

Exploration in Reinforcement Learning (theory)

Lecturers: *M. Pirotta*

(December 16, 2021)

Solution by **FILL fullname** command at the beginning of latex document

Instructions

- The deadline is **January 16, 2022. 23h59**
- By doing this homework you agree to the *late day policy, collaboration and misconduct rules* reported on Piazza.
- **Mysterious or unsupported answers will not receive full credit.** A correct answer, unsupported by calculations, explanation, or algebraic work will receive no credit; an incorrect answer supported by substantially correct calculations and explanations might still receive partial credit.
- Answers should be provided in **English**.

1 Best Arm Identification

In best arm identification (BAI), the goal is to identify the best arm in as few samples as possible. We will focus on the fixed-confidence setting where the goal is to identify the best arm with high probability $1 - \delta$ in as few samples as possible. A player is given k arms with expected reward μ_i . At each timestep t , the player selects an arm to pull (I_t), and they observe some reward ($X_{I_t, t}$) for that sample. At any timestep, once the player is confident that they have identified the best arm, they may decide to stop.

δ -correctness and fixed-confidence objective. Denote by τ_δ the stopping time associated to the stopping rule, by i^* the best arm and by \hat{i} an estimate of the best arm. An algorithm is δ -correct if it predicts the correct answer with probability at least $1 - \delta$. Formally, if $\mathbb{P}_{\mu_1, \dots, \mu_k}(\hat{i} \neq i^*) \leq \delta$ and $\tau_\delta < \infty$ almost surely for any μ_1, \dots, μ_k . Our goal is to find a δ -correct algorithm that minimizes the sample complexity, that is, $\mathbb{E}[\tau_\delta]$ the expected number of sample needed to predict an answer. Assume that the best arm i^* is *unique* (i.e., there exists only one arm with maximum mean reward).

Notation

- I_t : the arm chosen at round t .
- $X_{i,t} \in [0, 1]$: reward observed for arm i at round t .
- μ_i : the expected reward of arm i .
- $\mu^* = \max_i \mu_i$.
- $\Delta_i = \mu^* - \mu_i$: suboptimality gap.

Consider the following algorithm

The algorithm maintains an active set S and an estimate of the empirical reward of each arm $\hat{\mu}_{i,t} = \frac{1}{t} \sum_{j=1}^t X_{i,j}$.

- Compute the function $U(t, \delta)$ that satisfy the any-time confidence bound. Let

$$\mathcal{E} = \bigcup_{i=1}^k \bigcup_{t=1}^{\infty} \{ |\hat{\mu}_{i,t} - \mu_i| > U(t, \delta') \}.$$

```

Input:  $k$  arms, confidence  $\delta$ 
 $S = \{1, \dots, k\}$ 
for  $t = 1, \dots$  do
    Pull all arms in  $S$ 
     $S = S \setminus \left\{ i \in S : \exists j \in S, \hat{\mu}_{j,t} - U(t, \delta') \geq \hat{\mu}_{i,t} + U(t, \delta') \right\}$ 
    if  $|S| = 1$  then
        STOP
        return  $S$ 
    end
end

```

Using Hoeffding's inequality and union bounds, shows that $\mathbb{P}(\mathcal{E}) \leq \delta$ for a particular choice of δ' . This is called "bad event" since it means that the confidence intervals do not hold.

- Show that with probability at least $1 - \delta$, the optimal arm $i^* = \arg \max_i \{\mu_i\}$ remains in the active set S . Use your definition of δ' and start from the condition for arm elimination. From this, use the definition of $\neg\mathcal{E}$.
- Under event $\neg\mathcal{E}$, show that an arm $i \neq i^*$ will be removed from the active set when $\Delta_i \geq C_1 U(t, \delta')$ for some constant $C_1 \in \mathbb{N}$. Compute the time required to have such condition for each non-optimal arm. Use the condition of arm elimination applied to arm i^* .¹
- Compute a bound on the sample complexity (after how many *pulls* the algorithm stops) for identifying the optimal arm w.p. $1 - \delta$.
- We assumed that the optimal arm i^* is unique. Would the algorithm still work if there exist multiple best arms? Why?

Note that also a variations of UCB are effective in pure exploration.

