

A. Analysis: finite-horizon MDPs

In this section, we present the proof of Theorem 3 — a more general version of Theorem 2 — which accounts for the full ε -range.

Theorem 3. Consider any $\varepsilon \in (0, \frac{1}{1-\gamma}]$. Theorem 2 continues to hold if

$$T \geq \frac{c_3 H^4 (\log^3 T) (\log \frac{|S||\mathcal{A}|T}{\delta})}{\min\{\varepsilon^2, \varepsilon\}} \quad (27)$$

for some sufficiently large universal constant $c_3 > 0$.

A.1. Preliminaries

Let us first introduce several vector and matrix notations that are adopted for the finite-horizon case.

Vector and matrix notation. We use vector $\mathbf{r}_h \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ to represent the reward function r_h at step h . The vectors $\mathbf{V}_h^\pi \in \mathbb{R}^{|\mathcal{S}|}$, $\mathbf{V}_h^* \in \mathbb{R}^{|\mathcal{S}|}$, $\mathbf{Q}_h^\pi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and $\mathbf{Q}_h^* \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ are defined in an analogous manner. Let $\mathbf{Q}_{t,h} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ (resp. $\mathbf{V}_{t,h} \in \mathbb{R}^{|\mathcal{S}|}$) be the estimate $Q_{t,h}$ (resp. $V_{t,h}$) in the t -th iteration at step h , namely

$$\mathbf{Q}_{t,h} = (1 - \eta_t) \mathbf{Q}_{t-1,h} + \eta_t (\mathbf{r}_h + \mathbf{P}_{t,h} \mathbf{V}_{t,h+1}), \quad \forall 1 \leq h \leq H, \quad (28a)$$

$$\mathbf{V}_{t,h} = \max_a \mathbf{Q}_{t,h}, \quad \forall 1 \leq h \leq H. \quad (28b)$$

Here, the maximum in (28b) is taken in an entry-wise manner (cf. (8)). We also use matrix $\mathbf{P}_h \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|}$ to represent the probability transition kernel P_h at step h . Moreover, let the matrix $\mathbf{P}_{t,h} \in \{0, 1\}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|}$ be

$$\mathbf{P}_{t,h}((s, a), s') := \begin{cases} 1, & \text{if } s' = s_{t,h}(s, a), \\ 0, & \text{otherwise.} \end{cases} \quad (29)$$

Similar to the infinite-horizon case, let us first collect a couple of basic facts that will be useful in the proof.

Ranges of $\mathbf{Q}_{t,h}$ and $\mathbf{V}_{t,h}$. We shall start with some simple bounds for $\mathbf{Q}_{t,h}$ and $\mathbf{V}_{t,h}$. Lemma 3 (below) demonstrates that the estimates for the Q-function and the value function are bounded as long as they are properly initialized.

Lemma 3. Suppose that $0 \leq \eta_t \leq 1$ for all $t \geq 0$. Assume that $\mathbf{Q}_{H+1}^0 = \mathbf{V}_{H+1}^0 = \mathbf{0}$. Then for all $t \geq 0$ and $1 \leq h \leq H+1$, one has

$$\mathbf{0} \leq \mathbf{Q}_{t,h} \leq (H+1-h)\mathbf{1} \quad \text{and} \quad \mathbf{0} \leq \mathbf{V}_{t,h} \leq (H+1-h)\mathbf{1}. \quad (30)$$

Proof. We can use the induction argument to prove this. First, our initialization obeys (30) for $t = 0$ and $h = H+1$. Next, suppose that (30) is true for $t-1$ and $h+1$. By the update rule (19), it is straightforward to compute

$$\mathbf{Q}_{t,h} = (1 - \eta_t) \mathbf{Q}_{t-1,h} + \eta_t (\mathbf{r}_h + \mathbf{P}_{t,h} \mathbf{V}_{t,h+1}) \geq \mathbf{0},$$

and

$$\begin{aligned} \mathbf{Q}_{t,h} &= (1 - \eta_t) \mathbf{Q}_{t-1,h} + \eta_t (\mathbf{r}_h + \mathbf{P}_{t,h} \mathbf{V}_{t,h+1}) \\ &\leq (1 - \eta_t) \|\mathbf{Q}_{t-1,h}\|_\infty \mathbf{1} + \eta_t (\|\mathbf{r}_h\|_\infty + \|\mathbf{P}_{t,h}\|_1 \|\mathbf{V}_{t,h+1}\|_\infty) \mathbf{1} \\ &\leq (1 - \eta_t) (H+1-h) \mathbf{1} + \eta_t (1 + (H-h)) \mathbf{1} = (H+1-h) \mathbf{1}. \end{aligned}$$

where we use the facts $\mathbf{r}_h \leq \mathbf{1}$ and \mathbf{P}_h^t is a probability transition kernel. In addition, since $V_h^t(s) := \max_a Q_h^t(s, a)$ for all $t \geq 0$, $1 \leq h \leq H+1$ and $s \in \mathcal{S}$, it is easy to see that $\mathbf{0} \leq \mathbf{V}_{t,h} \leq (H+1-h)\mathbf{1}$. This completes the proof for (30). \square

It immediately follows from Lemma 3 that for all $t \geq 0$ and $1 \leq h \leq H+1$,

$$-H\mathbf{1} \leq -\mathbf{Q}_h^* \leq \mathbf{Q}_{t,h} - \mathbf{Q}_h^* \leq \mathbf{Q}_{t,h} \leq H\mathbf{1}.$$

