqrt{\frac{4 \log(8K/((1-\gamma)\delta))}{N}} \left\| \gamma \left( \mathbf{I} - \gamma \mathbf{P}^{\hat{\pi}} \right)^{-1} \sqrt{\text{Var}_{\mathcal{P}_{s,a}} \left( \hat{\mathbf{V}}^{\star} \right)} \right\|_{\infty} \\ &\quad + \frac{\gamma}{1-\gamma} \left[ 4 \sqrt{\frac{2 \log(4K/\delta)}{N}} + \frac{4 \log(8K/((1-\gamma)\delta))}{(1-\gamma)N} \right] \\ &\quad + \frac{10\gamma\xi}{(1-\gamma)^2} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}. \end{aligned} \quad (23)$$

For the first term, one has

$$\begin{aligned} \sqrt{\text{Var}_{\mathcal{P}_{s,a}} \left( \hat{\mathbf{V}}^{\star} \right)} &\leq \sqrt{\text{Var}_{\mathcal{P}_{s,a}} \left( \mathbf{V}^{\hat{\pi}} \right)} + \sqrt{\text{Var}_{\mathcal{P}_{s,a}} \left( \mathbf{V}^{\hat{\pi}} - \hat{\mathbf{V}}^{\hat{\pi}} \right)} + \sqrt{\text{Var}_{\mathcal{P}_{s,a}} \left( \hat{\mathbf{V}}^{\hat{\pi}} - \hat{\mathbf{V}}^{\star} \right)} \\ &\leq \sqrt{\text{Var}_{\mathcal{P}_{s,a}} \left( \mathbf{V}^{\hat{\pi}} \right)} + \left\| \mathbf{V}^{\hat{\pi}} - \hat{\mathbf{V}}^{\hat{\pi}} \right\|_{\infty} + \varepsilon_{\text{opt}} \\ &\leq \sqrt{\text{Var}_{\mathcal{P}_{s,a}} \left( \mathbf{V}^{\hat{\pi}} \right)} + \left\| \mathbf{Q}^{\hat{\pi}} - \hat{\mathbf{Q}}^{\hat{\pi}} \right\|_{\infty} + \varepsilon_{\text{opt}}, \end{aligned}$$

where the first inequality comes from the fact that  $\sqrt{\text{Var}(X+Y)} \leq \sqrt{\text{Var}(X)} + \sqrt{\text{Var}(Y)}$  for any random variables  $X$  and  $Y$ . It follows that

$$\begin{aligned} &\left\| \gamma \left( \mathbf{I} - \gamma \mathbf{P}^{\hat{\pi}} \right)^{-1} \sqrt{\text{Var}_{\mathcal{P}_{s,a}} \left( \hat{\mathbf{V}}^{\star} \right)} \right\|_{\infty} \\ &\leq \left\| \gamma \left( \mathbf{I} - \gamma \mathbf{P}^{\hat{\pi}} \right)^{-1} \sqrt{\text{Var}_{\mathcal{P}_{s,a}} \left( \mathbf{V}^{\hat{\pi}} \right)} \right\|_{\infty} + \frac{\gamma}{1-\gamma} \left( \left\| \mathbf{Q}^{\hat{\pi}} - \hat{\mathbf{Q}}^{\hat{\pi}} \right\|_{\infty} + \varepsilon_{\text{opt}} \right) \\ &\leq \gamma \sqrt{\frac{2}{(1-\gamma)^3}} + \frac{\gamma}{1-\gamma} \left( \left\| \mathbf{Q}^{\hat{\pi}} - \hat{\mathbf{Q}}^{\hat{\pi}} \right\|_{\infty} + \varepsilon_{\text{opt}} \right), \end{aligned} \quad (24)$$

where the second inequality utilizes [3, Lemma 7].

Plugging (24) into (23) yields

$$\begin{aligned} \left\| \mathbf{Q}^{\hat{\pi}} - \hat{\mathbf{Q}}^{\hat{\pi}} \right\|_{\infty} &\leq \sqrt{\frac{4 \log(8K/((1-\gamma)\delta))}{N}} \left[ \gamma \sqrt{\frac{2}{(1-\gamma)^3}} + \frac{\gamma}{1-\gamma} \left( \left\| \mathbf{Q}^{\hat{\pi}} - \hat{\mathbf{Q}}^{\hat{\pi}} \right\|_{\infty} + \varepsilon_{\text{opt}} \right) \right] \\ &\quad + \frac{\gamma}{1-\gamma} \left[ 4 \sqrt{\frac{2 \log(4K/\delta)}{N}} + \frac{4 \log(8K/((1-\gamma)\delta))}{(1-\gamma)N} \right] + \frac{10\gamma\xi}{(1-\gamma)^2} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}. \end{aligned}$$

Then we can rearrange terms to obtain

$$\left\| \mathbf{Q}^{\hat{\pi}} - \hat{\mathbf{Q}}^{\hat{\pi}} \right\|_{\infty} \leq 10\gamma \sqrt{\frac{\log(8K/((1-\gamma)\delta))}{N(1-\gamma)^3}} + \frac{11\gamma\xi}{(1-\gamma)^2} + \frac{3\gamma\varepsilon_{\text{opt}}}{1-\gamma} \quad (25)$$

as long as  $N \geq \tilde{C} \log(8K/((1-\gamma)\delta))/(1-\gamma)^2$  for some sufficiently large constant  $\tilde{C} > 0$ .

In a similar vein, we can use (20) and (22) to obtain that

$$\left\| \hat{\mathbf{Q}}^{\pi^{\star}} - \mathbf{Q}^{\star} \right\|_{\infty} \leq 10\gamma \sqrt{\frac{\log(8K/((1-\gamma)\delta))}{N(1-\gamma)^3}} + \frac{11\gamma\xi}{(1-\gamma)^2}. \quad (26)$$

Finally, we can substitute (25) and (26) into (18) to achieve

$$\mathbf{Q}^{\hat{\pi}} - \mathbf{Q}^{\star} \geq - \left( 20\gamma \sqrt{\frac{\log(8K/((1-\gamma)\delta))}{N(1-\gamma)^3}} + \frac{22\gamma\xi}{(1-\gamma)^2} + \frac{3\gamma\varepsilon_{\text{opt}}}{1-\gamma} + \varepsilon_{\text{opt}} \right) \mathbf{1}.$$

This result implies that

$$\mathbf{Q}^{\hat{\pi}} \geq \mathbf{Q}^{\star} - \left( \varepsilon + \frac{22\xi}{(1-\gamma)^2} + \frac{4\varepsilon_{\text{opt}}}{1-\gamma} \right) \mathbf{1},$$

as long as

$$N \geq \frac{C \log(8K/((1-\gamma)\delta))}{(1-\gamma)^3 \varepsilon^2},$$

for some sufficiently large constant  $C > 0$ .

## B.2 Proof of Lemma 1

To prove this theorem, we invoke the idea of  $s$ -absorbing MDP proposed by [1]. For a state  $s \in \mathcal{S}$  and a scalar  $u$ , we define a new MDP  $M_{s,u}$  to be identical to  $M$  on all the other states except  $s$ ; on state  $s$ ,  $M_{s,u}$  is absorbing such that  $P_{M_{s,u}}(s|s, a) = 1$  and  $r_{M_{s,u}}(s, a) = (1 - \gamma)u$  for all  $a \in \mathcal{A}$ . More formally, we define  $P_{M_{s,u}}$  and  $r_{M_{s,u}}$  as

$$\begin{aligned} P_{M_{s,u}}(s|s, a) &= 1, & r_{M_{s,u}}(s, a) &= (1 - \gamma)u, & \text{for all } a \in \mathcal{A}, \\ P_{M_{s,u}}(\cdot|s', a') &= P(\cdot|s', a'), & r_{M_{s,u}}(s, a) &= r(s, a), & \text{for all } s' \neq s \text{ and } a' \in \mathcal{A}. \end{aligned}$$

To streamline notations, we will use  $\mathbf{V}_{s,u}^\pi \in \mathbb{R}^{|\mathcal{S}|}$  and  $\mathbf{V}_{s,u}^* \in \mathbb{R}^{|\mathcal{S}|}$  to denote the value function of  $M_{s,u}$  under policy  $\pi$  and the optimal value function of  $M_{s,u}$  respectively. Furthermore, we denote by  $\widehat{M}_{s,u}$  the MDP whose probability transition kernel is identical to  $\widehat{P}$  at all states except that state  $s$  is absorbing. Similar as before, we use  $\widehat{\mathbf{V}}_{s,u}^* \in \mathbb{R}^{|\mathcal{S}|}$  to denote the optimal value function under  $\widehat{M}_{s,u}$ . The construction of this collection of auxiliary MDPs will facilitate our analysis by decoupling the statistical dependency between  $\widehat{P}$  and  $\widehat{\pi}^*$ .

