

## Exploration in Reinforcement Learning (theory)

Lecturers: *M. Pirotta*( *December 16, 2021* )

Solution by Megi Dervishi

**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  $\mu_i$ . 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.

**$\delta$ -correctness and fixed-confidence objective.** Denote by  $\tau_\delta$  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  $\delta$ -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  $\mu_1, \dots, \mu_k$ . Our goal is to find a  $\delta$ -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  $t$ .
- $\mu_i$ : 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  $\delta'$ . This is called "bad event" since it means that the confidence intervals do not hold.

===== **Answer:** =====

First we want to find  $U(t, \delta)$  such that:

$$\mathbb{P} \left( \bigcup_{t=1}^{+\infty} |\hat{\mu}_{i,t} - \mu_i| > U(t, \delta') \right) \leq \delta' \quad (1)$$

Note that:

$$\hat{\mu}_{i,t} = \frac{1}{t} \sum_{j=1}^t X_{i,j} \quad \text{and} \quad \mathbb{E}[\hat{\mu}_{i,t}] = \frac{1}{t} \sum_{j=1}^t \mathbb{E}[X_{i,j}] = \frac{1}{t} \sum_{j=1}^t \mu_i = \mu_i \quad (2)$$

Hence we can re-write inequality (1) as:

$$\begin{aligned} \mathbb{P} \left( \bigcup_{t=1}^{+\infty} |\hat{\mu}_{i,t} - \mu_i| > U(t, \delta') \right) &\leq \sum_{t=1}^{+\infty} \mathbb{P} (|\hat{\mu}_{i,t} - \mu_i| > U(t, \delta')) \\ &= \sum_{t=1}^{+\infty} \mathbb{P} (|\hat{\mu}_{i,t} - \mathbb{E}[\hat{\mu}_{i,t}]| > U(t, \delta')) \\ &= \sum_{t=1}^{+\infty} \mathbb{P} \left( \left| \sum_{j=1}^t X_{i,j} - \mathbb{E} \left[ \sum_{j=1}^t X_{i,j} \right] \right| > t U(t, \delta') \right) \end{aligned} \quad (3)$$

Using Hoeffding's inequality we can express (3) as:

$$\sum_{t=1}^{+\infty} 2 \exp \left( -2 \frac{t^2 U(t, \delta')^2}{\sum_{j=1}^t (b_j - a_j)^2} \right) \quad \text{where } a_j < X_{i,j} < b_j \quad (4)$$

We know  $X_{i,j} \in [0, 1]$  which gives  $a_j = 0$  and  $b_j = 1$ , which results to:

$$\sum_{t=1}^{+\infty} 2 \exp (-2tU(t, \delta')^2) \leq \delta' \quad (5)$$

For the series to be convergent we can use the Riemann criterion and the dominated-convergence

theorem which constrain  $U(t, \delta')$  as follows

$$\begin{aligned} 2 \exp(-2tU(t, \delta')^2) &< \frac{\delta'}{t} \\ 2tU(t, \delta')^2 &> \ln\left(\frac{2t}{\delta'}\right) \\ U(t, \delta') &> \sqrt{\frac{1}{2t} \ln\left(\frac{2t}{\delta'}\right)} \end{aligned} \quad (6)$$

(7)

Finally we have:

$$\mathbb{P}(\mathcal{E}) \leq \sum_{i=1}^k \mathbb{P}\left(\bigcup_{t=1}^{+\infty} |\hat{\mu}_{i,t} - \mu_i| > U(t, \delta')\right) \leq \sum_{i=1}^k \delta' \leq k\delta' = \delta \quad (8)$$

- 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  $\delta'$  and start from the condition for arm elimination. From this, use the definition of  $\neg \mathcal{E}$ .

=====Answer=====

For arm  $i^* = \arg \max_i \{\mu_i\}$  to remain in the active set means that

$$\bigcap_{t=1}^{+\infty} \bigcap_{j \in S_t} \{\hat{\mu}_{j,t} - U(t, \delta') < \hat{\mu}_{i^*,t} + U(t, \delta')\} \Leftrightarrow \bigcap_{t=1}^{+\infty} \bigcap_{j \in S_t} \{\hat{\mu}_{j,t} - \hat{\mu}_{i^*,t} < 2U(t, \delta')\} \quad (9)$$

Since  $i^* = \arg \max_i \{\mu_i\}$  we know that  $\forall j, \mu_j \leq \mu_{i^*}$  hence:

$$\bigcap_{t=1}^{+\infty} \bigcap_{j \in S_t} \{\hat{\mu}_{j,t} - \mu_j + \mu_j - \hat{\mu}_{i^*,t} < 2U(t, \delta')\} \supseteq \bigcap_{t=1}^{+\infty} \bigcap_{j \in S_t} \{\hat{\mu}_{j,t} - \mu_j + \mu_{i^*} - \hat{\mu}_{i^*,t} < 2U(t, \delta')\} \quad (10)$$