Combined with the fact $\|Q^*\|_\infty \leq H$, one further has

$$\max_{t \geq 0, 1 \leq h \leq H+1} \|Q_{t,h} - Q^*\|_\infty \leq H.$$

This suggests we can focus on the case where $\varepsilon \leq H$ and the claimed iteration number in (27) satisfies

$$T = \frac{c_3 H^4 \log^3 T \log \frac{|S||A|T}{\delta}}{\min\{\varepsilon^2, \varepsilon\}} \geq \frac{c_3 H^4 \log^3 T \log \frac{|S||A|T}{\delta}}{\varepsilon} \geq c_3 H^3 \log^3 T \log \frac{|S||A|T}{\delta}. \quad (31)$$

Several facts regarding the learning rates. Next, we present a few useful bounds regarding the learning rates $\eta_i^{(t)}$ defined in the same way as (25). From the assumption (11a) and the bound (31), it is easily seen that the step size obeys

$$\frac{H \log^2 T}{2c_1 T} \leq \frac{1}{1 + \frac{c_1 T}{H \log^2 T}} \leq \eta_t \leq \frac{1}{1 + \frac{c_2 t}{H \log^2 T}} \leq \frac{H \log^2 T}{c_2 t}. \quad (32)$$

Let us set

$$\beta := \frac{c_4}{H} \quad (33)$$

for some sufficiently small constant $c_4 > 0$. In what follows, we present two upper bounds of $\eta_i^{(t)}$ for any $t \geq \frac{T}{c_2 \log H}$.

- For any $0 \leq i \leq (1 - \beta)t$, one know from (32) and $T \geq H^2$ (cf. (11b)) that

$$\eta_i^{(t)} \leq \left(1 - \frac{H \log^2 T}{2c_1 T}\right)^{\beta t} \leq \left(1 - \frac{H \log^2 T}{2c_1 T}\right)^{\frac{c_4 T}{c_2 H (\log H)}} < \frac{1}{2T^2} \quad (34)$$

where the last step holds provided $c_1 c_2 \leq c_4/4$.

- Turning to $i > (1 - \beta)t \geq t/2$, we can use the condition $t \geq \frac{T}{c_2 \log H}$ to bound

$$\eta_i^{(t)} \leq \eta_i \leq \frac{H \log^2 T}{c_2 i} < \frac{2H \log^2 T}{c_2 t} \leq \frac{2H \log^2 T}{T / \log H} \leq \frac{2H \log^3 T}{T}. \quad (35)$$

Moreover, the sum of $\eta_i^{(t)}$ continues to satisfy

$$\sum_{i=0}^t \eta_i^{(t)} = \prod_{j=1}^t (1 - \eta_j) + \eta_1 \prod_{j=2}^t (1 - \eta_j) + \eta_2 \prod_{j=3}^t (1 - \eta_j) + \cdots + \eta_{t-1} (1 - \eta_t) + \eta_t = 1. \quad (36)$$

A.2. Proof of Theorem 3

We are now in the position to prove Theorem 3. For convenience of notation, we shall define

$$\Delta_{t,h} := Q_{t,h} - Q_h^*.$$

In addition, let π_t denote the policy such that for any state-action-horizon pair (s, a, h) ,¹

$$\pi_t(a | s, h) := \begin{cases} 1, & \text{if } a = \min \{a' \mid Q_{t,h}(s, a') = \max_{a''} Q_{t,h}(s, a'')\}, \\ 0, & \text{else.} \end{cases} \quad (37)$$

Namely, for any $s \in \mathcal{S}$ and $1 \leq h \leq H + 1$, the policy π_t chooses the smallest indexed action that achieves the largest Q-value in the estimate $Q_{t,h}(s, \cdot)$. It immediately follows that

$$Q_{t,h}(s, \pi_t(s, h)) = V_{t,h}(s) \quad \text{and} \quad P_h V_{t,h+1} = P_h^{\pi_t} Q_{t,h+1} \geq P_h^\pi Q_{t,h+1} \quad \text{for any } \pi, \quad (38)$$

where P^π is defined in (14).

¹If there is only a single action that satisfies $Q_{t,h}(s, a') = \max_{a''} Q_{t,h}(s, a'')$, then $\pi_t(a | s, h) = 1$ if and only if $a = \arg \min_{a'} Q_{t,h}(s, a')$ and 0 otherwise.

A.2.1. KEY DECOMPOSITION

We first make the following elementary decomposition:

$$\begin{aligned}
 \Delta_{t,h} &= \mathbf{Q}_{t,h} - \mathbf{Q}_h^* = (1 - \eta_t)\mathbf{Q}_{t-1,h} + \eta_t(\mathbf{r}_h + \mathbf{P}_{t,h}\mathbf{V}_{t,h+1}) - \mathbf{Q}_h^* \\
 &= (1 - \eta_t)(\mathbf{Q}_{t-1,h} - \mathbf{Q}_h^*) + \eta_t(\mathbf{r}_h + \mathbf{P}_{t,h}\mathbf{V}_{t,h+1} - \mathbf{Q}_h^*) \\
 &= (1 - \eta_t)\Delta_{t-1,h} + \eta_t(\mathbf{P}_{t,h}\mathbf{V}_{t,h+1} - \mathbf{P}_h\mathbf{V}_{h+1}^*) \\
 &= (1 - \eta_t)\Delta_{t-1,h} + \eta_t\{\mathbf{P}_h(\mathbf{V}_{t,h+1} - \mathbf{V}_{h+1}^*) + (\mathbf{P}_{t,h} - \mathbf{P}_h)\mathbf{V}_{t,h+1}\}.
 \end{aligned} \tag{39}$$