To begin with, we can decompose the quantity of interest as

$$\begin{aligned} \left| \left( \mathbf{P} - \widehat{\mathbf{P}} \right)_{s,a} \widehat{\mathbf{V}}^* \right| &= \left| \left( \mathbf{P} - \widehat{\mathbf{P}} \right)_{s,a} \left( \widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^* + \widehat{\mathbf{V}}_{s,u}^* \right) \right| \\ &\leq \left| \left( \mathbf{P} - \widehat{\mathbf{P}} \right)_{s,a} \widehat{\mathbf{V}}_{s,u}^* \right| + \left| \left( \mathbf{P} - \widehat{\mathbf{P}} \right)_{s,a} \left( \widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^* \right) \right| \\ &\stackrel{(i)}{\leq} \left| \left( \mathbf{P} - \widetilde{\mathbf{P}} \right)_{s,a} \widehat{\mathbf{V}}_{s,u}^* \right| + \left| \lambda(s, a) \left( \widetilde{\mathbf{P}}_{\mathcal{K}} - \mathbf{P}_{\mathcal{K}} \right) \widehat{\mathbf{V}}_{s,u}^* \right| \\ &\quad + \left| \lambda(s, a) \left( \mathbf{P}_{\mathcal{K}} - \widetilde{\mathbf{P}}_{\mathcal{K}} \right) \widehat{\mathbf{V}}_{s,u}^* \right| + \left( \|\mathbf{P}_{s,a}\|_1 + \|\widehat{\mathbf{P}}_{s,a}\|_1 \right) \|\widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^*\|_\infty \\ &\leq \left\| \left( \mathbf{P} - \widetilde{\mathbf{P}} \right)_{s,a} \right\|_1 \|\widehat{\mathbf{V}}_{s,u}^*\|_\infty + \|\lambda(s, a)\|_1 \cdot \left\| \left( \widetilde{\mathbf{P}}_{\mathcal{K}} - \mathbf{P}_{\mathcal{K}} \right) \widehat{\mathbf{V}}_{s,u}^* \right\|_\infty \\ &\quad + \|\lambda(s, a)\|_1 \cdot \left\| \left( \mathbf{P}_{\mathcal{K}} - \widetilde{\mathbf{P}}_{\mathcal{K}} \right) \widehat{\mathbf{V}}_{s,u}^* \right\|_\infty + 2 \|\widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^*\|_\infty \\ &\stackrel{(ii)}{\leq} \frac{2\xi}{1 - \gamma} + \max_{(s,a) \in \mathcal{K}} \left| \left( \mathbf{P} - \widehat{\mathbf{P}} \right)_{s,a} \widehat{\mathbf{V}}_{s,u}^* \right| + 2 \|\widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^*\|_\infty, \end{aligned} \quad (27)$$

where (i) makes use of  $\widetilde{\mathbf{P}}_{s,a} = \lambda(s, a)\widetilde{\mathbf{P}}_{\mathcal{K}}$  and  $\widehat{\mathbf{P}}_{s,a} = \lambda(s, a)\widehat{\mathbf{P}}_{\mathcal{K}}$ ; (ii) depends on  $\|\mathbf{P} - \widetilde{\mathbf{P}}\|_1 \leq \xi$ ,  $\|\lambda(s, a)\|_1 = 1$  and  $\|\widehat{\mathbf{V}}_{s,u}^*\|_\infty \leq (1 - \gamma)^{-1}$ . For each state  $s$ , the value of  $u$  will be selected from a set  $\mathcal{U}_s$ . The choice of  $\mathcal{U}_s$  will be specified later. Then for some fixed  $u$  in  $\mathcal{U}_s$  and fixed state-action pair  $(s, a) \in \mathcal{K}$ , due to the independence between  $\widehat{\mathbf{P}}_{s,a}$  and  $\widehat{\mathbf{V}}_{s,u}^*$ , we can apply Bernstein's inequality (cf. [5, Theorem 2.8.4]) conditional on  $\widehat{\mathbf{V}}_{s,u}^*$  to reveal that with probability greater than  $1 - \delta/2$ ,

$$\left| \left( \mathbf{P} - \widehat{\mathbf{P}} \right)_{s,a} \widehat{\mathbf{V}}_{s,u}^* \right| \leq \sqrt{\frac{2 \log(4/\delta)}{N} \text{Var}_{\mathbf{P}_{s,a}} \left( \widehat{\mathbf{V}}_{s,u}^* \right)} + \frac{2 \log(4/\delta)}{3(1 - \gamma)N}. \quad (28)$$

Invoking the union bound over all the  $K$  state-action pairs of  $\mathcal{K}$  and all the possible values of  $u$  in  $\mathcal{U}_s$  demonstrate that with probability greater than  $1 - \delta/2$ ,

$$\left| \left( \mathbf{P} - \widehat{\mathbf{P}} \right)_{s,a} \widehat{\mathbf{V}}_{s,u}^* \right| \leq \sqrt{\frac{2 \log(4K|\mathcal{U}_s|/\delta)}{N} \text{Var}_{\mathbf{P}_{s,a}} \left( \widehat{\mathbf{V}}_{s,u}^* \right)} + \frac{2 \log(4K|\mathcal{U}_s|/\delta)}{3(1 - \gamma)N}, \quad (29)$$

holds for all state-action pair  $(s, a) \in \mathcal{K}$  and all  $u \in \mathcal{U}_s$ . Here,  $\text{Var}_{\mathbf{P}_{s,a}}(\cdot)$  is defined in (13). Then we observe that

$$\begin{aligned} \sqrt{\text{Var}_{\mathbf{P}_{s,a}} \left( \widehat{\mathbf{V}}_{s,u}^* \right)} &\leq \sqrt{\text{Var}_{\mathbf{P}_{s,a}} \left( \widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^* \right)} + \sqrt{\text{Var}_{\mathbf{P}_{s,a}} \left( \widehat{\mathbf{V}}^* \right)} \\ &\leq \|\widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^*\|_\infty + \sqrt{\text{Var}_{\mathbf{P}_{s,a}} \left( \widehat{\mathbf{V}}^* \right)} \\ &\leq \left| \widehat{\mathbf{V}}^*(s) - u \right| + \sqrt{\text{Var}_{\mathbf{P}_{s,a}} \left( \widehat{\mathbf{V}}^* \right)}, \end{aligned} \quad (30)$$

where (i) is due to  $\sqrt{\text{Var}_{\mathcal{P}_{s,a}}(\mathbf{V}_1 + \mathbf{V}_2)} \leq \sqrt{\text{Var}_{\mathcal{P}_{s,a}}(\mathbf{V}_1)} + \sqrt{\text{Var}_{\mathcal{P}_{s,a}}(\mathbf{V}_2)}$  and (ii) holds since

$$\left\| \widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^* \right\|_{\infty} = \left\| \widehat{\mathbf{V}}_{s,\widehat{\mathbf{V}}^*(s)}^* - \widehat{\mathbf{V}}_{s,u}^* \right\|_{\infty} \leq \left| \widehat{\mathbf{V}}^*(s) - u \right|, \quad (31)$$

whose proof can be found in [1, Lemma 8 and 9].

By substituting (29), (30) and (31) into (27), we arrive at

$$\begin{aligned} \left| (\mathbf{P} - \widehat{\mathbf{P}})_{s,a} \widehat{\mathbf{V}}^* \right| &\leq \frac{2\xi}{1-\gamma} + \left| \widehat{\mathbf{V}}^*(s) - u \right| \left( 2 + \sqrt{\frac{2 \log(4K |\mathcal{U}_s| / \delta)}{N}} \right) \\ &\quad + \sqrt{\frac{2 \log(4K |\mathcal{U}_s| / \delta)}{N}} \sqrt{\text{Var}_{\mathcal{P}_{s,a}}(\widehat{\mathbf{V}}^*)} + \frac{2 \log(4K |\mathcal{U}_s| / \delta)}{3(1-\gamma)N}. \end{aligned} \quad (32)$$

Then it boils down to determining  $\mathcal{U}_s$ . The coarse bounds of  $\widehat{\mathbf{Q}}^{\pi^*}$  and  $\widehat{\mathbf{Q}}^*$  in the following lemma provide a guidance on the choice of  $\mathcal{U}_s$ .