2 Regret Minimization in RL

Consider a finite-horizon MDP $M^* = (S, A, p_h, r_h)$ with stage-dependent transitions and rewards. Assume rewards are bounded in $[0, 1]$. We want to prove a regret upper-bound for UCBVI. We will aim for the suboptimal regret bound ($T = KH$)

$$R(T) = \sum_{k=1}^K V_1^*(s_{1,k}) - V_1^{\pi_k}(s_{1,k}) = \tilde{O}(H^2 S \sqrt{AK})$$

Define the set of plausible MDPs as

$$\mathcal{M}_k = \{M = (S, A, p_{h,k}, r_{h,k}) : r_{h,k}(s, a) \in \beta_{h,k}^r(s, a), p_{h,k}(\cdot | s, a) \in \beta_{h,k}^p(s, a)\}$$

Confidence intervals can be anytime or not.

- Define the event $\mathcal{E} = \{\forall k, M^* \in \mathcal{M}_k\}$. Prove that $\mathbb{P}(\neg\mathcal{E}) \leq \delta/2$. First step, construct a confidence interval for rewards and transitions for each (s, a) using Hoeffding and Weissman inequality (see appendix), respectively. So, we want that

$$\mathbb{P}\left(\forall k, h, s, a : |\hat{r}_{hk}(s, a) - r_h(s, a)| \leq \beta_{hk}^r(s, a) \wedge \|\hat{p}_{hk}(\cdot | s, a) - p_h(\cdot | s, a)\|_1 \leq \beta_{hk}^p(s, a)\right) \geq 1 - \delta/2$$

¹Note that $at \geq \log(bt)$ can be solved using Lambert W function. We thus have $t \geq \frac{-W_{-1}(-a/b)}{a}$ since, given $a = \Delta_i^2$ and $b = 2k/\delta$, $-a/b \in (-1/e, 0)$. We can make the bound more explicit by noticing that $-1 - \sqrt{2u} - u \leq W_{-1}(-e^{-u-1}) \leq -1 - \sqrt{2u} - 2u/3$ for $u > 0$ [Chatzigeorgiou, 2016]. Then $t \geq \frac{1+\sqrt{2u}+u}{a}$ with $u = \log(b/a) - 1$.

- Define the bonus function and consider the Q-function computed at episode k

$$Q_{h,k}(s, a) = \hat{r}_{h,k}(s, a) + b_{h,k}(s, a) + \sum_{s'} \hat{p}_{h,k}(s'|s, a) V_{h+1,k}(s')$$

with $V_{h,k}(s) = \min\{H, \max_a Q_{h,k}(s, a)\}$. Recall that $V_{H+1,k}(s) = V_{H+1}^*(s) = 0$. Prove that under event \mathcal{E} , Q_k is optimistic, i.e.,

$$Q_{h,k}(s, a) \geq Q_h^*(s, a), \forall s, a$$

where Q^* is the optimal Q-function of the unknown MDP M^* . Note that $\hat{r}_{H,k}(s, a) + b_{H,k}(s, a) \geq r_{H,k}(s, a)$ and thus $Q_{H,k}(s, a) \geq Q_H^*(s, a)$ (for a properly defined bonus). Then use induction to prove that this holds for all the stages h .

- In class we have seen that

$$\delta_{1k}(s_{1,k}) \leq \sum_{h=1}^H Q_{hk}(s_{hk}, a_{hk}) - r(s_{hk}, a_{hk}) - \mathbb{E}_{Y \sim p(\cdot|s_{hk}, a_{hk})}[V_{h+1,k}(Y)] + m_{hk} \quad (1)$$

where $\delta_{hk}(s) = V_{hk}(s) - V_h^{\pi_k}(s)$ and $m_{hk} = \mathbb{E}_{Y \sim p(\cdot|s_{hk}, a_{hk})}[\delta_{h+1,k}(Y)] - \delta_{h+1,k}(s_{h+1,k})$. We now want to prove this result. Denote by a_{hk} the action played by the algorithm (you will have to use the greedy property).

