

### 7.3 Analysis

We will prove the following theorem in this section.

[Theorem 7.1] Regret Bound of UCBVI.

UCBVI achieves the following regret bound:

$$\begin{aligned} \text{Regret} &:= E \left[ \sum_{k=0}^{K-1} (V^*(s_0) - V^{\pi^k}(s_0)) \right] \\ &\leq 20H^2s\sqrt{AK \cdot \ln(SAH^2K^2)} = \hat{O}(H^2s\sqrt{AK}) \end{aligned}$$

Note: the dependency on  $H$  is not tight. By modifying the reward bonus using Bernstein inequality, we could tighten it.

[Lemma 7.2] State-action wise model error.

Fix  $\delta \in (0, 1)$ . For all  $k \in \{0, \dots, K-1\}$ ,  $s \in S$ ,  $a \in A$ ,  $h \in \{0, \dots, H-1\}$ , with prob.  $\geq 1 - \delta$ , we have that for any  $f: S \rightarrow [0, H]$ :

$$|(\hat{P}_h^k(\cdot|s, a) - P_h(\cdot|s, a))'f| \leq 8H\sqrt{S \ln(SAHK/\delta)/N_h^k(s, a)}$$

Pf: ① We consider a fixed tuple  $(s, a, k, h, f)$  first.

Recall the definition of  $\hat{P}_h^k(s, a)$ ,

$$\hat{P}_h^k(\cdot|s, a)'f = \frac{\sum_{i=0}^{k-1} \mathbb{1}\{(S_h^i, a_h^i) = (s, a)\} f(S_{h+1}^i)}{N_h^k(s, a)}$$

Define  $\mathcal{H}_{h,i}$  as the history from iter 0 to iter  $i$  including time step  $h$ .

$$\text{Define } X_i = \mathbb{1}\{(S_h^i, a_h^i) = (s, a)\} [f(S_{h+1}^i) - E_{S' \sim P_h(s, a)}[f(S')]],$$

We now show  $X_i$  is a martingale difference sequence.

$$\begin{aligned} (i) \quad E[X_i | \mathcal{H}_{h,i}] &= E[f(S_{h+1}^i) - E_{S' \sim P_h(s, a)}[f(S')] | \mathcal{H}_{h,i}] \text{ or } 0 \\ &= 0 \end{aligned}$$

$\therefore$   $X_i$  is a martingale difference sequence determined by  $\mathcal{H}_{h,i}$ .

Since  $\mathbb{1}_{\{f(S_h, a_h) = (s, a)\}}$  is determined by  $S_{h+1}$ .

(ii) We have  $|\bar{X}_i| = 0$  for  $(S_h^i, a_h^i) \neq (s, a)$

$$|\bar{X}_i| \leq H \text{ for } (S_h^i, a_h^i) = (s, a).$$

Then  $\bar{X}_i$  is a martingale difference sequence. By Azuma-Hoeffding's inequality: with prob.  $\geq 1 - \delta$ , we have:

$$\begin{aligned} \left| \sum_{i=0}^{k-1} \bar{X}_i \right| &= \left| \sum_{i=0}^{k-1} \mathbb{1}_{\{f(S_h^i, a_h^i) = (s, a)\}} f(S_{h+1}^i) - N_h^s(s, a) E_{S' \sim p_h(s, a)}[f(S')] \right| \\ &\leq 2H \sqrt{N_h^k(s, a) \ln(1/\delta)} \end{aligned}$$

Apply union bound over  $s \in S, a \in A, h \in \{0, \dots, H-1\}, k \in \{0, \dots, K-1\}$ , with prob  $\geq 1 - \delta$ ,

$$\begin{aligned} \left| \sum_{i=0}^{k-1} \mathbb{1}_{\{f(S_h^i, a_h^i) = (s, a)\}} f(S_{h+1}^i) - N_h^s(s, a) E_{S' \sim p_h(s, a)}[f(S')] \right| \\ \leq 2H \sqrt{N_h^k(s, a) \ln(SAKH/\delta)} \end{aligned}$$

② next to cover all  $f: S \rightarrow [0, H]$ .