**Lemma 2.** For  $\delta \in (0, 1)$ , with probability exceeding  $1 - \delta/2$  one has

$$\left\| \mathbf{Q}^* - \widehat{\mathbf{Q}}^{\pi^*} \right\|_{\infty} \leq \frac{\gamma}{1-\gamma} \sqrt{\frac{\log(4K/\delta)}{2N(1-\gamma)^2}} + \frac{2\gamma\xi}{(1-\gamma)^2}, \quad (33)$$

$$\left\| \mathbf{Q}^* - \widehat{\mathbf{Q}}^* \right\|_{\infty} \leq \frac{\gamma}{1-\gamma} \sqrt{\frac{\log(4K/\delta)}{2N(1-\gamma)^2}} + \frac{2\gamma\xi}{(1-\gamma)^2}. \quad (34)$$

*Proof.* See Appendix B.3. □

This inspires us to choose  $\mathcal{U}_s$  to be the set consisting of equidistant points in  $[\mathbf{V}^*(s) - R(\delta), \mathbf{V}^*(s) + R(\delta)]$  with  $|\mathcal{U}_s| = \lceil 1/(1-\gamma)^2 \rceil$  and

$$R(\delta) := \frac{\gamma}{1-\gamma} \sqrt{\frac{\log(4K/\delta)}{2N(1-\gamma)^2}} + \frac{2\gamma\xi}{(1-\gamma)^2}.$$

Since  $\|\mathbf{V}^* - \widehat{\mathbf{V}}^*\|_{\infty} \leq \|\mathbf{Q}^* - \widehat{\mathbf{Q}}^*\|_{\infty}$ , Lemma 2 implies that  $\widehat{\mathbf{V}}^*(s) \in [\mathbf{V}^*(s) - R(\delta), \mathbf{V}^*(s) + R(\delta)]$  with probability over  $1 - \delta/2$ . Hence, we have

$$\min_{u \in \mathcal{U}_s} \left| \widehat{\mathbf{V}}^*(s) - u \right| \leq \frac{2R(\delta)}{|\mathcal{U}_s| + 1} \leq 2\gamma \sqrt{\frac{2 \log(4K/\delta)}{N}} + 4\gamma\xi. \quad (35)$$

Consequently, with probability exceeding  $1 - \delta$ , one has

$$\begin{aligned} \left| (\mathbf{P} - \widehat{\mathbf{P}})_{s,a} \widehat{\mathbf{V}}^* \right| &\stackrel{(i)}{\leq} \frac{2\xi}{1-\gamma} + \min_{u \in \mathcal{U}_s} \left| \widehat{\mathbf{V}}^*(s) - u \right| \left( 2 + \sqrt{\frac{2 \log(4K |\mathcal{U}_s| / \delta)}{N}} \right) \\ &\quad + \sqrt{\frac{2 \log(4K |\mathcal{U}_s| / \delta)}{N}} \sqrt{\text{Var}_{\mathcal{P}_{s,a}}(\widehat{\mathbf{V}}^*)} + \frac{2 \log(4K |\mathcal{U}_s| / \delta)}{3(1-\gamma)N} \\ &\stackrel{(ii)}{\leq} \frac{2\xi}{1-\gamma} + \left( 2\gamma \sqrt{\frac{2 \log(4K/\delta)}{N}} + 4\gamma\xi \right) \left( 2 + \sqrt{\frac{4 \log(8K / ((1-\gamma)\delta))}{N}} \right) \\ &\quad + \sqrt{\frac{4 \log(8K / ((1-\gamma)\delta))}{N}} \sqrt{\text{Var}_{\mathcal{P}_{s,a}}(\widehat{\mathbf{V}}^*)} + \frac{2 \log(8K / ((1-\gamma)\delta))}{3(1-\gamma)N} \\ &\leq \frac{10\xi}{1-\gamma} + 4\sqrt{\frac{2 \log(4K/\delta)}{N}} + \frac{4 \log(8K / ((1-\gamma)\delta))}{(1-\gamma)N} \\ &\quad + \sqrt{\frac{4 \log(8K / ((1-\gamma)\delta))}{N}} \sqrt{\text{Var}_{\mathcal{P}_{s,a}}(\widehat{\mathbf{V}}^*)}, \end{aligned}$$

where (i) follows from (32) and (ii) utilizes (35). This finishes the proof for the first inequality. The second inequality can be proved in a similar way and is omitted here for brevity.

### B.3 Proof of Lemma 2

To begin with, one has

$$\begin{aligned}
\|(\hat{\mathbf{P}} - \mathbf{P}) \mathbf{V}^*\|_\infty &\leq \|\mathbf{\Lambda}(\hat{\mathbf{P}}_{\mathcal{K}} - \mathbf{P}_{\mathcal{K}}) \mathbf{V}^*\|_\infty + \|\mathbf{\Lambda}(\mathbf{P}_{\mathcal{K}} - \tilde{\mathbf{P}}_{\mathcal{K}}) \mathbf{V}^*\|_\infty + \|(\tilde{\mathbf{P}} - \mathbf{P}) \mathbf{V}^*\|_\infty \\
&\leq \|\mathbf{\Lambda}\|_1 \|(\hat{\mathbf{P}}_{\mathcal{K}} - \mathbf{P}_{\mathcal{K}}) \mathbf{V}^*\|_\infty + \|\mathbf{\Lambda}\|_1 \|(\mathbf{P}_{\mathcal{K}} - \tilde{\mathbf{P}}_{\mathcal{K}}) \mathbf{V}^*\|_\infty + \|\tilde{\mathbf{P}} - \mathbf{P}\|_1 \|\mathbf{V}^*\|_\infty \\
&\leq \|(\hat{\mathbf{P}}_{\mathcal{K}} - \mathbf{P}_{\mathcal{K}}) \mathbf{V}^*\|_\infty + \frac{2\xi}{1-\gamma},
\end{aligned} \tag{36}$$

where the first line uses  $\hat{\mathbf{P}} = \mathbf{\Lambda}\hat{\mathbf{P}}_{\mathcal{K}}$  and  $\tilde{\mathbf{P}} = \mathbf{\Lambda}\tilde{\mathbf{P}}_{\mathcal{K}}$ ; the last inequality comes from the facts that  $\|\tilde{\mathbf{P}} - \mathbf{P}\|_1 \leq \xi$ ,  $\|\mathbf{\Lambda}\|_1 = 1$  and  $\|\mathbf{V}^*\|_\infty \leq (1-\gamma)^{-1}$ . Then we turn to bound  $\|(\hat{\mathbf{P}}_{\mathcal{K}} - \mathbf{P}_{\mathcal{K}}) \mathbf{V}^*\|_\infty$ . In view of (4), Hoeffding's inequality (cf. [5, Theorem 2.2.6]) implies that for  $(s, a) \in \mathcal{K}$ ,

$$\mathbb{P}\left(\left|(\hat{\mathbf{P}} - \mathbf{P})_{s,a} \mathbf{V}^*\right| \geq t\right) \leq 2 \exp\left(-\frac{2t^2}{\|\mathbf{V}^*\|_\infty^2 / N}\right).$$

Hence by the standard union bound argument we have

$$\|(\hat{\mathbf{P}}_{\mathcal{K}} - \mathbf{P}_{\mathcal{K}}) \mathbf{V}^*\|_\infty \leq \sqrt{\frac{\|\mathbf{V}^*\|_\infty^2 \log(4K/\delta)}{2N}} \leq \sqrt{\frac{\log(4K/\delta)}{2N(1-\gamma)^2}}, \tag{37}$$

with probability over  $1 - \delta/2$ .

1. Now we are ready to bound  $\mathbf{Q}^{\pi^*} - \hat{\mathbf{Q}}^{\pi^*}$ . One has

$$\begin{aligned}
\mathbf{Q}^{\pi^*} - \hat{\mathbf{Q}}^{\pi^*} &= (\mathbf{I} - \gamma \mathbf{P}^{\pi^*})^{-1} \mathbf{r} - (\mathbf{I} - \gamma \hat{\mathbf{P}}^{\pi^*})^{-1} \mathbf{r} \\
&= (\mathbf{I} - \gamma \hat{\mathbf{P}}^{\pi^*})^{-1} \left( (\mathbf{I} - \gamma \hat{\mathbf{P}}^{\pi^*}) - (\mathbf{I} - \gamma \mathbf{P}^{\pi^*}) \right) \mathbf{Q}^{\pi^*} \\
&= \gamma (\mathbf{I} - \gamma \hat{\mathbf{P}}^{\pi^*})^{-1} (\mathbf{P}^{\pi^*} - \hat{\mathbf{P}}^{\pi^*}) \mathbf{Q}^{\pi^*} \\
&= \gamma (\mathbf{I} - \gamma \hat{\mathbf{P}}^{\pi^*})^{-1} (\mathbf{P} - \hat{\mathbf{P}}) \mathbf{V}^{\pi^*},
\end{aligned}$$

where the first equality makes use of (11). Then we take (36) and (37) collectively to achieve

$$\begin{aligned}
\left\| \gamma (\mathbf{I} - \gamma \hat{\mathbf{P}}^{\pi^*})^{-1} (\mathbf{P} - \hat{\mathbf{P}}) \mathbf{V}^* \right\|_\infty &\leq \gamma \sum_{i=0}^{\infty} \left\| \gamma^i (\hat{\mathbf{P}}^{\pi^*})^i (\mathbf{P} - \hat{\mathbf{P}}) \mathbf{V}^* \right\|_\infty \\
&\leq \gamma \sum_{i=0}^{\infty} \gamma^i \left\| (\hat{\mathbf{P}}^{\pi^*})^i \right\|_1 \left\| (\mathbf{P} - \hat{\mathbf{P}}) \mathbf{V}^* \right\|_\infty \\
&\leq \frac{\gamma}{1-\gamma} \sqrt{\frac{\log(4K/\delta)}{2N(1-\gamma)^2}} + \frac{2\gamma\xi}{(1-\gamma)^2},
\end{aligned}$$

where the last line comes from the fact that for all  $i \geq 1$ ,  $(\hat{\mathbf{P}}^{\pi^*})^i$  is a probability transition matrix so that  $\|(\hat{\mathbf{P}}^{\pi^*})^i\|_1 = 1$ . This justifies the first inequality (33).

