

Susan's version of Auer's UCB

To reduce notational burden!

$$T=2, |\mathcal{A}|=2, |\mathcal{S}|=2 \quad \mathcal{S} = \{0, 1\}$$

M is number of episodes.

From the prior $m-1$ episodes we will form $\pi^m = \begin{pmatrix} \pi_0 \\ \pi_m \end{pmatrix}$ e.g. π^m is a function of

$$\mathcal{H}_{m-1} = \{ S_{i_0} A_{i_0} S_{i_1} A_{i_1} \dots A_{i_{m-1}} S_{i_m} \mid i < m \}$$

UCB Algorithm (here $T=2$)

for $j=1, \dots, M$ episodes.

learner sees $S_{j_0} = s_0$

for $t=0, \dots, T-1$

learner selects A_{j_t} using π^j

learner receives $S_{j_{t+1}}$

$$R_{j_t} = R(S_{j_t}, A_{j_t}, S_{j_{t+1}})$$

end for

form π^{j+1} from j prior episodes
(UCB part)

end for

Prior to $j=m_0$, episodes are run to ensure that we see at least 1 (state, action) pair for each possible $(s, a) \in S \times A$ (maybe select actions uniformly in A) We ignore this in following.

How do we form π_j from \mathcal{H}_{j-1} ?

Some Statistics:

For $t=0, 1$

$$N_{sa}^{t,m-1} = \sum_{j=1}^{m-1} \mathbb{1}\{S_{jt}=s, A_{jt}=a\}$$

$$\text{where } N_{sa}^{m-1} = N_{sa}^{0,m-1} + N_{sa}^{1,m-1}$$

$$P_{sa}^{t,m-1} = \sum_{j=1}^{m-1} \mathbb{1}\{S_{jt}=s, A_{jt}=a\} R_{jt}$$

$$\text{where } P_{sa}^{m-1} = P_{sa}^{0,m-1} + P_{sa}^{1,m-1}$$

$$P_{sas'}^{m-1} = \sum_{j=1}^{m-1} \sum_{t=0}^1 \mathbb{1}\{S_{jt}=s, A_{jt}=a\} \mathbb{1}\{S_{j+t+1}=s'\}$$

For $t=0, 1$

$$\tilde{r}_a^{t,m-1}(s) = \frac{P_{sa}^{t,m-1}}{N_{sa}^{t,m-1}} + \sqrt{\frac{c \log m}{N_{sa}^{t,m-1}}} \quad \begin{cases} \text{UCB} \\ \text{version} \\ \text{of average} \\ \text{reward} \end{cases}$$

$$\tilde{r}_a^{m-1} = \frac{P_{sa}^{m-1}}{N_{sa}^{m-1}} + \sqrt{\frac{c \log m}{N_{sa}^{m-1}}}$$

Form a UCB version of the transition probabilities:

$$\tilde{P}_a^{m-1}(s, s') = \frac{P_{sa}^{m-1}}{N_{sa}^{m-1}} + \sqrt{\frac{C \log m}{N_{sa}^{m-1}}} \quad N_{sa}^{m-1} = N_{sa}^{0, m-1} + N_{sa}^{1, m-1}$$

Note that $\tilde{P}_a^{m-1}(s, s')$ may be > 1 .

(I will drop the superscript $m-1$ on $\tilde{r}_a(s)$, $\tilde{P}_a(s, s')$:

For any given π_i , we fix π_i as follows

$$\begin{aligned} \tilde{P}_a^{\pi_i, m-1}(s, 1) &= \min\left(1, \tilde{P}_a(s, 1)\right) \\ &\quad + \min\left(1 - \min\left(1, \tilde{P}_a(s, 0)\right), \right) \end{aligned}$$

$$\begin{aligned} \tilde{P}_a^{\pi_i, m-1}(s, 0) &= 1 - \tilde{P}_a^{\pi_i, m-1}(s, 1) \\ &= \min\left(1 - \min\left(1, \tilde{P}_a(s, 1)\right), \right) \\ &\quad + \min\left(1 - \min\left(1, \tilde{P}_a(s, 0)\right), \right) \end{aligned}$$