1. Show that $V_h^{\pi_k}(s_{hk}) = r(s_{hk}, a_{hk}) + \mathbb{E}_p[V_{h+1,k}(s')] - \delta_{h+1,k}(s_{h+1,k}) - m_{hk}$
 2. Show that $V_{h,k}(s_{hk}) \leq Q_{h,k}(s_{hk}, a_{hk})$.
 3. Putting everything together prove Eq. 1.
- Since $(m_{hk})_{hk}$ is an MDS, using Azuma-Hoeffding we show that with probability at least $1 - \delta/2$

$$\sum_{k,h} m_{hk} \leq 2H\sqrt{KH\log(2/\delta)}$$

Show that the regret is upper bounded with probability $1 - \delta$ by

$$R(T) \leq 2 \sum_{kh} b_{hk}(s_{hk}, a_{hk}) + 2H\sqrt{KH\log(2/\delta)}$$

- Finally, we have that [Domingues et al., 2021]

$$\sum_{h,k} \frac{1}{\sqrt{N_{hk}(s_{hk}, a_{hk})}} \lesssim H^2 S^2 A + 2 \sum_{h=1}^H \sum_{s,a} \sqrt{N_{hK}(s, a)}$$

Complete this by showing an upper-bound of $H\sqrt{SAK}$, which leads to $R(T) \lesssim H^2 S \sqrt{AK}$

A Weissman inequality

Denote by $\hat{p}(\cdot|s, a)$ the estimated transition probability build using n samples drawn from $p(\cdot|s, a)$. Then we have that

$$\mathbb{P}(\|\hat{p}_h(\cdot|s, a) - p_h(\cdot|s, a)\|_1 \geq \epsilon) \leq (2^S - 2) \exp\left(-\frac{n\epsilon^2}{2}\right)$$

References

Ioannis Chatzigeorgiou. Bounds on the Lambert function and their application to the outage analysis of user cooperation. *CoRR*, abs/1601.04895, 2016.

Omar Darwiche Domingues, Pierre Ménard, Matteo Pirotta, Emilie Kaufmann, and Michal Valko. Kernel-based reinforcement learning: A finite-time analysis. In *ICML*, volume 139 of *Proceedings of Machine Learning Research*, pages 2783–2792. PMLR, 2021.

```

Initialize  $Q_{h1}(s, a) = 0$  for all  $(s, a) \in S \times A$  and  $h = 1, \dots, H$ 

for  $k = 1, \dots, K$  do
    Observe initial state  $s_{1k}$  (arbitrary)
    Estimate empirical MDP  $\widehat{M}_k = (S, A, \widehat{p}_{hk}, \widehat{r}_{hk}, H)$  from  $\mathcal{D}_k$ 
        
$$\widehat{p}_{hk}(s'|s, a) = \frac{\sum_{i=1}^{k-1} \mathbb{1}\{(s_{hi}, a_{hi}, s_{h+1,i}) = (s, a, s')\}}{N_{hk}(s, a)}, \quad \widehat{r}_{hk}(s, a) = \frac{\sum_{i=1}^{k-1} r_{hi} \cdot \mathbb{1}\{(s_{hi}, a_{hi}) = (s, a)\}}{N_{hk}(s, a)}$$

    Planning (by backward induction) for  $\pi_{hk}$  using  $\widehat{M}_k$ 
    for  $h = H, \dots, 1$  do
        
$$Q_{h,k}(s, a) = \widehat{r}_{hk}(s, a) + b_{h,k}(s, a) + \sum_{s'} \widehat{p}_{hk}(s'|s, a) V_{h+1,k}(s')$$

        
$$V_{h,k}(s) = \min\{H, \max_a Q_{h,k}(s, a)\}$$

    end
    Define  $\pi_{h,k}(s) = \arg \max_a Q_{h,k}(s, a), \forall s, h$ 
    for  $h = 1, \dots, H$  do
        Execute  $a_{hk} = \pi_{hk}(s_{hk})$ 
        Observe  $r_{hk}$  and  $s_{h+1,k}$ 
        
$$N_{h,k+1}(s_{hk}, a_{hk}) = N_{h,k}(s_{hk}, a_{hk}) + 1$$

    end
end