2. In terms of the second one, [1, Section A.4] implies that

$$\left\| \mathbf{Q}^* - \hat{\mathbf{Q}}^* \right\|_\infty \leq \frac{\gamma}{1-\gamma} \left\| (\mathbf{P} - \hat{\mathbf{P}}) \mathbf{V}^* \right\|_\infty.$$

Substitution of (36) and (37) into the above inequality yields

$$\left\| \mathbf{Q}^* - \hat{\mathbf{Q}}^* \right\|_\infty \leq \frac{\gamma}{1-\gamma} \sqrt{\frac{\log(4K/\delta)}{2N(1-\gamma)^2}} + \frac{2\gamma\xi}{(1-\gamma)^2}.$$

## C Analysis of Q-learning (Proof of Theorem 2)

In this section, we will provide complete proof for Theorem 2. We actually prove a more general version of Theorem 2 that takes model misspecification into consideration, as stated below.

**Theorem 4.** Consider any  $\delta \in (0, 1)$  and  $\varepsilon \in (0, 1]$ . Suppose that there exists a probability transition model  $\tilde{\mathbf{P}}$  obeying Definition 1 and Assumption 1 with feature vectors  $\{\phi(s, a)\}_{(s,a) \in \mathcal{S} \times \mathcal{A}} \subset \mathbb{R}^K$  and anchor state-action pairs  $\mathcal{K}$  such that

$$\|\tilde{\mathbf{P}} - \mathbf{P}\|_1 \leq \xi$$

for some  $\xi \geq 0$ . Assume that the initialization obeys  $0 \leq Q_0(s, a) \leq \frac{1}{1-\gamma}$  for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$  and for any  $0 \leq t \leq T$ , the learning rates satisfy

$$\frac{1}{1 + \frac{c_1(1-\gamma)T}{\log^2 T}} \leq \eta_t \leq \frac{1}{1 + \frac{c_2(1-\gamma)t}{\log^2 T}}, \quad (38)$$

for some sufficiently small universal constants  $c_1 \geq c_2 > 0$ . Suppose that the total number of iterations  $T$  exceeds

$$T \geq \frac{C_3 \log(KT/\delta) \log^4 T}{(1-\gamma)^4 \varepsilon^2}, \quad (39)$$

for some sufficiently large universal constant  $C_3 > 0$ . If there exists a linear probability transition model  $\tilde{\mathbf{P}}$  satisfying Assumption 1 with feature vectors  $\{\phi(s, a)\}_{(s,a) \in \mathcal{S} \times \mathcal{A}}$  such that  $\|\tilde{\mathbf{P}} - \mathbf{P}\|_1 \leq \xi$ , then with probability exceeding  $1 - \delta$ , the output  $Q_T$  of Algorithm 2 satisfies

$$\max_{(s,a) \in \mathcal{S} \times \mathcal{A}} |Q_T(s, a) - Q^*(s, a)| \leq \varepsilon + \frac{6\gamma\xi}{(1-\gamma)^2}, \quad (40)$$

for some constant  $C_4 > 0$ . In addition, let  $\pi_T$  (resp.  $V_T$ ) to be the policy (resp. value function) induced by  $Q_T$ , then one has

$$\max_{s \in \mathcal{S}} |V^{\pi_T}(s) - V^*(s)| \leq \frac{2\gamma}{1-\gamma} \left( \varepsilon + \frac{6\gamma\xi}{(1-\gamma)^2} \right). \quad (41)$$

Theorem 4 subsumes Theorem 2 as a special case with  $\xi = 0$ . The remainder of this section is devoted to proving Theorem 4.

### C.1 Proof of Theorem 4

First we show that (41) can be easily obtained from (40). Since [49] gives rise to

$$\|V^{\pi_T} - V^*\|_\infty \leq \frac{2\gamma\|V_T - V^*\|_\infty}{1-\gamma},$$

we have

$$\|V^{\pi_T} - V^*\|_\infty \leq \frac{2\gamma\|Q_T - Q^*\|_\infty}{1-\gamma},$$

due to  $\|V_T - V^*\|_\infty \leq \|Q_T - Q^*\|_\infty$ . Then (41) follows directly from (40).

Therefore, we are left to justify (40). To start with, we consider the update rule

$$\mathbf{Q}_t = (1 - \eta_t) \mathbf{Q}_{t-1} + \eta_t \left( \mathbf{r} + \gamma \hat{\mathbf{P}}_t \mathbf{V}_{t-1} \right).$$

By defining the error term  $\Delta_t := \mathbf{Q}_t - \mathbf{Q}^*$ , we can decompose  $\Delta_t$  into

$$\begin{aligned} \Delta_t &= (1 - \eta_t) \mathbf{Q}_{t-1} + \eta_t \left( \mathbf{r} + \gamma \hat{\mathbf{P}}_t \mathbf{V}_{t-1} \right) - \mathbf{Q}^* \\ &= (1 - \eta_t) (\mathbf{Q}_{t-1} - \mathbf{Q}^*) + \eta_t \left( \mathbf{r} + \gamma \hat{\mathbf{P}}_t \mathbf{V}_{t-1} - \mathbf{Q}^* \right) \\ &= (1 - \eta_t) (\mathbf{Q}_{t-1} - \mathbf{Q}^*) + \gamma \eta_t \left( \hat{\mathbf{P}}_t \mathbf{V}_{t-1} - \mathbf{P} \mathbf{V}^* \right) \\ &= (1 - \eta_t) \Delta_{t-1} + \gamma \eta_t \Lambda \left( \hat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{t-1} + \gamma \eta_t \Lambda \mathbf{P}_{\mathcal{K}} (\mathbf{V}_{t-1} - \mathbf{V}^*) \\ &\quad + \gamma \eta_t (\Lambda \mathbf{P}_{\mathcal{K}} - \mathbf{P}) \mathbf{V}^*. \end{aligned} \quad (42)$$

Here in the penultimate equality, we make use of  $\mathbf{Q}^* = \mathbf{r} + \gamma \mathbf{P} \mathbf{V}^*$ ; and the last equality comes from  $\widehat{\mathbf{P}}_t = \Lambda \widehat{\mathbf{P}}_{\mathcal{K}}^{(t)}$  which is defined in (15). It is straightforward to check that  $\Lambda \mathbf{P}_{\mathcal{K}}$  is also a probability transition matrix. We denote by  $\overline{\mathbf{P}} = \Lambda \mathbf{P}_{\mathcal{K}}$  hereafter. The third term in the decomposition above can be upper and lower bounded by

$$\overline{\mathbf{P}}(\mathbf{V}_{t-1} - \mathbf{V}^*) = \overline{\mathbf{P}}^{\pi_{t-1}} \mathbf{Q}_{t-1} - \overline{\mathbf{P}}^{\pi^*} \mathbf{Q}^* \leq \overline{\mathbf{P}}^{\pi_{t-1}} \mathbf{Q}_{t-1} - \overline{\mathbf{P}}^{\pi_{t-1}} \mathbf{Q}^* = \overline{\mathbf{P}}^{\pi_{t-1}} \Delta_{t-1},$$

and

$$\overline{\mathbf{P}}(\mathbf{V}_{t-1} - \mathbf{V}^*) = \overline{\mathbf{P}}^{\pi_{t-1}} \mathbf{Q}_{t-1} - \overline{\mathbf{P}}^{\pi^*} \mathbf{Q}^* \geq \overline{\mathbf{P}}^{\pi^*} \mathbf{Q}_{t-1} - \overline{\mathbf{P}}^{\pi^*} \mathbf{Q}^* = \overline{\mathbf{P}}^{\pi^*} \Delta_{t-1}.$$

Plugging these bounds into (42) yields

$$\Delta_t \leq (1 - \eta_t) \Delta_{t-1} + \gamma \eta_t \Lambda \left( \widehat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{t-1} + \gamma \eta_t \overline{\mathbf{P}}^{\pi_{t-1}} \Delta_{t-1} + \gamma \eta_t (\Lambda \mathbf{P}_{\mathcal{K}} - \mathbf{P}) \mathbf{V}^*,$$

$$\Delta_t \geq (1 - \eta_t) \Delta_{t-1} + \gamma \eta_t \Lambda \left( \widehat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{t-1} + \gamma \eta_t \overline{\mathbf{P}}^{\pi^*} \Delta_{t-1} + \gamma \eta_t (\Lambda \mathbf{P}_{\mathcal{K}} - \mathbf{P}) \mathbf{V}^*.$$

Repeatedly invoking these two recursive relations leads to

$$\Delta_t \leq \eta_0^{(t)} \Delta_0 + \sum_{i=1}^t \eta_i^{(t)} \gamma \left( \overline{\mathbf{P}}^{\pi_{t-1}} \Delta_{t-1} + \Lambda \left( \widehat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{t-1} + (\Lambda \mathbf{P}_{\mathcal{K}} - \mathbf{P}) \mathbf{V}^* \right), \quad (43)$$

$$\Delta_t \geq \eta_0^{(t)} \Delta_0 + \sum_{i=1}^t \eta_i^{(t)} \gamma \left( \overline{\mathbf{P}}^{\pi^*} \Delta_{t-1} + \Lambda \left( \widehat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{t-1} + (\Lambda \mathbf{P}_{\mathcal{K}} - \mathbf{P}) \mathbf{V}^* \right), \quad (44)$$

where

$$\eta_i^{(t)} := \begin{cases} \prod_{j=1}^t (1 - \eta_j), & \text{if } i = 0, \\ \eta_i \prod_{j=i+1}^t (1 - \eta_j), & \text{if } 0 < i < t, \\ \eta_t, & \text{if } i = t. \end{cases}$$

Here we adopt the same notations as [4].