After $m-1$ episodes we form

$$\pi_i^{m-1}(s) = \arg \max_a \tilde{r}_a^i(s) \quad (\text{like VCB})$$

$$= \arg \max_{\pi_i} \tilde{V}_{\pi_i}^{m-1}(s) \quad \text{where } \tilde{V}_{\pi_i}^{m-1}(s) \triangleq \tilde{r}_{\pi_i^{m-1}(s)}^i$$

$$\pi_0^{m-1}(s) = \arg \max_a \left\{ \tilde{r}_a(s) + \sum_{s'} \tilde{P}_a^{\pi_i^{m-1}, m-1}(s, s') \tilde{r}_{\pi_i^{m-1}(s')}^i \right\}$$

$$= \arg \max_{\pi_0} \tilde{V}_{\pi_0}^{2, m-1}(s)$$

$$\text{where } \tilde{V}_{\pi_0}^{2, m-1}(s) = \tilde{r}_{\pi_0(s)}^i + \sum_{s'} \tilde{P}_{\pi_0(s)}^{\pi_i^{m-1}, m-1}(s, s') \tilde{r}_{\pi_i^{m-1}(s')}^i$$

(12)

So in the m^{th} episode

$$A_{m_0} = \pi_0^{m-1}(s_0), \quad A_{m_i} = \pi_i^{m-1}(s_{m_i}).$$

(drop superscript of $m-1$ on $\tilde{V}, \tilde{\pi}^{\hat{\pi}_i}$)

Recall the expected regret is

$$R_M(L, D, \pi) = E \left[\sum_{j=1}^M V_{\pi^*}^2(s_0) - \sum_{j=1}^M V_{\hat{\pi}^j}^2(s_0) \right]$$

$$T=2 \quad |S|=2 \\ |A|=2$$

We prove

$$R_M(L, D, \pi) = O\left(\frac{\log M}{\Delta}\right)$$

where $\Delta = \min_{\pi \neq \pi^*} V_{\pi^*}^2(s_0) - V_{\pi}^2(s_0)$

$$\pi = \begin{pmatrix} \pi_0 \\ \pi_1 \end{pmatrix} \text{ 13}$$

regret

$$= \sum_{m=1}^M V_{\pi^*}^2(s_0) - V_{\pi^{m-1}}^2(s_0)$$

$$= \sum_{m=1}^M \left(V_{\pi^*}^2(s_0) - V_{\pi^{m-1}}^2(s_0) \right)$$

$$* \underbrace{\mathbb{1}\{A_{m0} = \pi_0^{m-1}(s_0), A_{m1} = \pi_1^{m-1}(s_m)\}}$$

$$* \underbrace{\mathbb{1}\{\tilde{V}_{\pi^{m-1}}^{2,m}(s_0) - \tilde{V}_{\pi^*}^{2,m}(s_0) > 0\}}$$

$$\leq \sum_{\pi \neq \pi^*} \Delta_\pi \sum_{m=1}^M \underbrace{\mathbb{1}\{A_{m0} = \pi_0(s_0), A_{m1} = \pi_1(s_m)\}}_{* \underbrace{\mathbb{1}\{\tilde{V}_{\pi_1}^{2,m}(s_0) - \tilde{V}_{\pi^*}^{2,m}(s_0) > 0\}}}$$

$$\text{where } \Delta_\pi = V_{\pi^*}^2(s_0) - V_\pi^2(s_0)$$

$$\leq \sum_{\pi \neq \pi^*} \Delta_\pi \sum_{m=1}^M \underbrace{\mathbb{1}\{A_{m0} = \pi_0(s_0), A_{m1} = \pi_1(s_m)\}}_{* \underbrace{\mathbb{1}\{\tilde{V}_{\pi_1}^{2,m}(s_0) - \tilde{V}_{\pi^*}^{2,m}(s_0) > 0\}}} * \underbrace{\mathbb{1}\{N_{S_0 \pi_0(s_0)}^{m-1} > l, N_{\pi_1}^{1,m-1} > l_{\pi_1}\}}$$

$$+ \sum_{\pi \neq \pi^*} \Delta_\pi \cdot l_\pi$$

since

$$N_{\pi_1}^{1,m-1} = \sum_{j=1}^{m-1} \mathbb{1}\{A_{j1} = \pi_1(s_{j1})\}$$

(14)