Similar to (22), one can use (38) to control the quantity $\mathbf{P}_h(\mathbf{V}_{t,h+1} - \mathbf{V}_{h+1}^*)$ by

$$\mathbf{P}_h(\mathbf{V}_{t,h+1} - \mathbf{V}_{h+1}^*) = \mathbf{P}_h^{\pi_t}\mathbf{Q}_{t,h+1} - \mathbf{P}_h^{\pi^*}\mathbf{Q}_{h+1}^* \leq \mathbf{P}_h^{\pi_t}\mathbf{Q}_{t,h+1} - \mathbf{P}_h^{\pi_t}\mathbf{Q}_{h+1}^* = \mathbf{P}_h^{\pi_t}\Delta_{t,h+1}, \tag{40a}$$

$$\mathbf{P}_h(\mathbf{V}_{t,h+1} - \mathbf{V}_{h+1}^*) = \mathbf{P}_h^{\pi_t}\mathbf{Q}_{t,h+1} - \mathbf{P}_h^{\pi^*}\mathbf{Q}_{h+1}^* \geq \mathbf{P}_h^{\pi^*}\mathbf{Q}_{t,h+1} - \mathbf{P}_h^{\pi^*}\mathbf{Q}_{h+1}^* = \mathbf{P}_h^{\pi^*}\Delta_{t,h+1}, \tag{40b}$$

Combining (40) with (39) yields

$$\begin{aligned}
 \Delta_{t,h} &\leq (1 - \eta_t)\Delta_{t-1,h} + \eta_t\{\mathbf{P}_h^{\pi_t}\Delta_{t,h+1} + (\mathbf{P}_{t,h} - \mathbf{P}_h)\mathbf{V}_{t,h+1}\}; \\
 \Delta_{t,h} &\geq (1 - \eta_t)\Delta_{t-1,h} + \eta_t\{\mathbf{P}_h^{\pi^*}\Delta_{t,h+1} + (\mathbf{P}_{t,h} - \mathbf{P}_h)\mathbf{V}_{t,h+1}\}.
 \end{aligned}$$

We can then apply this relation recursively to reach

$$\Delta_{t,h} \leq \eta_0^{(t)}\Delta_{0,h} + \sum_{i=1}^t \eta_i^{(t)}(\mathbf{P}_h^{\pi_i}\Delta_{i,h+1} + (\mathbf{P}_{i,h} - \mathbf{P}_h)\mathbf{V}_{i,h+1}), \tag{41a}$$

$$\Delta_{t,h} \geq \eta_0^{(t)}\Delta_{0,h} + \sum_{i=1}^t \eta_i^{(t)}(\mathbf{P}_h^{\pi_i^*}\Delta_{i,h+1} + (\mathbf{P}_{i,h} - \mathbf{P}_h)\mathbf{V}_{i,h+1}). \tag{41b}$$

In the following, we shall use (41) to upper and lower bound $\Delta_{t,h}$ individually.

 A.2.2. UPPER BOUNDING $\Delta_{t,h}$

Let us first upper bound $\Delta_{t,h}$ for $t \geq \frac{T}{c_2 \log H}$. In view of (41a), we further decompose its right-hand side as

$$\begin{aligned}
 \Delta_{t,h} &\leq \underbrace{\eta_0^{(t)}\Delta_{0,h} + \sum_{i=1}^{(1-\beta)t} \eta_i^{(t)}(\mathbf{P}_h^{\pi_i}\Delta_{i,h+1} + (\mathbf{P}_{i,h} - \mathbf{P}_h)\mathbf{V}_{i,h+1})}_{=: \zeta_{t,h}} \\
 &\quad + \underbrace{\sum_{i=(1-\beta)t+1}^t \eta_i^{(t)}(\mathbf{P}_{i,h} - \mathbf{P}_h)\mathbf{V}_{i,h+1}}_{=: \xi_{t,h}} + \sum_{i=(1-\beta)t+1}^t \eta_i^{(t)}\mathbf{P}_h^{\pi_i}\Delta_{i,h+1}
 \end{aligned} \tag{42}$$

where we recall that $\beta := \frac{c_4}{H}$ defined in (33).

Step 1: bounding $\zeta_{t,h}$. From the upper bounds (34) for $\eta_i^{(t)}$, it is straightforward to control $\zeta_{t,h}$ as follows:

$$\begin{aligned}
 \|\zeta_{t,h}\|_\infty &\leq \eta_0^{(t)}\|\Delta_{0,h}\|_\infty + t \max_{i \leq (1-\beta)t} \eta_i^{(t)} \max_{1 \leq i \leq (1-\beta)t} (\|\mathbf{P}_h^{\pi_i}\Delta_{i,h+1}\|_\infty + \|\mathbf{P}_{i,h}\mathbf{V}_{i,h+1}\|_\infty + \|\mathbf{P}_h\mathbf{V}_{i,h+1}\|_\infty) \\
 &\leq \eta_0^{(t)}\|\Delta_{0,h}\|_\infty + t \max_{i \leq (1-\beta)t} \eta_i^{(t)} \max_{1 \leq i \leq (1-\beta)t} \left\{ \|\mathbf{P}_h^{\pi_i}\|_1 \|\Delta_{i,h+1}\|_\infty + (\|\mathbf{P}_{i,h}\|_1 + \|\mathbf{P}_h\|_1) \|\mathbf{V}_{i,h+1}\|_\infty \right\} \\
 &\stackrel{(i)}{=} \eta_0^{(t)}\|\Delta_{0,h}\|_\infty + t \max_{i \leq (1-\beta)t} \eta_i^{(t)} \max_{1 \leq i \leq (1-\beta)t} (\|\Delta_{i,h+1}\|_\infty + 2\|\mathbf{V}_{i,h+1}\|_\infty) \\
 &\stackrel{(ii)}{\leq} \frac{1}{2T^2} \cdot H + \frac{1}{2T^2} \cdot t \cdot 3H \\
 &\leq \frac{2H}{T}.
 \end{aligned}$$