Using the triangular inequality and the fact that  $|a| > a$  we can further simplify to obtain

$$\bigcap_{t=1}^{+\infty} \bigcap_{j \in S_t} \{\hat{\mu}_{j,t} - \mu_j + \mu_{i^*} - \hat{\mu}_{i^*,t} < 2U(t, \delta')\} \supseteq \bigcap_{t=1}^{+\infty} \bigcap_{j \in S_t} \{|\hat{\mu}_{j,t} - \mu_j| + |\mu_{i^*} - \hat{\mu}_{i^*,t}| < 2U(t, \delta')\} \quad (11)$$

Finally noting that  $S_t \subseteq [1, k] \cap \mathbb{N}$  we can obtain the final expression

$$\bigcap_{t=1}^{+\infty} \bigcap_{j \in S_t} \{|\hat{\mu}_{j,t} - \mu_j| + |\mu_{i^*} - \hat{\mu}_{i^*,t}| < 2U(t, \delta')\} \supseteq \bigcap_{t=1}^{+\infty} \bigcap_{j=1}^k \{|\hat{\mu}_{j,t} - \mu_j| < U(t, \delta')\} = \neg \mathcal{E} \quad (12)$$

where we used that

$$\forall t \in [1, +\infty[ \cap \mathbb{N}, \forall j \in [1, k[ \cap \mathbb{N}, |\hat{\mu}_{j,t} - \mu_j| < U(t, \delta') \Rightarrow \begin{cases} |\hat{\mu}_{j,t} - \mu_j| < U(t, \delta') \\ \forall t, j, |\mu_{i^*} - \hat{\mu}_{i^*,t}| < U(t, \delta') \end{cases} \quad (13)$$

$$\Rightarrow \forall t \in [1, +\infty[ \cap \mathbb{N}, \forall j \in S_t \cap \mathbb{N}, |\hat{\mu}_{j,t} - \mu_j| + |\mu_{i^*} - \hat{\mu}_{i^*,t}| < 2U(t, \delta') \quad (14)$$

Hence we have that

$$\mathbb{P}\left(\bigcap_{t=1}^{+\infty} \bigcap_{j \in S_t} \{\hat{\mu}_{j,t} - \hat{\mu}_{i^*,t} < 2U(t, \delta')\}\right) \geq \mathbb{P}(\neg \mathcal{E}) = 1 - \delta \quad (15)$$

- 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  $i^*$ .<sup>1</sup>

=====Answer:=====

Suppose we are in the event  $\neg\mathcal{E}$ . Hence

$$\forall t \in [1, +\infty[ \cap \mathbb{N}, \forall j \in [1, k] \cap \mathbb{N}, \quad |\hat{\mu}_{j,t} - \mu_j| < U(t, \delta') \quad (16)$$

Hence

$$\Delta_i \geq C_1 U(t, \delta') \Leftrightarrow \mu^* - \mu_i \geq C_1 U(t, \delta') \quad (17)$$

Now from the previous question we know that  $\neg\mathcal{E} \subseteq \{\text{The optimal arm is never removed.}\}$ . Hence a non-optimal arm  $j$  can be removed if for some  $t$  we have that

$$\hat{\mu}_{i^*,t} - U(t, \delta') \geq \hat{\mu}_{j,t} + U(t, \delta') \Leftrightarrow \hat{\mu}_{i^*,t} - \hat{\mu}_{j,t} \geq 2U(t, \delta') \quad (18)$$

Now starting from the separation hypothesis

$$\Delta_j \geq C_1 U(t, \delta') \Leftrightarrow (\mu_{i^*} - \hat{\mu}_{i^*,t}) + (\hat{\mu}_{j,t} - \mu_j) + (\hat{\mu}_{i^*,t} - \hat{\mu}_{j,t}) \geq C_1 U(t, \delta') \quad (19)$$

$$\Leftrightarrow \hat{\mu}_{i^*,t} - \hat{\mu}_{j,t} \geq C_1 U(t, \delta') - (\mu_{i^*} - \hat{\mu}_{i^*,t}) - (\hat{\mu}_{j,t} - \mu_j) \quad (20)$$

$$\Rightarrow \hat{\mu}_{i^*,t} - \hat{\mu}_{j,t} \geq C_1 U(t, \delta') - |\mu_{i^*} - \hat{\mu}_{i^*,t}| - |\hat{\mu}_{j,t} - \mu_j| \quad (21)$$

$$\Rightarrow \hat{\mu}_{i^*,t} - \hat{\mu}_{j,t} \geq C_1 U(t, \delta') - 2U(t, \delta') = (C_1 - 2)U(t, \delta') \quad (22)$$