Note

$$\begin{aligned}
 N_{\pi_i}^{1, m-1} &= \sum_{j=1}^{m-1} \mathbb{1}\{A_{ji} = \pi_i(1), S_{ji} = 1\} \\
 &\quad + \sum_{j=1}^{m-1} \mathbb{1}\{A_{ji} = \pi_i(0), S_{ji} = 0\} \\
 &= N_{i, \pi_i}^{1, m-1} + N_{0, \pi_i}^{1, m-1}
 \end{aligned}$$

(we will select l_{π_i} later) S_0

regret

$$\leq \sum_{\pi \neq \pi^*} \Delta_{\pi} \sum_{m=1}^M \mathbb{1}\{ \tilde{V}_{\pi}^{2, m-1}(s_0) - \tilde{V}_{\pi^*}^{2, m-1}(s_0) > 0 \} * \mathbb{1}\{ N_{S_0, \pi_i(s_0)}^{m-1} > l_{\pi}, N_{\pi_i}^{m-1} > l_{\pi} \}$$

$$+ \sum_{\pi \neq \pi^*} \Delta_{\pi} l_{\pi}$$

$$\mathbb{1}\{ \tilde{V}_{\pi}^{2, m-1}(s_0) - \tilde{V}_{\pi^*}^{2, m-1}(s_0) \geq 0 \} =$$

$$\left\{ \left(\tilde{V}_{\pi_i}^{2, 1}(s_0) - \tilde{V}_{\pi_i}^2(s_0) \right) - \left(\tilde{V}_{\pi^*}^{2, 1}(s_0) - \tilde{V}_{\pi^*}^2(s_0) \right) \right\} \\
 + \left(\tilde{V}_{\pi_i}^2(s_0) - \tilde{V}_{\pi^*}^2(s_0) \right) \geq 0$$

$$\text{Recall } \tilde{V}_{\pi}^{2,m-1}(s_0) = \tilde{r}_{\pi_0(s_0)} + \sum_{s'} \tilde{p}_{\pi_0(s_0), s'}^{\pi_1} \tilde{r}_{\pi_1(s')}^{m-1}$$

$$V_{\pi}^2(s_0) = r_{\pi_0(s_0)} + \sum_{s'} p_{\pi_0(s_0), s'} \tilde{r}_{\pi_1(s')}^{m-1}$$

(I am often leaving off the superscript)
of $m-1$
as in $\tilde{r}_{\pi_1(s')}^{m-1}$ or $\tilde{p}_{\pi_0(s_0), s'}^{\pi_1, m-1}$)

$$\begin{aligned} s_0 \quad \tilde{V}_{\pi}^{2,m-1}(s_0) - V_{\pi}^2(s_0) &= \tilde{r}_{\pi_0(s_0)} - r_{\pi_0(s_0)} \\ &\quad + \sum_{s'} p_{\pi_0(s_0), s'} (\tilde{r}_{\pi_1(s')}^{m-1} - r_{\pi_1(s')}) \\ &\quad + \sum_s (p_{\pi_0(s_0), s}^{\pi_1} - p_{\pi_0(s_0), s}) \tilde{r}_{\pi_1(s')}^{m-1} \end{aligned}$$

Because we only have 2 states

$$\begin{aligned} &\sum_{s'} (p_{\pi_0(s_0), s'}^{\pi_1} - p_{\pi_0(s_0), s'}) \tilde{r}_{\pi_1(s')}^{m-1} \\ &= (\tilde{p}_{\pi_0(s_0), 1}^{\pi_1} - p_{\pi_0(s_0), 1}) (\tilde{r}_{\pi_1(1)}^{m-1} - \tilde{r}_{\pi_1(0)}^{m-1}) \\ &= (\min(1, \tilde{p}_{\pi_0(s_0), 1}) - p_{\pi_0(s_0), 1}) (\tilde{r}_{\pi_1(1)}^{m-1} - \tilde{r}_{\pi_1(0)}^{m-1})^+ \\ &\quad + (\min(1, \tilde{p}_{\pi_0(s_0), 0}) - p_{\pi_0(s_0), 0}) (\tilde{r}_{\pi_1(0)}^{m-1} - \tilde{r}_{\pi_1(1)}^{m-1})^+ \end{aligned}$$