Here, (i) relies on $\|\mathbf{P}_h^{\pi_i}\|_1 = \|\mathbf{P}_{i,h}\|_1 = \|\mathbf{P}_h\|_1 = 1$ since they are all probability transition matrices; (ii) holds due to (34).

Step 2: bounding $\xi_{t,h}$. Observe that $\xi_{t,h}$ is a sum of martingale differences, namely

$$\xi_{t,h} = \sum_{i=(1-\beta)t+1}^t z_{i,h} \quad \text{with } z_{i,h} := \eta_i^{(t)} (\mathbf{P}_{i,h} - \mathbf{P}_h) \mathbf{V}_{i,h+1},$$

where the $z_{i,h}$'s satisfy

$$\mathbb{E} [z_{i,h} \mid \mathbf{V}_{i,h+1}, \dots, \mathbf{V}_{0,h+1}] = \mathbf{0}.$$

This suggests we can invoke Freedman's inequality (see Lemma 4) to control $\xi_{t,h}$ for any t such that $\frac{T}{c_2 \log H} \leq t \leq T$.

- First, it is straightforward to bound

$$\begin{aligned} B &:= \max_{(1-\beta)t < i \leq t} \|z_{i,h}\|_\infty \leq \max_{(1-\beta)t < i \leq t} \|\eta_i^{(t)} (\mathbf{P}_{i,h} - \mathbf{P}_h) \mathbf{V}_{i,h+1}\|_\infty \\ &\leq \max_{(1-\beta)t < i \leq t} \eta_i^{(t)} (\|\mathbf{P}_{i,h}\|_1 + \|\mathbf{P}_h\|_1) \|\mathbf{V}_{i,h+1}\|_\infty \leq \frac{4H^2 \log^3 T}{T}. \end{aligned} \quad (43)$$

where the last step arises from (35), Lemma 3, and the fact $\|\mathbf{P}_{i,h}\|_1 = \|\mathbf{P}_h\|_1 = 1$.

- Next, recall the notation $\text{Var}_{\mathcal{P}}(\mathbf{z})$ in (16). One can compute

$$\begin{aligned} \mathbf{W}_t &:= \sum_{i=(1-\beta)t+1}^t \text{Var}(z_{i,h} \mid \mathbf{V}_{i,h+1}, \dots, \mathbf{V}_{0,h+1}) = \sum_{i=(1-\beta)t+1}^t (\eta_i^{(t)})^2 \text{Var}_{\mathcal{P}_h}(\mathbf{V}_{i,h+1}) \\ &\leq \left(\max_{(1-\beta)t < i \leq t} \eta_i^{(t)} \right) \left(\sum_{i=(1-\beta)t+1}^t \eta_i^{(t)} \right) \max_{(1-\beta)t < i < t} \text{Var}_{\mathcal{P}_h}(\mathbf{V}_{i,h+1}) \\ &\leq \frac{2H \log^3 T}{T} \max_{(1-\beta)t < i < t} \text{Var}_{\mathcal{P}_h}(\mathbf{V}_{i,h+1}), \end{aligned} \quad (44)$$

where the last inequality relies on (35) and (36).

- Additionally, we can use Lemma 3 to further bound \mathbf{W}_t

$$|\mathbf{W}_t| \leq \frac{2H \log^3 T}{T} \cdot H^2 \mathbf{1} = \frac{2H^3 \log^3 T}{T} \mathbf{1} =: \sigma^2 \mathbf{1}.$$

In particular, we know that

$$\frac{\sigma^2}{2^K} \leq \frac{2H \log^3 T}{T} \quad (45)$$

where $K := \lceil 2 \log H \rceil$.

With the above bounds in place, we apply the Freedman inequality in Lemma 4 and the union bound to find: with probability at least $1 - \frac{\delta}{TH}$,

$$\begin{aligned} |\xi_{t,h}| &\leq \sqrt{8 \left(\mathbf{W}_t + \frac{\sigma^2}{2^K} \mathbf{1} \right) \log \frac{|\mathcal{S}| |\mathcal{A}| THK}{\delta}} + \left(\frac{4}{3} B \log \frac{|\mathcal{S}| |\mathcal{A}| THK}{\delta} \right) \cdot \mathbf{1} \\ &\stackrel{(i)}{\leq} \sqrt{16 \left(\mathbf{W}_t + \frac{2H \log^3 T}{T} \mathbf{1} \right) \log \frac{|\mathcal{S}| |\mathcal{A}| T}{\delta}} + \left(3B \log \frac{|\mathcal{S}| |\mathcal{A}| T}{\delta} \right) \cdot \mathbf{1} \\ &\stackrel{(ii)}{\leq} \sqrt{\frac{32H (\log^3 T) (\log \frac{|\mathcal{S}| |\mathcal{A}| T}{\delta})}{T} \left(\max_{(1-\beta)t < i < t} \text{Var}_{\mathcal{P}_h}(\mathbf{V}_{i,h+1}) + 1 \right)} + \frac{12H^2 (\log^3 T) (\log \frac{|\mathcal{S}| |\mathcal{A}| T}{\delta})}{T} \mathbf{1} \end{aligned}$$

Here, (i) arises from (45) and $\log \frac{|\mathcal{S}| |\mathcal{A}| THK}{\delta} \leq 2 \log \frac{|\mathcal{S}| |\mathcal{A}| T}{\delta}$ (which holds due to (31)); (ii) makes use of the relation (43) and (44).