To begin with, we consider the upper bound (43). It can be further decomposed as

$$\begin{aligned} \Delta_t &\leq \underbrace{\eta_0^{(t)} \Delta_0 + \sum_{i=1}^{(1-\alpha)t} \eta_i^{(t)} \gamma \left( \overline{\mathbf{P}}^{\pi_{t-1}} \Delta_{t-1} + \Lambda \left( \widehat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{t-1} \right)}_{=:\boldsymbol{\theta}_t} \\ &\quad + \underbrace{\sum_{i=(1-\alpha)t+1}^t \eta_i^{(t)} \gamma \Lambda \left( \widehat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{i-1}}_{=:\boldsymbol{\nu}_t} \\ &\quad + \underbrace{\sum_{i=1}^t \eta_i^{(t)} \gamma (\Lambda \mathbf{P}_{\mathcal{K}} - \mathbf{P}) \mathbf{V}^*}_{=:\boldsymbol{\omega}_t} + \sum_{i=(1-\alpha)t+1}^t \eta_i^{(t)} \gamma \overline{\mathbf{P}}^{\pi_{t-1}} \Delta_{i-1}, \end{aligned} \quad (45)$$

where we define  $\alpha := C_4(1 - \gamma)/\log T$  for some constant  $C_4 > 0$ . Next, we turn to bound  $\boldsymbol{\theta}_t$  and  $\boldsymbol{\nu}_t$  respectively for any  $t$  satisfying  $\frac{T}{c_2 \log \frac{1}{1-\gamma}} \leq t \leq T$  with stepsize choice (8).

**Bounding  $\boldsymbol{\omega}_t$ .** It is straightforward to bound

$$\begin{aligned} \|\boldsymbol{\omega}_t\|_{\infty} &\stackrel{(i)}{=} \|\gamma (\Lambda \mathbf{P}_{\mathcal{K}} - \mathbf{P}) \mathbf{V}^*\|_{\infty} \\ &\stackrel{(ii)}{\leq} \gamma \left( \|\Lambda\|_1 \left\| \left( \mathbf{P}_{\mathcal{K}} - \tilde{\mathbf{P}}_{\mathcal{K}} \right) \mathbf{V}^* \right\|_{\infty} + \left\| \left( \tilde{\mathbf{P}} - \mathbf{P} \right) \mathbf{V}^* \right\|_{\infty} \right) \\ &\stackrel{(iii)}{\leq} \frac{2\gamma\xi}{1-\gamma}, \end{aligned}$$

where the first equality comes from the fact that  $\sum_{i=1}^t \eta_i^{(t)} = 1$  [4, Equation (40)]; the second inequality utilizes  $\tilde{\mathbf{P}} = \Lambda \tilde{\mathbf{P}}_{\mathcal{K}}$ ; the last line uses the facts that  $\|\Lambda\|_1 = 1$ ,  $\|\mathbf{V}^*\|_{\infty} \leq (1 - \gamma)^{-1}$  and  $\|\tilde{\mathbf{P}}_{\mathcal{K}} - \mathbf{P}_{\mathcal{K}}\|_1 \leq \|\tilde{\mathbf{P}} - \mathbf{P}\|_1 \leq \xi$ .

**Bounding  $\theta_t$ .** By similar derivation as Step 1 in [4, Appendix A.2], we have

$$\begin{aligned}
\|\theta_t\|_\infty &\leq \eta_0^{(t)} \|\Delta_0\|_\infty + t \max_{1 \leq i \leq (1-\alpha)t} \eta_i^{(t)} \max_{1 \leq i \leq (1-\alpha)t} \left( \left\| \overline{\mathbf{P}}^{\pi_{t-1}} \Delta_{i-1} \right\|_\infty + \left\| \Lambda \widehat{\mathbf{P}}_{\mathcal{K}}^{(t)} \mathbf{V}_{i-1} \right\|_\infty + \left\| \Lambda \mathbf{P}_{\mathcal{K}} \mathbf{V}_{i-1} \right\|_\infty \right) \\
&\stackrel{(i)}{\leq} \eta_0^{(t)} \|\Delta_0\|_\infty + t \max_{1 \leq i \leq (1-\alpha)t} \eta_i^{(t)} \max_{1 \leq i \leq (1-\alpha)t} (\|\Delta_{i-1}\|_\infty + 2 \|\mathbf{V}_{i-1}\|_\infty) \\
&\stackrel{(ii)}{\leq} \frac{1}{2T^2} \cdot \frac{1}{1-\gamma} + \frac{1}{2T^2} \cdot t \cdot \frac{3}{1-\gamma} \\
&\leq \frac{2}{(1-\gamma)T},
\end{aligned} \tag{46}$$

where (i) is due to the fact that  $\|\overline{\mathbf{P}}^{\pi_{t-1}}\|_1 = \|\Lambda \widehat{\mathbf{P}}_{\mathcal{K}}^{(t)}\|_1 = \|\Lambda \mathbf{P}_{\mathcal{K}}\|_1 = 1$  and (ii) comes from [4, Equation (39a)].

**Bounding  $\nu_t$ .** To control the second term, we apply the following Freedman's inequality.

**Lemma 3** (Freedman's Inequality). *Consider a real-valued martingale  $\{Y_k : k = 0, 1, 2, \dots\}$  with difference sequence  $\{X_k : k = 1, 2, 3, \dots\}$ . Assume that the difference sequence is uniformly bounded:*

$$|X_k| \leq R \quad \text{and} \quad \mathbb{E}[X_k | \{X_j\}_{j=1}^{k-1}] = 0 \quad \text{for all } k \geq 1.$$

Let

$$S_n := \sum_{k=1}^n X_k, \quad T_n := \sum_{k=1}^n \text{Var} \{X_k | \{X_j\}_{j=1}^{k-1}\}.$$

Then for any given  $\sigma^2 \geq 0$ , one has

$$\mathbb{P}(|S_n| \geq \tau \text{ and } T_n \leq \sigma^2) \leq 2 \exp\left(-\frac{\tau^2/2}{\sigma^2 + R\tau/3}\right).$$

In addition, suppose that  $W_n \leq \sigma^2$  holds deterministically. For any positive integer  $K \geq 1$ , with probability at least  $1 - \delta$  one has

$$|S_n| \leq \sqrt{8 \max\left\{T_n, \frac{\sigma^2}{2K}\right\} \log \frac{2K}{\delta}} + \frac{4}{3} R \log \frac{2K}{\delta}.$$

*Proof.* See [4, Theorem 4]. □

To apply this inequality, we can express  $\nu_t$  as

$$\nu_t := \sum_{i=(1-\alpha)t+1}^t \mathbf{x}_i,$$

with

$$\mathbf{x}_i := \eta_i^{(t)} \gamma \Lambda \left( \widehat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{i-1}, \quad \text{and} \quad \mathbb{E}[\mathbf{x}_i | \mathbf{V}_{i-1}, \dots, \mathbf{V}_0] = \mathbf{0}. \tag{47}$$

1. In order to calculate bound  $R$  in Lemma 3, one has

$$\begin{aligned}
B &:= \max_{(1-\alpha)t < t \leq t} \|\mathbf{x}_i\|_\infty \leq \max_{(1-\alpha)t < t \leq t} \eta_i^{(t)} \Lambda \left( \widehat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{i-1} \Big\|_\infty \\
&\leq \max_{(1-\alpha)t < t \leq t} \eta_i^{(t)} \left( \left\| \Lambda \widehat{\mathbf{P}}_{\mathcal{K}}^{(t)} \right\|_1 + \|\Lambda \mathbf{P}_{\mathcal{K}}\|_1 \right) \|\mathbf{V}_{i-1}\|_\infty \\
&\leq \max_{(1-\alpha)t < t \leq t} \eta_i^{(t)} \cdot \frac{2}{1-\gamma} \leq \frac{4 \log^4 T}{(1-\gamma)^2 T},
\end{aligned}$$

where the last inequality comes from [4, Eqn (39b)] and the fact that  $\|\mathbf{V}_{i-1}\|_\infty \leq \frac{1}{1-\gamma}$ .