(Recall that $F_a^2(s) = \frac{F_a^1(s)}{N_{s,a}^{m-1}} + \sqrt{\frac{c \log m}{N_{s,a}^{m-1}}}$)

On the set $\tilde{p}_{\pi_0(s_0), 1} \leq 1 \wedge \tilde{p}_{\pi_0(s_0), 0} \leq 1$

we have

$$\begin{aligned} \tilde{V}_{\pi}^{2,m-1}(s_0) - V_{\pi}^2(s_0) &= \frac{p_{s_0, \pi_0(s_0)}^{m-1}}{N_{s_0, \pi_0}} - r_{\pi_0(s_0)} - \sqrt{\frac{c \log m}{N_{s_0, \pi_0}}} \\ &\quad + \sum_s p_{\pi_0(s_0), s} \left(\frac{p_{s, \pi_1(s)}^{m-1}}{N_{s, \pi_1(s)}} - r_{\pi_1(s)} \right) \\ &\quad + \left(\frac{\tilde{p}_{s_0, \pi_0(s_0), 1}}{N_{s_0, \pi_0}} - p_{\pi_0(s_0), 1} - \sqrt{\frac{c \log m}{N_{s_0, \pi_0}}} \right)^+ \\ &\quad \quad \quad \times (\tilde{r}_{\pi_1(1)}^{m-1} - \tilde{r}_{\pi_1(0)}^{m-1}) \\ &\quad + \left(\frac{\tilde{p}_{s_0, \pi_0(s_0), 0}}{N_{s_0, \pi_0}} - p_{\pi_0(s_0), 0} - \sqrt{\frac{c \log m}{N_{s_0, \pi_0}}} \right)^+ \\ &\quad \quad \quad \times (\tilde{r}_{\pi_1(0)}^{m-1} - \tilde{r}_{\pi_1(1)}^{m-1}) \\ &\quad + 2 \sqrt{\frac{c \log m}{N_{s_0, \pi_0}}} + 2 \sqrt{\frac{c \log m}{N_{s_0, \pi_0}}} |\tilde{r}_{\pi_1(0)}^{m-1} - \tilde{r}_{\pi_1(1)}^{m-1}| \end{aligned}$$

$$\leq \frac{P_{S_0, \pi_0(s_0)}^{m-1}}{N_{S_0 \pi_0}} - r_{\pi_0(s_0)} - \sqrt{\frac{c \log m}{N_{S_0 \pi_0}}} + \\ + \sum_S P_{\pi_0(s_0), S} \left(\frac{P_{S_0, \pi_1(S)}^{m-1}}{N_{S_0 \pi_1(S)}} - r_{\pi_1(S)} - \sqrt{\frac{c \log m}{N_{S_0 \pi_1(S)}}} \right) +$$

$$+ \left(\frac{\tilde{P}_{S_0 \pi_0(s_0), 1}}{N_{S_0 \pi_0}} - \tilde{P}_{\pi_0(s_0), 1} - \sqrt{\frac{c \log m}{N_{S_0 \pi_0}}} \right) + \\ * \left(\tilde{r}_{\pi_0(1)} - \tilde{r}_{\pi_0(0)} \right)$$

$$+ \left(\frac{\tilde{P}_{S_0 \pi_0(s_0), 0}}{N_{S_0 \pi_0}} - \tilde{P}_{\pi_0(s_0), 0} - \sqrt{\frac{c \log m}{N_{S_0 \pi_0}}} \right) + \\ * \left(\tilde{r}_{\pi_0(0)} - \tilde{r}_{\pi_0(1)} \right)$$