To pass from 20 to 21 we simply use that  $a < |a|$  and to pass from 21 to 22 we use 16. Hence we see that for  $C_1 \geq 4$  then if  $\Delta_j \geq C_1 U(t, \delta')$  the arm  $j$  will be removed. In order to have the tightest bound we take  $C_1 = 4$ . Now we take the previously computed explicit form of  $U(t, \delta')$

$$U(t, \delta') > \sqrt{\frac{1}{2t} \ln \left( \frac{2t}{\delta'} \right)} \quad (23)$$

where  $\delta' = \delta/k$ , then the condition for arm  $j$  is verified if

$$\sqrt{\frac{1}{2t} \ln \left( \frac{2kt}{\delta} \right)} \leq \frac{\Delta_j}{4} \Leftrightarrow \ln \left( \frac{2kt}{\delta} \right) \leq \frac{\Delta_j^2}{8} t \quad (24)$$

Note that this is the inequality detailed in the footnote where

$$a = \frac{\Delta_j^2}{8} = \tilde{\Delta}_j^2 \quad \text{and} \quad b = \frac{2k}{\delta} \quad (25)$$

Hence we know that it can be written as

$$t \geq \frac{-W_{-1}(a/b)}{a} \quad (26)$$

and we can get a more explicit bound by setting

$$u = \log(b/a) - 1 = \log \left( \frac{2k}{\delta \tilde{\Delta}_j^2} \right) - 1 \quad (27)$$

Then

$$t \geq \frac{1 + \sqrt{2u} + u}{a} = \left( \tilde{\Delta}_j^{-2} \log \left( \frac{2k}{\delta \tilde{\Delta}_j^2} \right) \right) + \frac{\sqrt{2u}}{a} \gtrsim \left( \tilde{\Delta}_j^{-2} \log \left( \frac{2k}{\delta \tilde{\Delta}_j^2} \right) \right) \quad (28)$$

---

<sup>1</sup>Note 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  $u = \log(b/a) - 1$ .

- Compute a bound on the sample complexity (after how many *pulls* the algorithm stops) for identifying the optimal arm w.p.  $1 - \delta$ .

=====Answer:=====

Since to eliminate arm  $j$  we need

$$t \geq \frac{1}{a_j} (1 + \sqrt{2u_j} + u_j) \sim \frac{1}{\bar{\Delta}_j^2} \left( 1 + \sqrt{2 \log \left( \frac{1}{\bar{\Delta}_j^2} \right) + \log \left( \frac{1}{\bar{\Delta}_j^2} \right)} \right) \quad (29)$$

which is a strictly decreasing function of  $\Delta_j$  hence for all non-optimal arms to be suppressed and to identify the optimal arm with probability  $1 - \delta$  we need at least

$$t \geq \frac{1}{a^*} (1 + \sqrt{2u^*} + u^*) \quad \text{where } a^*, u^* \text{ are defined as above with } \min_j \Delta_j \quad (30)$$

- We assumed that the optimal arm  $i^*$  is unique. Would the algorithm still work if there exist multiple best arms? Why?

=====Answer:=====

Yes it would also work however instead of being left with one optimal arm we might be left with a set of the multiple best arms even at long times.

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 Weissmain inequality (see appendix), respectively. So, we want that

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

=====Answer:=====

Similarly to what was done before we start by computing the any-time confidence interval for transitions. Indeed from the Weissmain inequality

$$\mathbb{P}(\|\hat{p}_{hk}(\cdot | s, a) - p_h(\cdot | s, a)\|_1 \geq \beta_{hk}^p(s, a)) \leq (2^S - 2) \exp\left(-\frac{N_{hk}(s, a) \beta_{hk}^p(s, a)^2}{2}\right) \quad (31)$$

and similarly as before

$$\mathbb{P}(\exists k, h, s, a : \|\hat{p}_{hk}(\cdot|s, a) - p_h(\cdot|s, a)\|_1 \geq \beta_{hk}^p(s, a)) \leq \sum_{k, h, s, a} (2^S - 2) \exp\left(-\frac{N_{hk}(s, a) \beta_{hk}^p(s, a)^2}{2}\right) \quad (32)$$

hence choosing  $\beta_{hk}^p(s, a)$  such that

$$(2^S - 2) \exp\left(-\frac{N_{hk}(s, a) \beta_{hk}^p(s, a)^2}{2}\right) = \frac{\delta'_p}{HSAK} \quad (33)$$

which is readily satisfied by

$$\beta_{hk}^p(s, a) = \sqrt{\frac{2 \log\left(\frac{(2^S - 2)AHKS}{\delta'_p}\right)}{N_{hk}(s, a)}} \quad (34)$$

leads to

$$\mathbb{P}(\exists k, h, s, a : \|\hat{p}_{hk}(\cdot|s, a) - p_h(\cdot|s, a)\|_1 \geq \beta_{hk}^p(s, a)) \leq \delta'_p \quad (35)$$