2. Then regarding the variance term, we claim for the moment that

$$\begin{aligned}\mathbf{W}_t &:= \sum_{i=(1-\alpha)t+1}^t \text{diag}(\text{Var}(\mathbf{x}_i | \mathbf{V}_{i-1}, \dots, \mathbf{V}_0)) \\ &\leq \gamma^2 \sum_{i=(1-\alpha)t+1}^t \left(\eta_i^{(t)}\right)^2 \text{Var}_{\overline{\mathcal{P}}}(\mathbf{V}_{i-1}).\end{aligned}\quad (48)$$

Then we have

$$\begin{aligned}\mathbf{W}_t &\leq \max_{(1-\alpha)t \leq i \leq t} \eta_i^{(t)} \left( \sum_{i=(1-\alpha)t+1}^t \eta_i^{(t)} \right) \max_{(1-\alpha)t \leq i < t} \text{Var}_{\overline{\mathcal{P}}}(\mathbf{V}_i) \\ &\leq \frac{2 \log^4 T}{(1-\gamma)T} \max_{(1-\alpha)t \leq i < t} \text{Var}_{\overline{\mathcal{P}}}(\mathbf{V}_i),\end{aligned}\quad (49)$$

where the second line comes from [4, Eqns (39b), (40)]. A trivial upper bound for  $\mathbf{W}_t$  is

$$|\mathbf{W}_t| \leq \frac{2 \log^4 T}{(1-\gamma)T} \cdot \frac{1}{(1-\gamma)^2} \mathbf{1} = \frac{2 \log^4 T}{(1-\gamma)^3 T} \mathbf{1},$$

which uses the fact that  $\text{Var}_{\mathcal{P}}(\mathbf{V}_i) \leq \|\mathbf{V}_i\|_\infty^2 \leq 1/(1-\gamma)^2$ .

Then, we invoke Lemma 3 with  $K = \lceil 2 \log_2 \frac{1}{1-\gamma} \rceil$  and apply the union bound argument over  $\mathcal{K}$  to arrive at

$$\begin{aligned}|\boldsymbol{\nu}_t| &\leq \sqrt{8 \left( \mathbf{W}_t + \frac{\sigma^2}{2K} \mathbf{1} \right) \log \frac{8KT \log \frac{1}{1-\gamma}}{\delta} + \frac{4}{3} B \log \frac{8KT \log \frac{1}{1-\gamma}}{\delta} \mathbf{1}} \\ &\leq \sqrt{8 \left( \mathbf{W}_t + \frac{2 \log^4 T}{(1-\gamma)T} \mathbf{1} \right) \log \frac{8KT}{\delta} + \frac{4}{3} B \log \frac{8KT \log \frac{1}{1-\gamma}}{\delta} \mathbf{1}} \\ &\leq \sqrt{\frac{32 \log^4 T}{(1-\gamma)T} \log \frac{8KT}{\delta} \left( \max_{(1-\alpha)t \leq i < t} \text{Var}_{\Lambda \mathcal{P}_{\mathcal{K}}}(\mathbf{V}_i) + \mathbf{1} \right) + \frac{12 \log^4 T}{(1-\gamma)^2 T} \log \frac{8KT}{\delta} \mathbf{1}}.\end{aligned}\quad (50)$$

Hence if we define

$$\boldsymbol{\varphi}_t := 64 \frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)T} \left( \max_{\frac{t}{2} \leq i \leq t} \text{Var}_{\overline{\mathcal{P}}}(\mathbf{V}_i) + \mathbf{1} \right),$$

then (46) and (50) implies that

$$|\boldsymbol{\theta}_t| + |\boldsymbol{\nu}_t| + |\boldsymbol{\omega}_t| \leq \sqrt{\boldsymbol{\varphi}_t} + \frac{2\gamma\xi}{1-\gamma} \mathbf{1},\quad (51)$$

with probability over  $1-\delta$  for all  $2t/3 \leq k \leq t$ , as long as  $T \gg \log^4 T \log \frac{KT}{\delta} / (1-\gamma)^3$ . Therefore, plugging (51) into (45), we arrive at the recursive relationship

$$\boldsymbol{\Delta}_t \leq \sqrt{\boldsymbol{\varphi}_t} + \frac{2\gamma\xi}{1-\gamma} \mathbf{1} + \sum_{i=(1-\alpha)k+1}^k \eta_i^{(k)} \gamma \overline{\mathcal{P}}^{\pi_{i-1}} \boldsymbol{\Delta}_{i-1} = \sqrt{\boldsymbol{\varphi}_t} + \frac{2\gamma\xi}{1-\gamma} \mathbf{1} + \sum_{i=(1-\alpha)k}^{k-1} \eta_i^{(k)} \gamma \overline{\mathcal{P}}^{\pi_{i-1}} \boldsymbol{\Delta}_i.$$

This recursion is expressed in a similar way as [4, Eqn. (46)] so we can invoke similar derivation in [4, Appendix A.2] to obtain that

$$\boldsymbol{\Delta}_t \leq 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \left( 1 + \max_{\frac{t}{2} \leq i < t} \|\boldsymbol{\Delta}_i\|_\infty \right)} \mathbf{1} + \frac{2\gamma\xi}{(1-\gamma)^2} \mathbf{1}.\quad (52)$$

Then we turn to (44). Applying a similar argument, we can deduce that

$$\boldsymbol{\Delta}_t \geq -30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \left( 1 + \max_{\frac{t}{2} \leq i < t} \|\boldsymbol{\Delta}_i\|_\infty \right)} \mathbf{1} - \frac{2\gamma\xi}{(1-\gamma)^2} \mathbf{1}.\quad (53)$$

For any  $t$  satisfying  $\frac{T}{c_2 \log \frac{1}{1-\gamma}} \leq t \leq T$ , taking (52) and (53) collectively gives rise to

$$\|\Delta_t\|_\infty \leq 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \left(1 + \max_{\frac{t}{2} \leq i < t} \|\Delta_i\|_\infty\right)} + \frac{2\gamma\xi}{(1-\gamma)^2}. \quad (54)$$

Let

$$u_k := \max \left\{ \|\Delta_t\|_\infty : 2^k \frac{T}{c_2 \log \frac{1}{1-\gamma}} \leq t \leq T \right\}.$$

By taking supremum over  $t \in \{[2^k T / (c_2 \log \frac{1}{1-\gamma})], \dots, T\}$  on both sides of (54), we have

$$u_k \leq 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} (1 + u_{k-1})} + \frac{2\gamma\xi}{(1-\gamma)^2} \quad \forall 1 \leq k \leq \log \left( c_2 \log \frac{1}{1-\gamma} \right). \quad (55)$$

It is straightforward to bound  $u_0 \leq \frac{1}{1-\gamma}$ . For  $k \geq 1$ , it is straightforward to obtain from (55) that

$$u_k \leq 3 \max \left\{ 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}}, 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} u_{k-1}}, \frac{2\gamma\xi}{(1-\gamma)^2} \right\}, \quad (56)$$

for  $1 \leq k \leq \log(c_2 \log \frac{1}{1-\gamma})$ . We analyze (56) under two different cases:

1. If there exists some integer  $k_0$  with  $1 \leq k_0 < \lceil \log(c_2 \log \frac{1}{1-\gamma}) \rceil$ , such that

$$u_{k_0} \leq \max \left\{ 1, \frac{6\gamma\xi}{(1-\gamma)^2} \right\},$$

then it is straightforward to check from (56) that

$$u_{k_0+1} \leq 3 \max \left\{ 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}}, \frac{2\gamma\xi}{(1-\gamma)^2} \right\} \quad (57)$$

as long as  $T \geq C_3 (1-\gamma)^{-4} \log^4 T \log(KT/\delta)$  for some sufficiently large constant  $C_3 > 0$ .

2. Otherwise we have  $u_k > \max\{1, \frac{6\gamma\xi}{(1-\gamma)^2}\}$  for all  $1 \leq k < \lceil \log(c_2 \log \frac{1}{1-\gamma}) \rceil$ . This together with (56) suggests that

$$\max \left\{ 1, \frac{6\gamma\xi}{(1-\gamma)^2} \right\} < 3 \max \left\{ 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}}, 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} u_{k-1}}, \frac{2\gamma\xi}{(1-\gamma)^2} \right\},$$

and therefore

$$\max \left\{ 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}}, 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} u_{k-1}}, \frac{2\gamma\xi}{(1-\gamma)^2} \right\} = 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} u_{k-1}}$$

for all  $1 \leq k \leq \log(c_2 \log \frac{1}{1-\gamma})$ . Let

$$v_k := 90 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} u_{k-1}}.$$

Then we know from (55) that

$$u_k \leq v_k \quad \forall 1 \leq k \leq \log \left( c_2 \log \frac{1}{1-\gamma} \right).$$

By applying the above two inequalities recursively, we know that

$$\begin{aligned}
u_k \leq v_k &= \left( \frac{8100 \log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \right)^{1/2} u_{k-1}^{1/2} \leq \left( \frac{8100 \log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \right)^{1/2} v_{k-1}^{1/2} \\
&\leq \left( \frac{8100 \log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \right)^{1/2+1/4} u_{k-2}^{1/4} \leq \left( \frac{8100 \log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \right)^{1/2+1/4} v_{k-2}^{1/4} \\
&\leq \dots \leq \left( \frac{8100 \log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \right)^{1-1/2^k} u_0^{1/2^k} \leq \sqrt{\frac{8100 \log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}} \left( \frac{1}{1-\gamma} \right)^{1/2^k},
\end{aligned}$$

where the last inequality holds as long as  $T \geq C_3 \log^4 T \log(KT/\delta)(1-\gamma)^{-4}$  for some sufficiently large constant  $C_3 > 0$ . Let  $k_0 = \tilde{c} \log \log \frac{1}{1-\gamma}$  for some properly chosen constant  $\tilde{c} > 0$  such that  $k_0$  is an integer between 1 and  $\log(c_2 \log \frac{1}{1-\gamma})$ , we have

$$u_{k_0} \leq \sqrt{\frac{8100 \log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}} \left( \frac{1}{1-\gamma} \right)^{1/2^{k_0}} = O\left( \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}} \right).$$

When  $T \geq C_3 \log^4 T \log(KT/\delta)(1-\gamma)^{-4}$  for some sufficiently large constant  $C_3 > 0$ , this implies that  $u_{k_0} < 1$ , which contradicts with the preassumption that  $u_k > \max\{1, \frac{6\gamma\xi}{(1-\gamma)^2}\}$  for all  $1 \leq k \leq c_2 \log \frac{1}{1-\gamma}$ .