$$+ 2 \sqrt{\frac{c \log m}{N_{S_0 \pi_0}}} + 2 \sqrt{\frac{c \log m}{N_{S_0 \pi_0}}} \left| \tilde{r}_{\pi_1(0)} - \tilde{r}_{\pi_1(1)} \right|$$

$$+ \sum_S P_{\pi_0(s_0), S} \sqrt{\frac{c \log m}{N_{S \pi_1}}} \left\{ \frac{P_{S, \pi_1(S)}^{m-1}}{N_{S \pi_1}} - r_{\pi_1(S)} \right\} > \sqrt{\frac{c \log m}{N_{S \pi_1}}}$$

In the following we ignore that fact that the above inequality was only shown to hold under

$\tilde{P}_{\pi_0(s_0), 1} \leq 1 \wedge \tilde{P}_{\pi_0(s_0), 0} \leq 1$ and we will assume that

$|\tilde{r}_a^1(s)| \leq 1$ as well. This can be fixed albeit with more technical, drawn out derivations!

Recall

$$1 \left\{ \tilde{V}_{\pi}^{z_1, m-1}(s_0) - \tilde{V}_{\pi^*}^{z_1, m-1}(s_0) \geq 0 \right\} = \\ \left\{ \left(\tilde{V}_{\pi}^{z_1}(s_0) - \tilde{V}_{\pi^*}^{z_1}(s_0) \right) - \left(\tilde{V}_{\pi^*}^{z_1}(s_0) - \tilde{V}_{\pi^*}^{z_1}(s_0) \right) \right\} \\ + \left(\tilde{V}_{\pi}^{z_1}(s_0) - \tilde{V}_{\pi^*}^{z_1}(s_0) \right) > 0$$

Using $\mathbb{1}\{A+B+C > 0\} \leq \mathbb{1}\{A > 0\} + \mathbb{1}\{B > 0\} + \mathbb{1}\{C > 0\}$: (17)

$$\mathbb{1}\left\{\tilde{V}_{\pi^*}^{2_{i^{m-1}}}(s_0) - \tilde{V}_{\pi^*}^{2_{i^{m-1}}}(s_0) > 0\right\}$$

$$\begin{aligned} &\leq \mathbb{1} \left\{ \frac{\tilde{P}_{s_0, \pi_0(s_0)}^{m-1}}{N_{s_0, \pi_0}} - r_{\pi_0(s_0)} - \sqrt{\frac{c \log m}{N_{s_0, \pi_0}}} + \right. \\ &\quad + \sum_s P_{\pi_0(s_0), s} \left(\frac{\tilde{P}_{s_0, \pi_1(s)}^{m-1}}{N_{s_0, \pi_1(s)}} - r_{\pi_1(s)} - \sqrt{\frac{c \log m}{N_{s_0, \pi_1}}} \right)^+ \\ &\quad + \left(\frac{\tilde{P}_{s_0, \pi_0(s_0), 1}}{N_{s_0, \pi_0}} - \tilde{P}_{\pi_0(s_0), 1} - \sqrt{\frac{c \log m}{N_{s_0, \pi_0}}} \right)^+ * \left(\tilde{r}_{\pi_1(1)} - \tilde{r}_{\pi_1(0)} \right)^+ \\ &\quad + \left(\frac{\tilde{P}_{s_0, \pi_0(s_0), 0}}{N_{s_0, \pi_0}} - \tilde{P}_{\pi_0(s_0), 0} - \sqrt{\frac{c \log m}{N_{s_0, \pi_0}}} \right)^+ * \left(\tilde{r}_{\pi_1(0)} - \tilde{r}_{\pi_1(1)} \right)^+ \geq 0 \left. \right\} \end{aligned}$$

$$+ \mathbb{1}\left\{-\left(\tilde{V}_{\pi^*}^{2_1}(s_0) - \tilde{V}_{\pi^*}^{2_1}(s_0)\right) > 0\right\}$$