```

Algorithm 1: UCBVI

Reinforcement Learning - framework 3

Nicolaas Uriolante

① Best Arm Identification

a) Let $\delta > 0$

First, by Hoeffding's inequality

$$P(|\hat{\mu}_{i,t} - \mu_i| > U) \leq 2e^{-\frac{2tU^2}{\delta}} = \delta \Leftrightarrow U(t, \delta) = \sqrt{\frac{1}{2t} \log\left(\frac{2}{\delta}\right)}$$

$$\text{Let } \delta = \frac{6}{\pi^2 K} \frac{1}{t^2}$$

$$\begin{aligned} P\left(\bigcup_{i=1}^K \bigcup_{t=1}^{\infty} \{|\hat{\mu}_{i,t} - \mu_i| > U(t, \delta)\}\right) &\leq \sum_{i=1}^K \sum_{t=1}^{\infty} \frac{6}{\pi^2 K} \frac{1}{t^2} \\ &= \frac{\delta \cdot 6}{\pi^2 K} \left(\sum_{i=1}^K 1\right) \left(\sum_{t=1}^{\infty} \frac{1}{t^2}\right) \\ &= \frac{\delta \cdot 6 \cdot K \cdot \pi^2}{\pi^2 K} = \delta \end{aligned}$$

b)

$\exists j \in S / \{|\hat{\mu}_{j,t} - U(t, \delta)| > |\hat{\mu}_{i,t} + U(t, \delta)|\}$
 then μ^* is outside the confidence interval
 $\Rightarrow P(i^* \text{ removed from } S) \leq P(\varepsilon_0) \leq \delta$
 $\Rightarrow P(i^* \text{ remains in } S) \geq 1 - \delta$

c) Under ε^c , $\Delta_i = \mu^* - \mu_i \geq C_1 U(t, \delta)$

The worst case is

then we need $C_1 = 4$ at least

Then, if $\Delta_i \geq 4 U(t, \delta)$ the arm i will be removed

Analytically, since ε^c holds: $\mu^* \leq \hat{\mu}^* + U(t, \delta)$ and $\mu_i \geq \hat{\mu}_i - U(t, \delta)$

The arm elimination condition is

$$\hat{\mu}^* - U(t, \delta) \geq \hat{\mu}_i + U(t, \delta)$$

$$\Leftrightarrow \mu^* - U(t, \delta) - U(t, \delta) \geq \mu_i + U(t, \delta) + U(t, \delta)$$

$$\Leftrightarrow \mu^* - \mu_i \geq 4 U(t, \delta)$$

(continue)

$$\textcircled{1} \quad (\text{c}) \quad \Delta_i \geq \frac{1}{4} \sqrt{\frac{1}{2t} \log(b^2 t^2)} \quad \text{with } b = \frac{9\pi^2 K}{68} = \frac{\pi^2 K}{38}$$

$$\Delta_i^2 \geq \frac{16}{D_2} \log(b^2 t^2) = \frac{16}{t} \log(bt) \Leftrightarrow \left(\frac{\Delta_i^2}{16} \right) \cdot t \geq \log(bt)$$

Then, using Lambert W $t \geq \frac{1 + \sqrt{1 + 2u}}{a} + u$ with $u = \log(b/a) - 1$

$$t \geq \frac{1 + \sqrt{1 + 2(\log(b/a) - 1)}}{a} + \log(b/a) - 1$$

$$t \geq \frac{16}{\Delta_i^2} \left\{ \sqrt{2 \log\left(\frac{16 \pi^2 \sqrt{K}}{38}\right) - 2} + \log\left(\frac{16 \pi^2 \sqrt{K}}{38}\right) \right\} \quad (*)$$

(d) The equation (*) is for the arm i and has complexity

$$\Theta\left(\frac{1}{\Delta_i^2} \log\left(\frac{1}{\Delta_i^2} \sqrt{\frac{K}{S}}\right)\right) \quad (\text{simplifying constants and lower-order terms})$$

If we consider all the pulls we have a sum over arms $i \neq i^*$, so

$$\Theta\left(\sum_{i \neq i^*} \frac{1}{\Delta_i^2} \log\left(\frac{1}{\Delta_i^2} \sqrt{\frac{K}{S}}\right)\right) \quad \text{is the sample complexity}$$

(e) The algorithm won't work correctly. Between two optimal arms we have $\Delta = 0$, hence an infinite time for the removal of the arm, as given by equation (*).