Step 3: using the bounds on $\zeta_{t,h}$ and $\xi_{t,h}$ to control $\Delta_{t,h}$. Let us define

$$\varphi_{t,h} := c_5 \frac{H \log^3 T \log \frac{|S||A|T}{\delta}}{T} \left(\max_{\frac{t}{2} \leq i < t} \text{Var}_{P_h}(\mathbf{V}_{i,h+1}) + 1 \right) \quad (46)$$

for some sufficiently large constant $c_5 > 0$. In view of bounds for $\zeta_{t,h}$ and $\xi_{t,h}$, the following holds with probability exceeding $1 - \delta$: for all $\frac{2t}{3} \leq k \leq t$ and $1 \leq h \leq H$,

$$|\zeta_{k,h}| + |\xi_{k,h}| \leq \sqrt{\varphi_{t,h}}. \quad (47)$$

Inserting (47) into (42) reveals

$$\Delta_{t,h} \leq \sqrt{\varphi_{t,h}} + \sum_{i=(1-\beta)t+1}^t \eta_i^{(t)} P_h^{\pi_i} \Delta_{i,h+1}. \quad (48)$$

We now define a sequence $\{\alpha_i^{(t)}\}_i$ as follows

$$\alpha_i^{(t)} := \frac{\eta_i^{(t)}}{\sum_{j=(1-\beta)t+1}^t \eta_j^{(t)}}, \quad 0 \leq i \leq t.$$

It is easy to check that for any t , the sequence $\{\alpha_i^{(t)}\}_i$ satisfies

$$\alpha_i^{(t)} \geq \eta_i^{(t)} \quad \text{and} \quad \sum_{i=(1-\beta)t+1}^t \alpha_i^{(t)} = 1 \quad (49)$$

where the first inequality results from (36). This enables us to rewrite (48) as

$$\Delta_{k,h} \leq \sqrt{\varphi_{t,h}} + \sum_{i_h=(1-\beta)k+1}^k \eta_{i_h}^{(k)} P_h^{\pi_{i_h}} \Delta_{i_h,h+1} = \sum_{i_h=(1-\beta)k+1}^k \left(\alpha_{i_h}^{(k)} \sqrt{\varphi_{t,h}} + \eta_{i_h}^{(k)} P_h^{\pi_{i_h}} \Delta_{i_h,h+1} \right). \quad (50)$$

for all $2t/3 \leq k \leq t$. By the definition of β (cf. (33)), one has $(1-\beta)t \geq 2t/3$. We can then exploit this relation recursively to obtain

$$\begin{aligned} \Delta_{t,h} &\leq \sum_{i_h=(1-\beta)t+1}^t \left(\alpha_{i_h}^{(t)} \sqrt{\varphi_{t,h}} + \eta_{i_h}^{(t)} P_h^{\pi_{i_h}} \Delta_{i_h,h+1} \right) \\ &\leq \sum_{i_h=(1-\beta)t+1}^t \left\{ \alpha_{i_h}^{(t)} \sqrt{\varphi_{t,h}} + \eta_{i_h}^{(t)} P_h^{\pi_{i_h}} \sum_{i_{h+1}=(1-\beta)i_h+1}^{i_h} \left(\alpha_{i_{h+1}}^{(i_h)} \sqrt{\varphi_{t,h+1}} + \eta_{i_{h+1}}^{(i_h)} P_{h+1}^{\pi_{i_{h+1}}} \Delta_{i_{h+1},h+2} \right) \right\} \\ &\stackrel{(i)}{\leq} \sum_{i_h=(1-\beta)t+1}^t \alpha_{i_h}^{(t)} \sqrt{\varphi_{t,h}} + \sum_{i_h=(1-\beta)t+1}^t \sum_{i_{h+1}=(1-\beta)i_h+1}^{i_h} \alpha_{i_h}^{(t)} \alpha_{i_{h+1}}^{(i_h)} P_h^{\pi_{i_h}} \sqrt{\varphi_{t,h+1}} \\ &\quad + \sum_{i_h=(1-\beta)t+1}^t \sum_{i_{h+1}=(1-\beta)i_h+1}^{i_h} \eta_{i_h}^{(t)} \eta_{i_{h+1}}^{(i_h)} \prod_{k=h}^{h+1} P_k^{\pi_{i_k}} \Delta_{i_{h+1},h+2} \\ &\stackrel{(ii)}{=} \sum_{i_h=(1-\beta)t+1}^t \sum_{i_{h+1}=(1-\beta)i_h+1}^{i_h} \alpha_{i_h}^{(t)} \alpha_{i_{h+1}}^{(i_h)} \left\{ \sqrt{\varphi_{t,h}} + P_h^{\pi_{i_h}} \sqrt{\varphi_{t,h+1}} \right\} \\ &\quad + \sum_{i_h=(1-\beta)t+1}^t \sum_{i_{h+1}=(1-\beta)i_h+1}^{i_h} \eta_{i_h}^{(t)} \eta_{i_{h+1}}^{(i_h)} \prod_{k=h}^{h+1} P_k^{\pi_{i_k}} \Delta_{i_{h+1},h+2}, \end{aligned} \quad (51)$$

where (i) relies on $\eta_{i_h}^{(t)} \leq \alpha_{i_h}^{(t)}$ in (49), and (ii) holds since $\sum_{i_{h+1}=(1-\beta)i_h+1}^{i_h} \alpha_{i_{h+1}}^{(i_h)} = 1$ by (49).