$$\begin{aligned} &+ \mathbb{1} \left(\tilde{V}_{\pi}^2(s_0) - \tilde{V}_{\pi^*}^2(s_0) \right. \\ &\quad \left. + 2 \sqrt{\frac{c \log m}{N_{s_0, \pi_0}}} + 2 \sqrt{\frac{c \log m}{N_{s_0, \pi_0}}} \left| \tilde{r}_{\pi_1(0)} - \tilde{r}_{\pi_1(1)} \right| \right) \end{aligned}$$

$$\sum_s P_{\pi_0(s_0, s)} \sqrt{\frac{c \log m}{N_{s\pi_1}}} \left[\left\{ \frac{P_{s, \pi_1(s)}^{m-1}}{N_{s\pi_1}} - r_{\pi_1}(s) \right\} > \sqrt{\frac{c \log m}{N_{s\pi_1}}} \right] > 0$$

Recall the regret is

$$\leq \sum_{\pi \neq \pi^*} \sum_{m=1}^M \left[\begin{array}{l} 1 \{ V_{\pi_1}(s_0) - \tilde{V}_{\pi^*}^{2^{m-1}}(s_0) > 0 \} \\ * 1 \{ N_{s_0\pi_0(s_0)}^{m-1} > l_{\pi_1}, N_{s\pi_1}^{m-1} > l_{\pi_1} \} \end{array} \right] + \sum_{\pi \neq \pi^*} \Delta_{\pi} l_{\pi}$$

Consider the contribution of the *** term**, pg 17 to the regret.

$$\leq 1 \left\{ 0 \leq -\Delta_{\pi} + 6 \sqrt{\frac{c \log m}{N_{s\pi_0}}} + \sum_s P_{\pi_0(s_0, s)} \sqrt{\frac{c \log m}{N_{s\pi_1}}} \left[\left\{ \frac{P_{s\pi_1}^{m-1}}{N_{s\pi_1}} - r_{\pi_1}(s) \right\} > \sqrt{\frac{c \log m}{N_{s\pi_1}}} \right] * 1 \{ N_{s_0\pi_0}^{m-1} > l_{\pi_1}, N_{s\pi_1}^{m-1} + N_{s\pi_0}^{m-1} > l_{\pi_1} \} \right\}$$

note $N_{s\pi_1}^{m-1} + N_{s\pi_0}^{m-1} \geq l_{\pi_1} \Rightarrow N_{s\pi_1}^{m-1} > \frac{1}{2} l_{\pi_1} \cup N_{s\pi_0}^{m-1} > \frac{1}{2} l_{\pi_1}$.

$$\leq 1 \left\{ 0 \leq -\Delta_{\pi} + 6 \sqrt{\frac{c \log m}{l_{\pi}}} + \sum_s P_{\pi_0(s_0, s)} \sqrt{\frac{c \log m}{N_{s\pi_1}}} \left[\left\{ \frac{P_{s\pi_1}^{m-1}}{N_{s\pi_1}} - r_{\pi_1}(s) \right\} > \sqrt{\frac{c \log m}{N_{s\pi_1}}} \right] * 1 \{ N_{s\pi_1}^{m-1} > \frac{1}{2} l_{\pi_1} \} \right\}$$

$$* 1 \{ N_{s\pi_1}^{m-1} > \frac{1}{2} l_{\pi_1} \}$$

(19)

$$+ 1 \left\{ \begin{array}{l} 0 \leq -\Delta_{\pi} + 6 \sqrt{\frac{c \log m}{l}} \\ + \sum_s p_{\pi_0}(s_0, s) \sqrt{\frac{c \log m}{N_{s\pi_1}}} \mathbb{1} \left\{ \frac{p_{s\pi_1}^{m-1}}{N_{s\pi_1}} - r_{\pi_1}(s) > \sqrt{\frac{c \log m}{N_{s\pi_1}}} \right\} \end{array} \right\}$$

$$\neq 1 \left\{ N_{0\pi_1}^{m-1} \geq \frac{1}{2} l_{\pi_1} \right\}$$

(consideration of the two terms is
the same - so consider only the 1st
term)