Consequently, (57) must hold true and then the definition of  $u_k$  immediately leads to

$$\|\Delta_T\|_\infty \leq 90 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}} + \frac{6\gamma\xi}{(1-\gamma)^2}.$$

Then for any  $\varepsilon \in (0, 1]$ , one has

$$\|\Delta_T\|_\infty \leq \varepsilon + \frac{6\gamma\xi}{(1-\gamma)^2},$$

as long as

$$90 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}} \leq \varepsilon.$$

Hence, if the total number of iterations  $T$  satisfies

$$T \geq C_3 \frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 \varepsilon^2}$$

for some sufficiently large constant  $C_3 > 0$ , (10) would hold for Algorithm 1 with probability over  $1 - \delta$ .

Finally, we are left to justify (48). Recall the definition of  $\mathbf{x}_i$  (cf. (47)), one has

$$\begin{aligned}
\text{diag}(\text{Var}(\mathbf{x}_i | \mathbf{V}_{i-1}, \dots, \mathbf{V}_0)) &= \gamma^2 \left( \eta_i^{(t)} \right)^2 \text{diag} \left( \text{Var} \left( \mathbf{\Lambda} \left( \hat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{i-1} | \mathbf{V}_{i-1} \right) \right) \\
&= \gamma^2 \left( \eta_i^{(t)} \right)^2 \text{diag} \left( \mathbf{\Lambda} \text{Var} \left( \left( \hat{\mathbf{P}}_{\mathcal{K}}^{(i)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{i-1} | \mathbf{V}_{i-1} \right) \mathbf{\Lambda}^\top \right) \\
&= \gamma^2 \left( \eta_i^{(t)} \right)^2 \left\{ \boldsymbol{\lambda}(s, a)^2 \text{Var}_{\mathbf{P}_{\mathcal{K}}}(\mathbf{V}_{i-1}) \right\}_{s,a},
\end{aligned}$$

where the notation  $\text{Var}_{\mathbf{P}_{\mathcal{K}}}(\mathbf{V}_{i-1})$  is defined in (12). Plugging this into the definition of  $\mathbf{W}_t$  leads to

$$\begin{aligned}
\mathbf{W}_t &= \gamma^2 \sum_{i=(1-\alpha)t+1}^t \left( \eta_i^{(t)} \right)^2 \left\{ \boldsymbol{\lambda}(s, a)^2 \text{Var}_{\mathbf{P}_{\mathcal{K}}}(\mathbf{V}_{i-1}) \right\}_{s,a} \\
&= \gamma^2 \sum_{i=(1-\alpha)t+1}^t \left( \eta_i^{(t)} \right)^2 \left\{ \boldsymbol{\lambda}(s, a)^2 \left( \mathbf{P}_{\mathcal{K}}(\mathbf{V}_{i-1} \circ \mathbf{V}_{i-1}) - (\mathbf{P}_{\mathcal{K}} \mathbf{V}_{i-1}) \circ (\mathbf{P}_{\mathcal{K}} \mathbf{V}_{i-1}) \right) \right\}_{s,a}. \quad (58)
\end{aligned}$$

Then we introduce a useful claim as follows. The proof is deferred to Appendix C.2.

*Claim 1.* For any state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$  and vector  $\mathbf{V} \in \mathbb{R}^{|\mathcal{S}|}$ , one has

$$\begin{aligned} & \lambda(s, a)^2 (\mathbf{P}_{\mathcal{K}}(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}_{\mathcal{K}}\mathbf{V}) \circ (\mathbf{P}_{\mathcal{K}}\mathbf{V})) \\ & \leq \lambda(s, a) \mathbf{P}_{\mathcal{K}}(\mathbf{V} \circ \mathbf{V}) - (\lambda(s, a) \mathbf{P}_{\mathcal{K}}\mathbf{V}) \circ (\lambda(s, a) \mathbf{P}_{\mathcal{K}}\mathbf{V}). \end{aligned} \quad (59)$$

By invoking this claim with  $\mathbf{V} = \mathbf{V}^{i-1}$  and taking collectively with (58), one has

$$\begin{aligned} \mathbf{W}_t & \leq \gamma^2 \sum_{i=(1-\beta)t+1}^t \left( \eta_i^{(t)} \right)^2 \{ \lambda(s, a) \mathbf{P}_{\mathcal{K}}(\mathbf{V}_{i-1} \circ \mathbf{V}_{i-1}) - (\lambda(s, a) \mathbf{P}_{\mathcal{K}}\mathbf{V}_{i-1}) \circ (\lambda(s, a) \mathbf{P}_{\mathcal{K}}\mathbf{V}_{i-1}) \}_{s,a} \\ & = \gamma^2 \sum_{i=(1-\beta)t+1}^t \left( \eta_i^{(t)} \right)^2 [ \Lambda \mathbf{P}_{\mathcal{K}}(\mathbf{V}_{i-1} \circ \mathbf{V}_{i-1}) - (\Lambda \mathbf{P}_{\mathcal{K}}\mathbf{V}_{i-1}) \circ (\Lambda \mathbf{P}_{\mathcal{K}}\mathbf{V}_{i-1}) ] \\ & = \gamma^2 \sum_{i=(1-\beta)t+1}^t \left( \eta_i^{(t)} \right)^2 \text{Var}_{\overline{\mathbf{P}}}(\mathbf{V}_{i-1}), \end{aligned}$$

which is the desired result.

## C.2 Proof of Claim 1

To simplify notations in this proof, we use  $[\lambda_i]_{i=1}^K$ ,  $[P_{i,j}]_{1 \leq i \leq K, 1 \leq j \leq |\mathcal{S}|}$  and  $[V_i]_{i=1}^{|\mathcal{S}|}$  to denote  $\lambda(s, a)$ ,  $\mathbf{P}_{\mathcal{K}}$  and  $\mathbf{V}$  respectively. Then one has

$$\begin{aligned} & \lambda(s, a) \mathbf{P}_{\mathcal{K}}(\mathbf{V} \circ \mathbf{V}) - (\lambda(s, a) \mathbf{P}_{\mathcal{K}}\mathbf{V}) \circ (\lambda(s, a) \mathbf{P}_{\mathcal{K}}\mathbf{V}) \\ & \quad - \lambda(s, a)^2 (\mathbf{P}_{\mathcal{K}}(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}_{\mathcal{K}}\mathbf{V}) \circ (\mathbf{P}_{\mathcal{K}}\mathbf{V})) \\ & = \sum_{i=1}^K \sum_{j=1}^{|\mathcal{S}|} \lambda_i P_{i,j} V_j^2 - \left( \sum_{i=1}^K \sum_{j=1}^{|\mathcal{S}|} \lambda_i P_{i,j} V_j \right)^2 - \sum_{i=1}^K \sum_{j=1}^{|\mathcal{S}|} \lambda_i^2 P_{i,j} V_j^2 + \sum_{i=1}^K \lambda_i^2 \left( \sum_{j=1}^{|\mathcal{S}|} P_{i,j} V_j \right)^2 \\ & = \sum_{i=1}^K \sum_{j=1}^{|\mathcal{S}|} \lambda_i P_{i,j} V_j \left[ (1 - \lambda_i) V_j - \sum_{i' \neq i} \sum_{j'=1}^{|\mathcal{S}|} \lambda_{i'} P_{i',j'} V_{j'} \right]. \\ & = \sum_{i=1}^K \sum_{j=1}^{|\mathcal{S}|} \lambda_i P_{i,j} V_j \left[ \left( \sum_{i'=1}^K \sum_{j'=1}^{|\mathcal{S}|} \lambda_{i'} P_{i',j'} - \lambda_i \right) V_j - \sum_{i' \neq i} \sum_{j'=1}^{|\mathcal{S}|} \lambda_{i'} P_{i',j'} V_{j'} \right] \\ & = \sum_{i=1}^K \sum_{j=1}^{|\mathcal{S}|} \sum_{i' \neq i} \sum_{j'=1}^{|\mathcal{S}|} \lambda_i P_{i,j} V_j \lambda_{i'} P_{i',j'} (V_j - V_{j'}) \end{aligned}$$

where in the penultimate equality, we use the fact that

$$\sum_{i'=1}^K \sum_{j'=1}^{|\mathcal{S}|} \lambda_{i'} P_{i',j'} = \lambda(s, a) \mathbf{P}_{\mathcal{K}} \mathbf{1} = 1.$$