Our proof strategy is applying (50) recursively to control $\Delta_{t,h}$ for all $1 \leq h \leq H$. Towards this, we need some preparation beforehand. First, let us define

$$\alpha_{\{i_k\}_{k=h}^H} := \alpha_{i_h}^{(t)} \alpha_{i_{h+1}}^{(i_h)} \dots \alpha_{i_H}^{(i_{H-1})} \geq 0, \quad 1 \leq h \leq H \quad (52)$$

for any $t \geq i_h \geq i_{h+1} \geq \dots \geq i_H$. By (49), one has

$$\alpha_{\{i_k\}_{k=h}^H} \geq \eta_{i_h}^{(t)} \eta_{i_{h+1}}^{(i_h)} \dots \eta_{i_H}^{(i_{H-1})}. \quad (53)$$

Next, let us define the index set

$$\mathcal{I}_{t,h} := \left\{ (i_h, \dots, i_H) \mid (1-\beta)t \leq i_h \leq t-1, (1-\beta)i_j \leq i_{j+1} \leq i_j-1, \forall h \leq j < H \right\}, \quad (54)$$

which satisfies

$$\sum_{(i_h, \dots, i_H) \in \mathcal{I}_{t,h}} \alpha_{\{i_k\}_{k=h}^H} = 1. \quad (55)$$

In addition, as $\beta := c_4/H$ for some sufficiently small constant $c_4 > 0$, we know that

$$(1-\beta)^H = \left(1 - \frac{c_4}{H}\right)^H \geq \frac{2}{3},$$

and consequently

$$i_h \geq i_{h+1} \geq \dots \geq i_H \geq (1-\beta)^H t \geq 2t/3 \quad \text{for all } (i_h, \dots, i_H) \in \mathcal{I}_{t,h}.$$

With these in place, we now invoke the relation (50) in a recursive manner to obtain

$$\begin{aligned} \Delta_{t,h} &\leq \sum_{(i_h, \dots, i_H) \in \mathcal{I}_{t,h}} \alpha_{\{i_k\}_{k=h}^H} \left\{ \sqrt{\varphi_{t,h}} + \sum_{j=h+1}^H \prod_{k=h}^{j-1} P_k^{\pi_{i_k}} \sqrt{\varphi_{t,j}} \right\} \\ &\leq \max_{(i_h, \dots, i_H) \in \mathcal{I}_{t,h}} \left\{ \sum_{j=h}^H \prod_{k=h}^{j-1} P_k^{\pi_{i_k}} \sqrt{\varphi_{t,j}} \right\} \end{aligned} \quad (56)$$

for all $t \geq \frac{T}{c_2 \log \frac{1}{1-\gamma}}$. Here, the last step arises from the fact that $\sum_{(i_h, \dots, i_H) \in \mathcal{I}_{t,h}} \alpha_{\{i_k\}_{k=h}^H} = 1$ (cf. (55)). One can upper bound the entrywise square of the quantity in curly braces as follows

$$\begin{aligned} \left| \sum_{j=h}^H \prod_{k=h}^{j-1} P_k^{\pi_{i_k}} \sqrt{\varphi_{t,j}} \right|^2 &\stackrel{(i)}{\leq} \left| \sum_{j=h}^H \sqrt{\prod_{k=h}^{j-1} P_k^{\pi_{i_k}} \varphi_{t,j}} \right|^2 \stackrel{(ii)}{\leq} H \sum_{j=h}^H \prod_{k=h}^{j-1} P_k^{\pi_{i_k}} \varphi_{t,j} \\ &\stackrel{(iii)}{\lesssim} H \sum_{j=h}^H \prod_{k=h}^{j-1} P_k^{\pi_{i_k}} \frac{H \log^3 T \log \frac{|S||A|T}{\delta}}{T} \left(\max_{\frac{t}{2} \leq i < t} \text{Var}_{P_j}(\mathbf{V}_{i,j+1}) + 1 \right) \\ &\stackrel{(iv)}{\leq} \frac{H^2 (\log^3 T) (\log \frac{|S||A|T}{\delta})}{T} \sum_{j=h}^H \prod_{k=h}^{j-1} P_k^{\pi_{i_k}} \max_{\frac{t}{2} \leq i < t} \text{Var}_{P_j}(\mathbf{V}_{i,j+1}) + \frac{H^3 (\log^3 T) (\log \frac{|S||A|T}{\delta})}{T} \mathbf{1} \\ &\stackrel{(v)}{\lesssim} \frac{H^4 (\log^3 T) (\log \frac{|S||A|T}{\delta})}{T} \left(1 + \max_{\frac{t}{2} \leq i \leq t, k > h} \|\Delta_{i,k}\|_{\infty} \right) \mathbf{1}. \end{aligned} \quad (57)$$

where (i) arises from the Jensen inequality and the fact $\prod_{k=h}^{j-1} P_k^{\pi_{i_k}}$ is a probability transition matrix; (ii) relies on the Cauchy-Schwarz inequality; (iii) is due to the definition of $\varphi_{t,h}$; (iv) holds since $\prod_{k=h}^{j-1} P_k^{\pi_{i_k}} \mathbf{1} = \mathbf{1}$; (v) is valid as long as Lemma 4 holds.