Reinforcement Learning - Homework 3

Nicola's Violante

② Regret Minimization in RL

$$1) \mathbb{P}(|\hat{r}_{hk}(s, a) - r_h(s, a)| \geq \beta_{hk}^r(s, a)) \leq 2e^{-2U_{hk}(s, a)\beta_{hk}^r(s, a)^2} = \delta_1$$

$$\Leftrightarrow \beta_{hk}^r(s, a) = \sqrt{\frac{1}{2N_{hk}(s, a)} \log\left(\frac{2}{\delta_1}\right)}, \quad (1)$$

$$\mathbb{P}(\|\hat{p}_{hk}(\cdot|s, a) - p_h(\cdot|s, a)\|_1 \geq \beta_{hk}^p(s, a)) \leq (2^S - 2) e^{-\frac{2U_{hk}(s, a)}{2}\beta_{hk}^p(s, a)^2} = \delta_2$$

$$\frac{1}{2} U_{hk}(s, a) \beta_{hk}^p = \log\left(\frac{2^S - 2}{\delta_2}\right)$$

$$\Rightarrow \beta_{hk}^p = \sqrt{\frac{2}{U_{hk}(s, a)} \log\left(\frac{2^S - 2}{\delta_2}\right)}, \quad (2)$$

• Let h, k, s, a be fixed and let $S' > 0$:

$$\mathbb{P}(|\hat{r}_{hk}(s, a) - r_h(s, a)| \geq \beta_{hk}^r(s, a) \text{ or } \|\hat{p}_{hk}(\cdot|s, a) - p_h(\cdot|s, a)\|_1 \geq \beta_{hk}^p(s, a))$$

$$\leq \mathbb{P}(|\hat{r}_{hk}(s, a) - r_h(s, a)| \geq \beta_{hk}^r(s, a)) + \mathbb{P}(\|\hat{p}_{hk}(\cdot|s, a) - p_h(\cdot|s, a)\|_1 \geq \beta_{hk}^p(s, a))$$

$$\leq \frac{\delta'}{2} + \frac{\delta'}{2} = S' \text{ for a proper } \beta_{hk}^r(s, a) \text{ and } \beta_{hk}^p(s, a)$$

$$\mathbb{P}(\forall K, h, s, a : |\hat{r}_{hk}(s, a) - r_h(s, a)| \geq \beta_{hk}^r(s, a) \text{ and } \|\hat{p}_{hk}(\cdot|s, a) - p_h(\cdot|s, a)\|_1 \geq \beta_{hk}^p(s, a))$$

$$= 1 - \mathbb{P}\left(\bigcup_{h=1}^K \bigcup_{s=1}^S \bigcup_{a=1}^A \{|\hat{r}_{hk}(s, a) - r_h(s, a)| \geq \beta_{hk}^r(s, a) \text{ or } \|\hat{p}_{hk}(\cdot|s, a) - p_h(\cdot|s, a)\|_1 \geq \beta_{hk}^p(s, a)\}\right)$$

$$\geq 1 - \sum_{h=1}^K \sum_{s=1}^S \sum_{a=1}^A \mathbb{P}(|\hat{r}_{hk}(s, a) - r_h(s, a)| \geq \beta_{hk}^r(s, a) \text{ or } \|\hat{p}_{hk}(\cdot|s, a) - p_h(\cdot|s, a)\|_1 \geq \beta_{hk}^p(s, a))$$

$$= 1 - \underbrace{KHS}_{:= \delta} \Rightarrow \delta = 2KHS\delta' \Rightarrow \frac{\delta'}{2} = \frac{\delta}{4KHS}$$

$$\Rightarrow \beta_{hk}^r(s, a) = \sqrt{\frac{1}{2U_{hk}(s, a)} \log(8KHS)}$$

} then $\mathbb{P}(\gamma \varepsilon) \geq 1 - \frac{\delta}{2}$

$$\beta_{hk}^p = \sqrt{\frac{2}{U_{hk}(s, a)} \left[\log\left(\frac{(2^S - 2)4KHS}{8}\right) \right]}$$

Reinforcement Learning - Homework 3

Nicola's Violante.

$$(2) b) Q_{n,k}(s,a) = \hat{r}_{n,k}(s,a) + b_{n,k}(s,a) + \sum_{s'} \hat{p}_{n,k}(s'|s,a) V_{n+1,k}(s') \text{ under } \mathcal{E}$$