Note  $\|f\|_2 \leq H\sqrt{S}$  for all  $f$ , there exist a  $\varepsilon$ -net  $N_\varepsilon$  with  $|N_\varepsilon| \leq (1 + 2H\sqrt{S}/\varepsilon)^S$  s.t. for any  $f \in [0, H]^S$ ,  $\exists f' \in N_\varepsilon$  s.t.  $\|f - f'\|_2 \leq \varepsilon$ . [This is obvious since  $[0, H]^S$  is obvious]

Then we have

$$\begin{aligned} &\left| \sum_{i=0}^{k-1} \mathbb{1}_{\{f(S_h^i, a_h^i) = (s, a)\}} f(S_{h+1}^i) / N_h^k(s, a) - E_{S' \sim p_h(s, a)}[f(S')] \right| \\ &\leq \left| \sum_{i=0}^{k-1} \mathbb{1}_{\{f(S_h^i, a_h^i) = (s, a)\}} f'(S_{h+1}^i) / N_h^k(s, a) - E_{S' \sim p_h(s, a)}[f'(S')] \right| \\ &\quad + \left| \sum_{i=0}^{k-1} \mathbb{1}_{\{f(S_h^i, a_h^i) = (s, a)\}} (f(S_{h+1}^i) - f'(S_{h+1}^i)) / N_h^k(s, a) \right| \\ &\quad + \left| E_{S' \sim p_h(s, a)}(f(S') - f'(S')) \right| \\ &\leq 2H \sqrt{S \ln(SAKH(1 + 2H\sqrt{S}/\varepsilon^2)/\delta) / N_h^k(s, a)} + 2\varepsilon \end{aligned}$$

Since  $\varepsilon^2$ -Net,  $|f(s) - f'(s)|^2 \leq \varepsilon^2 \Rightarrow |f(s) - f'(s)| \leq \varepsilon$ .

Now we set  $\varepsilon^2 = 1/K$  and use the fact that  $N_h^k(s, a) \leq K$

we have:

1.  $K-1$

$$\begin{aligned}
& \left| \sum_{i=0}^k \mathbb{1}\{(S_h^i, a_h^i) = (s, a)\} f(S_{h+1}^i) / N_h^k(s, a) - \mathbb{E}_{S' \sim P_h(s, a)} [f(S')] \right| \\
& \leq 2H \sqrt{\ln(\text{SAKH}(1+2HK\sqrt{S}/\delta)) / N_h^k(s, a)} + 2/\sqrt{K} \\
& \leq 4H \sqrt{\ln(4H^2 S^2 k^2 A / \delta) / N_h^k(s, a)} \\
& \leq 8H \sqrt{\ln(HSKA/\delta) / N_h^k(s, a)}
\end{aligned}$$

which completes the proof of Lemma 7.2  $\square$

[Lemma 7.3] State-action wise average model error under  $V^*$ .

Fix  $\delta \in (0, 1)$ . For all  $k \in \{0, \dots, K-1\}$ , with prob.  $\geq 1-\delta$ :

$$| \hat{P}_h^k(\cdot | s, a) \cdot V_{h+1}^* - P_h(\cdot | s, a) \cdot V_{h+1}^* | \leq 2H \sqrt{\ln(\text{SAHN}/\delta) / N_h^k(s, a)}$$

Pf: Although we can bound the LHS by Lemma 7.2, since  $V^*$  is independent with data collected during learning, we can get a tighter upper bound.

① Consider a fixed tuple  $s, a, k, h$  first

$$\hat{P}_h^k(\cdot | s, a) V_{h+1}^* = \frac{1}{N_h^k(s, a)} \sum_{i=0}^{k-1} \mathbb{1}\{(S_h^i, a_h^i) = (s, a)\} V_{h+1}^*(S_{h+1}^i)$$

We define

$$\mathbb{X}_i = \mathbb{1}\{(S_h^i, a_h^i) = (s, a)\} V_{h+1}^*(S_{h+1}^i) - \mathbb{E}[\mathbb{1}\{(S_h^i, a_h^i) = (s, a)\} V_{h+1}^*(S_{h+1}^i) | \mathcal{H}_{h,i}]$$

for  $i=0, \dots, k-1$ .