$$1 \left\{ \begin{array}{l} 0 \leq -\Delta_{\pi} + 6 \sqrt{\frac{c \log m}{l_{\pi}}} \\ + \sum_s p_{\pi_0}(s_0, s) \sqrt{\frac{c \log m}{N_{s\pi_1}}} \mathbb{1} \left\{ \frac{p_{s\pi_1}^{m-1}}{N_{s\pi_1}} - r_{\pi_1}(s) > \sqrt{\frac{c \log m}{N_{s\pi_1}}} \right\} \end{array} \right\}$$

$$\neq 1 \left\{ N_{1\pi_1}^{m-1} \geq \frac{1}{2} l_{\pi_1} \right\}$$

$$\leq 1 \left\{ \begin{array}{l} 0 \leq -\Delta_{\pi} + 6 \sqrt{\frac{c \log m}{l_{\pi}}} \\ + p_{\pi_0}(s_0, 1) \sqrt{\frac{c \log m}{\gamma_2 l_{\pi}}} \\ - + p_{\pi_0}(s_0, 0) \sqrt{\frac{c \log m}{N_{0\pi_1}}} \mathbb{1} \left\{ \frac{p_{0\pi_1}^{m-1}}{N_{0\pi_1}} - r_{\pi_1}(0) > \sqrt{\frac{c \log m}{N_{0\pi_1}}} \right\} \end{array} \right\}$$

(20)

$$+ 1 \left\{ N_{\pi_1}^{m-1} \geq \frac{1}{2} l_{\pi_1} \right\}$$

$$= 1 \left\{ 0 \leq -\Delta_{\pi_1} + 6 \sqrt{\frac{c \log m}{l_{\pi_1}}} + P_{\pi_0}(s_0, 1) \sqrt{\frac{c \log m}{l_{\pi_1}}} \right\}$$

$$+ 1 \left\{ \frac{P_{\pi_1}}{N_{\pi_1}} - r_{\pi_1}(0) > \sqrt{\frac{c \log m}{N_{\pi_1}}} \right\}$$

we will choose l_{π_1} so that 1st indicator
is zero $\forall m!$ e.g. $l_{\pi_1} = \frac{8^2 c \log M}{\Delta_{\pi_1}^2}$

So the regret is bounded above
by

$$\leq \sum_{\pi \neq \pi^*} \Delta_{\pi} \sum_{m=1}^M 1 \left\{ \tilde{V}_{\pi_1}^{2^{m-1}}(s_0) - \tilde{V}_{\pi^*}^{2^{m-1}}(s_0) > 0 \right\} * 1 \left\{ N_{s_0, \pi_0}^{m-1} > l_{\pi_1}, N_{\pi_1}^{m-1} > l_{\pi_1} \right\}$$

$$+ \sum_{\pi \neq \pi^*} \Delta_{\pi} l_{\pi_1}$$

$$\leq \sum_{\pi \neq \pi^*} \Delta_{\pi} l_{\pi_1} +$$

$$\sum_{\pi \neq \pi^*} \Delta_{\pi} \cdot \sum_{m=1}^M 1 \left\{ \frac{P_{s_0, \pi_0}^{m-1}}{N_{s_0, \pi_0}} - r_{\pi_0}(s_0) - \sqrt{\frac{c \log m}{N_{s_0, \pi_0}}} > 0 \right\}$$

$$+ 1 \left\{ \frac{P_{s_0, \pi_1}^{m-1}}{N_{s_0, \pi_1}(s_0)} - r_{\pi_1}(s_0) - \sqrt{\frac{c \log m}{N_{s_0, \pi_1}}} > 0 \right\}$$