• Base case $k=1$

$$Q_{1,k}(s,a) = \hat{r}_{1,k}(s,a) + b_{1,k}(s,a) + \sum_{s'} \hat{p}_{1,k}(s'|s,a) \underbrace{V_{(1+1),k}(s')}_{=0} \\ = \hat{r}_{1,k}(s,a) + b_{1,k}(s,a)$$

$$Q_1^*(s,a) = r_{1,k}(s,a)$$

$$\Rightarrow Q_{1,k}(s,a) \geq Q_1^*(s,a) \text{ if } b_{1,k}(s,a) \geq \beta_{1,k}^P(s,a) \Rightarrow$$

$$\hat{r}_{1,k}(s,a) + b_{1,k}(s,a) \geq \hat{r}_{1,k}(s,a)$$

• Inductive step

$$\text{Assume } Q_{n+1,k}(s,a) \geq Q_{n+1,k}^*(s,a) \forall s,a$$

$$Q_{n,k}(s,a) - Q_n^*(s,a) = \hat{r}_{n,k}(s,a) + b_{n,k}(s,a) + \sum_{s'} \hat{p}_{n,k}(s'|s,a) \min \left\{ 1, \max_a Q_{n+1,k}(s',a) \right\} \\ - r_n(s,a) - \sum_{s'} p_n(s'|s,a) \max_a Q_{n+1}^*(s',a) \\ \geq \hat{r}_{n,k}(s,a) + b_{n,k}(s,a) - r_n(s,a) \\ + \sum_{s'} \min \left\{ 1, \max_a Q_{n+1}^*(s',a) \right\} (\hat{p}_{n,k}(s'|s,a) - p_n(s'|s,a)) \\ \geq \hat{r}_{n,k}(s,a) + b_{n,k}(s,a) - r_n(s,a) \\ - \sum_{s'} \min \left\{ 1, \max_a Q_{n+1}^*(s',a) \right\} |\hat{p}_{n,k}(s'|s,a) - p_n(s'|s,a)| \\ \geq \hat{r}_{n,k}(s,a) + b_{n,k}(s,a) - r_n(s,a) - H \underbrace{\sum_{s'} |\hat{p}_{n,k}(s'|s,a) - p_n(s'|s,a)|}_{||\hat{p}_{n,k}(s'|s,a) - p_n(s'|s,a)||_1} \\ \geq \hat{r}_{n,k}(s,a) + b_{n,k}(s,a) - r_n(s,a) - H \beta_{n,k}^P(s,a) \\ \geq - \beta_{n,k}^P(s,a) + b_{n,k}(s,a) - H \beta_{n,k}^P(s,a) \geq 0 \\ \Rightarrow b_{n,k}(s,a) \geq \beta_{n,k}^P(s,a) + H \beta_{n,k}^P(s,a)$$

Reinforcement Learning - Homework 3 Niculae's Violante

② 3) 1.

$$V_n^{th}(s) = r(s_{hk}, a_{hk}) + \sum_{s'} p(s' | s_{hk}, a_{hk}) V_{n+1}^{th}(s')$$

$$= \delta_{hk, t}(s') + V_{n+1, t}(s')$$

$$\Rightarrow \sum_{s'} p(s' | s_{hk}, a_{hk}) (V_{n+1, t}(s') - \delta_{hk, t}(s')) =$$

$$= E_{s' \sim p} (V_{n+1, t}(s)) - E_{s' \sim p} [\delta_{hk, t}(s')] = E_{s' \sim p} [V_{n+1, t}(s)] - \delta_{hk, t}(s_{hk, t}) - m_{hk}$$

$$\Rightarrow V_n^{th}(s) = r(s_{hk}, a_{hk}) + E_{s' \sim p} (V_{n+1, t}(s')) - \delta_{hk, t}(s_{hk, t}) - m_{hk}$$

$$2. V_{n, t}(s_{hk}) = \min \left\{ H, \max_a Q_{n, t}(s_{hk}, a) \right\} = \min \left\{ H, Q_{n, t}(s_{hk}, a_{hk}) \right\} \\ \leq Q_{n, t}(s_{hk}, a_{hk})$$

$$3. S_{1k}(s_{1k}) = V_{1k}(s_{1k}) - V_1^{th}(s_{1k})$$

$$\leq Q_{1k}(s_{1k}, a_{1k}) - r(s_{1k}, a_{1k}) - E_p [V_{2k}(s')] + m_{1k} + \delta_{2k}(s_{2k})$$

$$\leq Q_{1k}(s_{1k}, a_{1k}) - r(s_{1k}, a_{1k}) - E_p [V_{2k}(s')] + m_{1k}$$

$$+ Q_{2k}(s_{2k}, a_{2k}) - r(s_{2k}, a_{2k}) - E_p [V_{2k}(s')] + m_{2k} + \delta_{3k}(s_{3k})$$