We have  $\mathbb{E}(\mathbb{X}_i | \mathcal{H}_{h,i}) = 0$ ,  $|\mathbb{X}_i| \leq H$ . Using Azuma-Hoeffding inequality, with prob.  $\geq 1-\delta$

$$\begin{aligned}
\left| \sum_{i=0}^{k-1} \mathbb{X}_i \right| &= \left| \sum_{i=0}^{k-1} \mathbb{1}\{(S_h^i, a_h^i) = (s, a)\} V_{h+1}^*(S_{h+1}^i) - N_h^k(s, a) \mathbb{E}_{S' \sim P_h(s, a)} V_{h+1}^*(S') \right| \\
&\leq 2H \sqrt{N_h^k(s, a) \ln(1/\delta)}
\end{aligned}$$

Divide  $N_h^k(s, a)$  on both side and use the fact  $P_h^* V = \mathbb{E}_{S' \sim P_h(s, a)} (V_{h+1}^*)$  we have:

$$| \hat{P}_h^k(\cdot | s, a)' V_{h+1}^* - P_h(\cdot | s, a)' V_{h+1}^* | \leq 2H \sqrt{\ln(1/\delta) / N_h^k(s, a)}$$

with union bound over  $S, A, [N], [H]$ , we conclude the pf  $\square$

Now we condition on  $\mathcal{E}_{\text{model}}$  (Lemma 7.2 & 7.3) being true.

We study the effect of reward bonus: we want  $\pi^*$  to be optimal under  $r_h + b_h^k$  and empirical  $\hat{P}_h^k$ , i.e. we want  $\hat{V}_0^k(s_0) \geq V_0^*(s_0)$  for all  $s_0$ .

[Lemma 7.4] Optimism.

Assume  $\mathcal{E}_{\text{model}}$  is true. For all episode  $k$  we have

$$\hat{V}_0^k(s_0) \geq V_0^*(s_0), \quad \forall s_0 \in S$$

where  $\hat{V}_0^k(s_0)$  follows VI.

Pf: By induction.

(i) At time step  $H$ :  $\hat{V}_H^k(s) = V_H^*(s) = 0$  for all  $s$

(ii) Starting at  $k+1$ , assuming that  $\hat{V}_{h+1}^k(s) \geq V_{h+1}^*(s)$  for  $\forall s$

Consider any  $s, a \in S \times A$ . First, if  $Q_h^k(s, a) = H$ , then

$$Q_h^k(s, a) \geq Q_h^*(s, a)$$

$$\begin{aligned} \hat{Q}_h^k(s, a) - Q_h^*(s, a) &= b_h^k(s, a) + \hat{P}_h^k(\cdot | s, a) \cdot \hat{V}_{h+1}^k - P_h^*(\cdot | s, a) V_{h+1}^* \\ &\geq b_h^k(s, a) + (\hat{P}_h^k(\cdot | s, a) - P_h^*(\cdot | s, a)) V_{h+1}^* \\ &\geq b_h^k(s, a) - 2H \sqrt{\frac{\ln(SAHK/8)}{N_h^k(s, a)}} \quad (\text{Lemma 7.3}) \\ &\geq 0 \end{aligned}$$

Then we have  $V_h^k(s) \geq V_h^*(s) \quad \forall s$ . □

[Lemma 7.5] Consider  $\mathcal{T}^k = \{s_h^k, a_h^k\}_{h=0}^{H-1}$  for  $k=0, \dots, K-1$ :

$$\sum_{k=0}^{K-1} \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^k(s_h^k, a_h^k)}} \leq 2H\sqrt{SAK}$$

Pf: LHS = 
$$\sum_{h=0}^{H-1} \sum_{k=0}^{K-1} \frac{1}{\sqrt{N_h^k(s_h^k, a_h^k)}}$$

$$\begin{aligned}
&= \sum_{h=0}^{H-1} \sum_{(s,a)} \frac{1}{i} \sqrt{i} \\
&\leq 2 \sum_{h=0}^{H-1} \sum_{(s,a)} \sqrt{N_h^K(s,a)} \quad \left( \frac{1}{\sqrt{i}} = \frac{2}{2\sqrt{i}} \leq \frac{2}{\sqrt{i} + \sqrt{i-1}} = 2(\sqrt{i} - \sqrt{i-1}) \right) \\
&\leq \sum_{h=0}^{H-1} \sqrt{SA \sum_{s,a} N_h^K(s,a)} \quad \text{Cauchy ineq.} \\
&= H \sqrt{SAK}
\end{aligned}$$

□

Now we turn back to prove Th. 7.1 :

$$\text{Regret} \leq 10H^2S\sqrt{AK \ln(SAH^2K)} = \tilde{O}(H^2S\sqrt{AK})$$

Pf [Th. 7.1]:

Consider episode  $k$  and history  $\mathcal{H}_{<k}$ .