Lemma 4. *One has*

$$\sum_{j=h}^H \prod_{k=h}^{j-1} P_k^{\pi_{i_k}} \max_{\frac{t}{2} \leq i < t} \text{Var}_{P_j}(\mathbf{V}_{i,j+1}) \lesssim H^2 \left(1 + \max_{\frac{t}{2} \leq i \leq t, k > h} \|\Delta_{i,k}\|_{\infty} \right) \mathbf{1}. \quad (58)$$

Proof. See Section A.2.5. □

Plugging (57) back into (56) reveals that the following holds simultaneously for all $t \geq \frac{T}{c_2 \log H}$ with probability at least $1 - \delta$:

$$\Delta_{t,h} \lesssim \sqrt{\frac{H^4 (\log^3 T) (\log \frac{|S||A|T}{\delta})}{T} \left(1 + \max_{\frac{t}{2} \leq i \leq t, k > h} \|\Delta_{i,k}\|_\infty\right) \mathbf{1}}. \quad (59)$$

A.2.3. LOWER BOUNDING $\Delta_{t,h}$

In this section, we proceed to lower bound $\Delta_{t,h}$. Invoking a similar argument for (56) and replacing π_i with π^* , we can derive

$$\Delta_{t,h} \geq - \max_{(i_h, \dots, i_H) \in \mathcal{I}_{t,h}} \left\{ \sum_{j=h}^H \prod_{k=h}^{j-1} P_k^{\pi^*} \sqrt{\varphi_{t,j}} \right\}.$$

One can further apply an analogous argument for (57) to bound the right-hand side as follows

$$\left| \sum_{j=h}^H \prod_{k=h}^{j-1} P_k^{\pi^*} \sqrt{\varphi_{t,j}} \right|^2 \lesssim \frac{H^4 (\log^3 T) (\log \frac{|S||A|T}{\delta})}{T} \left(1 + \max_{\frac{t}{2} \leq i \leq t, k > h} \|\Delta_{i,k}\|_\infty\right) \mathbf{1}.$$

Consequently, we find that with probability at least $1 - \delta$,

$$\Delta_{t,h} \gtrsim - \sqrt{\frac{H^4 (\log^3 T) (\log \frac{|S||A|T}{\delta})}{T} \left(1 + \max_{\frac{t}{2} \leq i \leq t, k > h} \|\Delta_{i,k}\|_\infty\right) \mathbf{1}} \quad (60)$$

holds simultaneously for all $t \geq \frac{T}{c_2 \log H}$.

A.2.4. COMBINING OUR UPPER AND LOWER BOUNDS ON $\Delta_{t,h}$

Taking (59) and (60) together, we know that with probability greater than $1 - 2\delta$,

$$\|\Delta_{t,h}\|_\infty \lesssim \sqrt{\frac{H^4 (\log^3 T) (\log \frac{|S||A|T}{\delta})}{T} \left(1 + \max_{\frac{t}{2} \leq i \leq t, k > h} \|\Delta_{i,k}\|_\infty\right)}, \quad (61)$$

holds simultaneously for all $t \geq \frac{T}{c_2 \log H}$. As a result, the claim in Theorem 2 immediately follows from the same argument for the infinite-horizon case in Appendix 4.2.3, which we omit for the sake of conciseness.

A.2.5. PROOF OF LEMMA 4

We shall invoke a similar argument for Li et al. (2021a, Lemma 5) to establish Lemma 4. For the sake of conciseness, we omit some details of proof.

To begin with, we can argue analogously as for Li et al. (2021a, (64)) to show that for any $1 \leq j \leq H$,

$$\max_{\frac{t}{2} \leq i < t} \text{Var}_{P_j}(\mathbf{V}_{i,j+1}) - \text{Var}_{P_j}(\mathbf{V}_{j+1}^*) \leq 4H \max_{\frac{t}{2} \leq i < t} \|\Delta_{i,j+1}\|_\infty \mathbf{1}. \quad (62)$$

As a consequence, one can bound the left-hand side of (58) by

$$\sum_{j=h}^H \prod_{k=h}^{j-1} P_k^{\pi_{i_k}} \max_{\frac{t}{2} \leq i < t} \text{Var}_{P_j}(\mathbf{V}_{i,j+1}) \leq \sum_{j=h}^H \prod_{k=h}^{j-1} P_k^{\pi_{i_k}} \text{Var}_{P_j}(\mathbf{V}_{j+1}^*) + 4H^2 \max_{\frac{t}{2} \leq i \leq t, k > h} \|\Delta_{i,k}\|_\infty \mathbf{1}. \quad (63)$$