⋮

⋮

$$\leq \sum_{h=1}^H (Q_{hk}(s_{hk}, a_{hk}) - r(s_{hk}, a_{hk}) - E_p [V_{hk}(s')]) + m_{1k}$$

$$4) R(T) = \sum_{h=1}^K V_1^*(s_{1k}) - V_1^{th}(s_{1k}) \leq \sum_{h=1}^T S_{1k}(s_{1k}) \quad \text{since } Q_{1k}(s, a) > Q_1^*(s, a) \forall s, a$$

$$\leq \sum_{h=2}^K \sum_{n_k} (Q_{hk}(s_{hk}, a_{hk}) - r(s_{hk}, a_{hk}) - E_p [V_{h+1, k}(s')]) + m_{1k}$$

$$= \sum_{h=1}^K \sum_{n_k} \left\{ \hat{r}_{hk}(s_{hk}, a_{hk}) - r(s_{hk}, a_{hk}) + \sum_{s'} |\hat{p}_{hk}(s' | s_{hk}, a_{hk}) - p_{hk}(s' | s_{hk}, a_{hk})| V_{h+1, k}(s') \right\} +$$

$$+ \sum_{n_k} m_{hk} + \sum_{n_k} b_{hk}(s_{hk}, a_{hk})$$

using the bounds $\sum_{h=1}^K \sum_{n_k} |\hat{r}_{hk}(s_{hk}, a_{hk}) - r(s_{hk}, a_{hk})| + H \|\hat{p}_{hk}(\cdot | s_{hk}, a_{hk}) - p_{hk}(\cdot | s_{hk}, a_{hk})\|_1$
 found at previous points:

$$\text{with prob at least } 1-\delta \quad + \sum_{n_k} m_{hk} + \sum_{n_k} b_{hk}(s_{hk}, a_{hk}) \leq b_{hk}(s_{hk}, a_{hk})$$

$$\leq \sum_{n_k} \beta_{hk}^*(s_{hk}, a_{hk}) + H \beta_{hk}^*(s_{hk}, a_{hk}) + b_{hk}(s_{hk}, a_{hk}) + 2H \sqrt{K H \log(\frac{1}{\delta})}$$

$$\leq \sum_{n_k} 2b_{hk}(s_{hk}, a_{hk}) + 2H \sqrt{K H \log(\frac{1}{\delta})}$$

5) Since V^* is concave, by Jensen's inequality

$$\begin{aligned} \sum_{h=1}^H \sum_{s,a} \sqrt{N_{hk}(s,a)} &= HSA \sum_{h=1}^H \sum_{s,a} \frac{1}{HSA} \sqrt{N_{hk}(s,a)} \\ &\leq HSA \sqrt{\frac{1}{HSA} \sum_{h=1}^H \sum_{s,a} N_{hk}(s,a)} \\ &\leq \sqrt{HSA} \sqrt{\sum_{h=1}^H K} = H\sqrt{SAK} \end{aligned}$$

Now we chose $b_{hk}(s_{hk}, a_{hk}) = \beta_{hk}^r(s_{hk}, a_{hk}) + H\beta_{hk}^p(s_{hk}, a_{hk})$

Also since $S \gg 1 \rightarrow \log(2^S - 2) \approx \log 2^S = S \log 2$

$$\begin{aligned} P(t) &\leq 2 \sum_{h_k} \sqrt{\frac{1}{2N_{hk}} \log(8KUSA)} + H \sqrt{\frac{2}{N_{hk}(s_{hk}, a_{hk})} S \log 2} \log\left(\frac{4KUSA}{8}\right) \\ &\quad + 2H \sqrt{Kt \log\left(\frac{D}{S}\right)} \\ &\leq H^2 S^2 A + H\sqrt{SAK} + t \sqrt{S} (H^2 S^2 A + H\sqrt{SAK}) + 2H \sqrt{KH \log\left(\frac{D}{8}\right)} \\ &\approx t^2 S \sqrt{AK} \end{aligned}$$