① We bound  $V^* - V^{\pi^k}$ .

$$\begin{aligned}
&V_0^*(s_0) - V_0^{\pi^k}(s_0) \\
&= \hat{V}_0^k(s_0) - V_0^{\pi^k}(s_0) \\
&= \hat{Q}_0^k(s_0, \pi^k(s_0)) - Q_0^{\pi^k}(s_0, \pi^k(s_0)) \\
&\leq r_0(s_0, \pi^k(s_0)) + b_h^k(s_0, \pi^k(s_0)) + \hat{P}_0^k(\cdot | s_0, \pi^k(s_0)) \cdot \hat{V}_1^k \quad \leq \text{comes from } Q_0^k \text{, which is min of } H, \sim 1. \\
&\quad - (r_0(s_0, \pi^k(s_0)) + P_0(\cdot | s_0, \pi^k(s_0)) \cdot V_1^{\pi^k}) \\
&= b_h^k(s_0, \pi^k(s_0)) + (\hat{P}_0^k(\cdot | s_0, \pi^k(s_0)) - P_0(\cdot | s_0, \pi^k(s_0))) \cdot \hat{V}_1^k \\
&\quad + P_0(\cdot | s_0, \pi^k(s_0)) \cdot (\hat{V}_1^k - V_1^{\pi^k}) \\
&= \sum_{h=0}^{H-1} E_{s,a \sim d_h^{\pi^k}} [ b_h^k(s,a) + (\hat{P}_h^k(\cdot | s,a) - P_h(\cdot | s,a)) \cdot \hat{V}_{h+1}^k ] \quad (*)
\end{aligned}$$

With Lemma 7.2:

$$\begin{aligned}
\text{RHS of } (*) &\leq \sum_{h=0}^{H-1} E_{s,a \sim d_h^{\pi^k}} [ b_h^k(s,a) + 8H \sqrt{S \ln(SAHK/\delta) / N_h^K(s,a)} ] \\
&\leq \sum_{h=0}^{H-1} E_{s,a \sim d_h^{\pi^k}} [ 10H \sqrt{SL / N_h^K(s,a)} ] \\
&= 10H \sqrt{SL} E \left[ \sum_{h=0}^{H-1} 1 / \sqrt{N_h^K(s,a)} \mid \mathcal{H}_{<k} \right]
\end{aligned}$$

where  $L = \ln(SAHK/\delta)$

② By summing all episodes together:

$$E \left[ \sum_{k=0}^{K-1} V_0^*(s_0) - V_0^{\pi^k}(s_0) \right]$$

$$= E \left[ \sum_{k=0}^{K-1} V_0^*(s_0) - V_0^{\pi^k}(s_0) \right] = E \left[ \sum_{k=0}^{K-1} V_0^*(s_0) - V_0^{\pi^k}(s_0) \right]$$

$$= E[ \mathbb{E}(\xi_{model} | \bigwedge_{k=0}^{K-1} V_0(s_0) - V_0(s_1)) ] + E[ \mathbb{E}(\xi_{model} | \bigwedge_{k=0}^{K-1} V(s_0) - V(s_1)) ] \\ \leq 10H\sqrt{SL} E[ \sum_{k=0}^{K-1} \sum_{h=0}^{H-1} 1/\sqrt{N_h^k(s_h^k, a_h^k)} ] + 2\delta KH \quad (**)$$

With Lemma 7.5,

$$\text{RHS of } (**) \leq 10H\sqrt{SL} \cdot 2H\sqrt{SAK} + 2\delta KH \\ = 20H^2S\sqrt{AK\ln(SAHK/\delta)} + 2\delta KH$$

set  $\delta = 1/KH$ , we get

$$E[ \sum_{k=0}^{K-1} V^*(s_0) - V^{\pi^k}(s_0) ] \leq 20H^2S\sqrt{AK\ln(SAH^2K^2)} + 2 \\ = O(H^2S\sqrt{AK\ln(SAH^2K^2)})$$

which completes the proof of Th. 7.1.  $\square$