Now it remains to control the first term on the right-hand side of (63). Towards this, we can first bound similarly as in Li et al. (2021a, (67)) to obtain

$$\begin{aligned} \text{Var}_{\mathcal{P}_j}(\mathbf{V}_{j+1}^*) &= \mathcal{P}_j(\mathbf{V}_{j+1}^* \circ \mathbf{V}_{j+1}^*) - (\mathcal{P}_j \mathbf{V}_{j+1}^*) \circ (\mathcal{P}_j \mathbf{V}_{j+1}^*) \\ &\leq (\mathcal{P}_j^{\pi_{i_j}}(\mathbf{Q}_{j+1}^* \circ \mathbf{Q}_{j+1}^*) - \mathbf{Q}_j^* \circ \mathbf{Q}_j^*) + 2\mathbf{Q}_j^* \circ \mathbf{r}_j + 4H \max_{\frac{t}{2} \leq i < t} \|\Delta_{i,j+1}\|_\infty \mathbf{1}. \end{aligned} \quad (64)$$

With the estimate for $\text{Var}_{\mathcal{P}_j}(\mathbf{V}_{j+1}^*)$ in place, one can invoke the same argument for Li et al. (2021a, (68)) to reach

$$\begin{aligned} \sum_{j=h}^H \prod_{k=h}^{j-1} \mathcal{P}_k^{\pi_{i_k}} \text{Var}_{\mathcal{P}_j}(\mathbf{V}_{j+1}^*) &\leq \sum_{j=h}^H \prod_{k=h}^{j-1} \mathcal{P}_k^{\pi_{i_k}} (\mathcal{P}_j^{\pi_{i_j}}(\mathbf{Q}_{j+1}^* \circ \mathbf{Q}_{j+1}^*) - \mathbf{Q}_j^* \circ \mathbf{Q}_j^*) \\ &\quad + \sum_{j=h}^H \prod_{k=h}^{j-1} \mathcal{P}_k^{\pi_{i_k}} (2\mathbf{Q}_j^* \circ \mathbf{r}_j + 4H \max_{\frac{t}{2} \leq i < t} \|\Delta_{i,j+1}\|_\infty \mathbf{1}) \\ &\leq 4H^2 (1 + 4 \max_{\frac{t}{2} \leq i \leq t, k > h} \|\Delta_{i,k}\|_\infty) \mathbf{1}. \end{aligned}$$

Plugging the above bounds back into (64) immediately establishes the claimed bound (58).

B. Freedman's inequality

The analysis of this work relies heavily on Freedman's inequality (Freedman, 1975), which is an extension of the Bernstein inequality and allows one to establish concentration results for martingales. For ease of presentation, we include a user-friendly version of Freedman's inequality as follows.

Theorem 4. *Suppose that $Y_n = \sum_{k=1}^n X_k \in \mathbb{R}$, where $\{X_k\}$ is a real-valued scalar sequence obeying*

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

Define

$$W_n := \sum_{k=1}^n \mathbb{E}_{k-1}[X_k^2],$$

where we write \mathbb{E}_{k-1} for the expectation conditional on $\{X_j\}_{j:j < k}$. Then for any given $\sigma^2 \geq 0$, one has

$$\mathbb{P}\{|Y_n| \geq \tau \text{ and } W_n \leq \sigma^2\} \leq 2 \exp\left(-\frac{\tau^2/2}{\sigma^2 + R\tau/3}\right). \quad (65)$$

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

$$|Y_n| \leq \sqrt{8 \max\left\{W_n, \frac{\sigma^2}{2K}\right\} \log \frac{2K}{\delta}} + \frac{4}{3} R \log \frac{2K}{\delta}. \quad (66)$$

Proof. See (Freedman, 1975; Tropp, 2011) for the proof of (65). As an immediate consequence of (65), one has

$$\mathbb{P}\left\{|Y_n| \geq \sqrt{4\sigma^2 \log \frac{2}{\delta}} + \frac{4}{3} R \log \frac{2}{\delta} \text{ and } W_n \leq \sigma^2\right\} \leq \delta. \quad (67)$$

Next, we turn attention to (66). Consider any positive integer K . As can be easily seen, the event

$$\mathcal{H}_K := \left\{|Y_n| \geq \sqrt{8 \max\left\{W_n, \frac{\sigma^2}{2K}\right\} \log \frac{2K}{\delta}} + \frac{4}{3} R \log \frac{2K}{\delta}\right\}$$

is contained within the union of the following K events

$$\mathcal{H}_K \subseteq \bigcup_{0 \leq k < K} \mathcal{B}_k,$$

where we define

$$\mathcal{B}_k := \left\{ |Y_n| \geq \sqrt{\frac{4\sigma^2}{2^{k-1}} \log \frac{2K}{\delta}} + \frac{4}{3}R \log \frac{2K}{\delta} \text{ and } \frac{\sigma^2}{2^k} \leq W_n \leq \frac{\sigma^2}{2^{k-1}} \right\}, \quad 1 \leq k \leq K-1,$$

$$\mathcal{B}_0 := \left\{ |Y_n| \geq \sqrt{\frac{4\sigma^2}{2^{K-1}} \log \frac{2K}{\delta}} + \frac{4}{3}R \log \frac{2K}{\delta} \text{ and } W_n \leq \frac{\sigma^2}{2^{K-1}} \right\}.$$

Invoking inequality (67) with σ^2 set to be $\frac{\sigma^2}{2^{k-1}}$ and δ set to be $\frac{\delta}{K}$, we arrive at $\mathbb{P}\{\mathcal{B}_k\} \leq \delta/K$. Taken this fact together with the union bound gives

$$\mathbb{P}\{\mathcal{H}_K\} \leq \sum_{k=0}^{K-1} \mathbb{P}\{\mathcal{B}_k\} \leq \delta.$$

This concludes the proof. □