It follows that

$$\begin{aligned}
& \lambda(s, a) \mathbf{P}_{\mathcal{K}}(\mathbf{V} \circ \mathbf{V}) - (\lambda(s, a) \mathbf{P}_{\mathcal{K}} \mathbf{V}) \circ (\lambda(s, a) \mathbf{P}_{\mathcal{K}} \mathbf{V}) \\
& \quad - \lambda(s, a)^2 (\mathbf{P}_{\mathcal{K}}(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}_{\mathcal{K}} \mathbf{V}) \circ (\mathbf{P}_{\mathcal{K}} \mathbf{V})) \\
& = \sum_{i=1}^K \sum_{1 \leq i' < i} \sum_{j=1}^{|\mathcal{S}|} \sum_{j'=1}^{|\mathcal{S}|} [\lambda_i P_{i,j} V_j \lambda_{i'} P_{i',j'} (V_j - V_{j'}) + \lambda_{i'} P_{i',j} V_j \lambda_i P_{i,j'} (V_j - V_{j'})] \\
& = \sum_{i=1}^K \sum_{1 \leq i' < i} \lambda_i \lambda_{i'} \left[ \sum_{j=1}^{|\mathcal{S}|} \sum_{j'=1}^{|\mathcal{S}|} P_{i,j} V_j P_{i',j'} (V_j - V_{j'}) + \sum_{j=1}^{|\mathcal{S}|} \sum_{j'=1}^{|\mathcal{S}|} P_{i',j} V_j P_{i,j'} (V_j - V_{j'}) \right] \\
& \stackrel{(i)}{=} \sum_{i=1}^K \sum_{1 \leq i' < i} \lambda_i \lambda_{i'} \left[ \sum_{j=1}^{|\mathcal{S}|} \sum_{j'=1}^{|\mathcal{S}|} P_{i,j} V_j P_{i',j'} (V_j - V_{j'}) + \sum_{j=1}^{|\mathcal{S}|} \sum_{j'=1}^{|\mathcal{S}|} P_{i',j'} V_{j'} P_{i,j} (V_{j'} - V_j) \right] \\
& = \sum_{i=1}^K \sum_{1 \leq i' < i} \lambda_i \lambda_{i'} \left[ \sum_{j=1}^{|\mathcal{S}|} \sum_{j'=1}^{|\mathcal{S}|} P_{i,j} P_{i',j'} (V_j - V_{j'})^2 \right] \\
& \geq 0,
\end{aligned}$$

where in (i), we exchange the indices  $j$  and  $j'$ .

## D Feature dimension and the number of anchor state-action pairs

The assumption that the feature dimension (denoted by  $K_d$ ) and the number of anchor state-action pairs (denoted by  $K_n$ ) are equal is actually non-essential. In what follows, we will show that if  $K_d \neq K_n$ , then we can modify the current feature mapping  $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^{K_d}$  to achieve a new feature mapping  $\phi' : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^{K_n}$  that does not change the transition model  $P$ . By doing so, the new feature dimension  $K_n$  equals to the number of anchor state-action pairs.

To begin with, we recall from Definition 1 that there exists  $K_d$  unknown functions  $\psi_1, \dots, \psi_{K_d} : \mathcal{S} \rightarrow \mathbb{R}$ , such that

$$P(s'|s, a) = \sum_{k=1}^{K_d} \phi_k(s, a) \psi_k(s'),$$

for every  $(s, a) \in \mathcal{S} \times \mathcal{A}$  and  $s' \in \mathcal{S}$ . In addition, we also recall from Assumption 1 that there exists  $\mathcal{K} \subseteq \mathcal{S} \times \mathcal{A}$  with  $|\mathcal{K}| = K_n$  such that for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\phi(s, a) = \sum_{i:(s_i, a_i) \in \mathcal{K}} \lambda_i(s, a) \phi(s_i, a_i) \in \mathbb{R}^{K_d} \quad \text{for} \quad \sum_{i=1}^{K_n} \lambda_i(s, a) = 1 \quad \text{and} \quad \lambda_i(s, a) \geq 0.$$

**Case 1:**  $K_d > K_n$ . In this case, the vectors in  $\{\phi(s, a) : (s, a) \in \mathcal{K}\}$  are linearly independent. For ease of presentation and without loss of generality, we assume that  $K_d = K_n + 1$ . This indicates that the matrix  $\Phi \in \mathbb{R}^{K_d \times (|\mathcal{S}| |\mathcal{A}|)}$  whose columns are composed of the feature vectors of all state-action pairs has rank  $K_n$  and is hence not full row rank. This suggests that there exists  $K_n$  linearly independent rows (without loss of generality, we assume they are the first  $K_n$  rows). We can remove the last row from  $\Phi$  to obtain  $\Phi' := \Phi_{1:K_n, :} \in \mathbb{R}^{K_n \times (|\mathcal{S}| |\mathcal{A}|)}$  such that  $\Phi'$  is full row rank. Then we show that we can actually use the columns of  $\Phi'$  as new feature mappings. To see why this is true, note that the last row  $\Phi_{K_n+1, :}$  can be represented as a linear combination of the first  $K_n$  rows, namely there must exist constants  $\{c_k\}_{k=1}^{K_n}$  such that for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\phi_{K_n+1}(s, a) = \sum_{k=1}^{K_n} c_k \phi_k(s, a).$$

Define  $\psi'_k = \psi_k + c_k \psi_{K_n+1}$  for  $k = 1, \dots, K_n$ , we have

$$\begin{aligned} P(s'|s, a) &= \sum_{k=1}^{K_d} \phi_k(s, a) \psi_k(s') = \phi_{K_n+1}(s, a) \psi_{K_n+1}(s') + \sum_{k=1}^{K_n} \phi_k(s, a) \psi_k(s') \\ &= \sum_{k=1}^{K_n} \phi_k(s, a) [\psi_k(s') + c_k \psi_{K_n+1}(s')] = \sum_{k=1}^{K_n} \phi_k(s, a) \psi'_k(s'), \end{aligned}$$

which is linear with respect to the new  $K_n$  dimensional feature vectors. It is also straightforward to check that the new feature mapping satisfies Assumption 1 with the original anchor state-action pairs  $\mathcal{K}$ .

**Case 2:**  $K_d < K_n$ . For ease of presentation and without loss of generality, we assume that  $K_n = K_d + 1$  and that the subspace spanned by the feature vectors of anchor state-action pairs is non-degenerate, i.e., has rank  $K_d$  (otherwise we can use similar method as in Case 1 to further reduce the feature dimension  $K_d$ ). In this case, the matrix  $\Phi_{\mathcal{K}} \in \mathbb{R}^{K_d \times K_n}$  whose columns are composed of the feature vectors of anchor state-action pairs has rank  $K_d$ . We can add  $K_n - K_d = 1$  new row to  $\Phi_{\mathcal{K}}$  to obtain  $\Phi'_{\mathcal{K}} \in \mathbb{R}^{K_n \times K_n}$  such that  $\Phi'_{\mathcal{K}}$  has full rank  $K_n$ . Then we let the columns of  $\Phi'_{\mathcal{K}} = [\phi'(s, a)]_{(s,a) \in \mathcal{K}}$  to be the new feature vectors of the anchor state-action pairs, and define the new feature vectors for all other state-action pairs  $(s, a) \notin \mathcal{K}$  by

$$\phi'(s, a) = \sum_{i:(s_i, a_i) \in \mathcal{K}} \lambda_i(s, a) \phi'(s_i, a_i).$$

We can check that the transition model  $P$  is not changed if we let  $\psi_{K_n}(s') = 0$  for every  $s' \in \mathcal{S}$ . It is also straightforward to check that Assumption 1 is satisfied.

To conclude, when  $K_d \neq K_n$ , we can always construct a new set of feature mappings with dimension  $K_n$  such that: (i) the feature dimension equals to the number of anchor state-action pairs (they are both  $K_n$ ); (ii) the transition model can still be linearly parameterized by this new set of feature mappings; and (iii) the anchor state-action pair assumption (Assumption 1) is satisfied with the original anchor state-action pairs.

## References

- [1] Alekh Agarwal, Sham Kakade, and Lin F Yang. Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pages 67–83. PMLR, 2020.
- [2] Sanjeev Arora, Rong Ge, and Ankur Moitra. Learning topic models—going beyond svd. In *2012 IEEE 53rd annual symposium on foundations of computer science*, pages 1–10. IEEE, 2012.
- [3] Mohammad Gheshlaghi Azar, Rémi Munos, and Hilbert J Kappen. Minimax pac bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, 91(3):325–349, 2013.
- [4] Gen Li, Changxiao Cai, Yuxin Chen, Yuantao Gu, Yuting Wei, and Yuejie Chi. Is q-learning minimax optimal? a tight sample complexity analysis. *arXiv preprint arXiv:2102.06548*, 2021.
- [5] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge University Press, 2018.