21

$$+ \mathbb{1} \left\{ \frac{\tilde{P}_{S_0 \pi_0(S_0), 1}}{N_{S_0 \pi_0}} - P_{\pi_0(S_0), 1} - \sqrt{\frac{c \log m}{N_{S_0 \pi_0}}} > 0 \right\}$$

$$+ \mathbb{1} \left\{ \frac{\tilde{P}_{S_0 \pi_0(S_0), 0}}{N_{S_0 \pi_0}} - P_{\pi_0(S_0), 0} - \sqrt{\frac{c \log m}{n N_{S_0 \pi_0}}} > 0 \right\}$$

$$+ \mathbb{1} \left\{ - (\tilde{V}_{\pi^*}^2(S_0) - V_{\pi^*}^2(S_0)) \geq 0 \right\}$$

$$+ \mathbb{1} \left\{ \frac{P_{0 \pi_1}^{m-1}}{N_{0 \pi_1}} - r_{\pi_1}(0) - \sqrt{\frac{c \log m}{N_{0 \pi_1}}} > 0 \right\}$$

$$+ \mathbb{1} \left\{ \frac{P_{1 \pi_1}^{m-1}}{N_{1 \pi_1}} - r_{\pi_1}(1) - \sqrt{\frac{c \log m}{N_{1 \pi_1}}} > 0 \right\}$$

The expectation of each term is treated the same: for example, consider

$$P \left\{ \frac{P_{S_0, \pi_0(S_0), S}^{m-1}}{N_{S_0 \pi_0}} - P_{\pi_0(S_0), S} - \sqrt{\frac{c \log m}{N_{S_0 \pi_0}}} > 0 \right\}$$

$$\leq P \left[\exists n | s_n \leq m-1 \text{ for which } \frac{P_{S_0, \pi_0(S_0), S}^{n-1}}{n} - P_{\pi_0(S_0), S} - \sqrt{\frac{c \log m}{n}} > 0 \right]$$

 \leq

(22)

$$\sum_{n=1}^m P \left[\frac{P_{S_0, \pi_0(S_0), S}^{n-1}}{n} - P_{\pi_0(S_0)}(S_0, S) - \sqrt{\frac{c \log m}{n}} > 0 \right]$$

$$\leq \sum_{n=1}^m P \left[\frac{1}{n} \sum_{j=1}^n \mathbb{1}_{\{A_{j0} = \pi_0(S_0), S_{j0} = S_0\}} \left(\mathbb{1}_{\{S_j = S\}} - P_{\pi_0(S_0)}(S_0, S) \right) \right] > \sqrt{\frac{c \log m}{n}}$$

now apply Hoeffding-Azuma (non-independent terms)

note that

$$E \left[\mathbb{1}_{\{A_{j0} = \pi_0(S_0), S_{j0} = S_0\}} \left(\mathbb{1}_{\{S_j = S\}} - P_{\pi_0(S_0)}(S_0, S) \right) \right] = 0$$

to see that

$$\begin{aligned} &\leq \exp \left\{ -2 \cdot n \cdot \frac{c \log m}{n} \right\} = \exp \{-2c \log m\} \\ &= m^{-2c} \end{aligned}$$

$$\sum_{\pi \neq \pi^*} \Delta_\pi \sum_{m=1}^M \sum_{n=1}^m m^{-2c}$$

$$= \sum_{\pi \neq \pi^*} \Delta_\pi \sum_{m=1}^M m^{-2(c-1)}$$

sums finitely for $c \geq 1$.

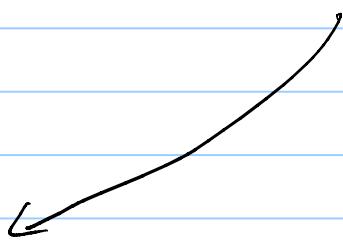
(23)

Thus the expected regret is bounded above by

$$\Delta_\pi = V_{\pi^*}^2(s_0) - V_\pi^2(s_0)$$

$$E \left[\sum_{m=1}^M V_{\pi^*}^2(s_0) - V_{\pi^{m-1}}^2(s_0) \right]$$

$$\leq \sum_{\pi \neq \pi^*} \Delta_\pi \frac{8^2 c \log M}{\Delta_\pi^2}$$



$$+ \text{some finite } * \sum_{\text{integer}} \sum_{\pi \neq \pi^*} \Delta_\pi \sum_{m=1}^M m^{-2(c-1)}$$

$$= O \left(\frac{\log M}{\min_{\pi \neq \pi^*} \Delta_\pi} |\pi|^{1/\zeta} \right)$$