The exact same reasoning can be repeated for the rewards using Hoeffding's inequality instead of Weissmain analogously to what was done in the previous section (note that rewards are bounded in  $[0, 1]$  so the denominator simplifies same as before).

$$\mathbb{P}(|\hat{r}_{hk}(s, a) - r_h(s, a)| \geq \beta_{hk}^r(s, a)) \leq 2 \exp\left(-\frac{2}{h} \beta_{hk}^r(s, a)^2\right) \quad (36)$$

and the equation

$$2 \exp\left(-\frac{2}{H} \beta_{hk}^r(s, a)^2\right) = \frac{\delta'_r}{HSAK} \quad (37)$$

is solved by

$$\beta_{hk}^r(s, a) = \sqrt{\frac{H}{2} \log\left(\frac{2AHKS}{\delta'_r}\right)} \quad (38)$$

which leads to

$$\mathbb{P}(\exists k, h, s, a : |\hat{r}_{hk}(s, a) - r_h(s, a)| \geq \beta_{hk}^r(s, a)) \leq \delta'_r \quad (39)$$

Then we have that

$$\mathbb{P}(\exists k, h, s, a : \|\hat{p}_{hk}(\cdot|s, a) - p_h(\cdot|s, a)\|_1 \geq \beta_{hk}^p(s, a) \vee |\hat{r}_{hk}(s, a) - r_h(s, a)| > \beta_{hk}^r(s, a)) \leq \delta'_p + \delta'_r = \delta/2 \quad (40)$$

where we simply took  $\delta'_p = \delta'_r = \delta/4$ . Hence

$$\mathbb{P}(\neg \mathcal{E}) \leq \delta/2 \quad (41)$$

- 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  $h$ .

=====Answer=====

Since we are in event  $\mathcal{E}$  we know that

$$\forall k, h, s, a, \quad |\hat{r}_{hk}(s, a) - r_h(s, a)| \leq \beta_{hk}^r(s, a) \quad (42)$$

Hence if we take

$$b_{h,k}(s, a) = \beta_{hk}^r(s, a) + f(h, k, \beta_{hk}^p(s, a)) \quad (43)$$

with only constraint that  $f(H, k, \beta_{hk}^p(s, a)) \geq 0$  then

$$Q_{Hk}(s, a) = \hat{r}_{Hk}(s, a) + b_{Hk}(s, a) \geq r_H(s, a) = Q_H^*(s, a) \quad (44)$$

Hence the property holds for  $h = H$ , now assume it holds for  $h + 1 \in [0, H] \cap \mathbb{N}$ , we want to show that it holds for  $h$ . We have that

$$Q_h^*(s, a) = r_h(s, a) + \sum_{s'} p_h(s'|s, a) V_{h+1}(s') \quad (45)$$

now we use our induction hypothesis

$$V_{h+1}(s) = \max_a Q_{h+1}(s, a) \leq \max_a Q_{h+1,k}(s, a) = V_{h+1,k}(s, a) \quad (46)$$

hence we have that

$$Q_{h,k}(s, a) - Q_h^*(s, a) \geq (\hat{r}_{h,k}(s, a) - r_h(s, a)) + \sum_{s'} (\hat{p}_{h,k}(s'|s, a) - p_h(s'|s, a)) V_{h+1,k}(s') + b_{h,k}(s, a) \quad (47)$$

Now notice that

$$\sum_{s'} (\hat{p}_{h,k}(s'|s, a) - p_h(s'|s, a)) V_{h+1,k}(s') \leq \|\hat{p}_{h,k}(\cdot|s, a) - p_h(\cdot|s, a)\|_1 \|V_{h+1,k}\|_\infty \quad (48)$$

$$\leq H \|\hat{p}_{h,k}(\cdot|s, a) - p_h(\cdot|s, a)\|_1 \quad (49)$$

and since we are in event  $\mathcal{E}$  we further have that

$$\sum_{s'} (\hat{p}_{h,k}(s'|s, a) - p_h(s'|s, a)) V_{h+1,k}(s') \leq H \beta_{hk}^p(s, a) \quad (50)$$

Hence taking

$$b_{h,k}(s, a) = \beta_{h,k}^r(s, a) + H \beta_{hk}^p(s, a) \quad (51)$$

ensures that if we are in  $\mathcal{E}$  then 47 reduces to

$$Q_{h,k}(s, a) - Q_h^*(s, a) \geq 0 \quad (52)$$

Hence we have proved by induction that

$$\forall h, k, s, a, \quad Q_{h,k}(s, a) \geq Q_h^*(s, a) \quad (53)$$

if we are in event  $\mathcal{E}$ .

- In class we have seen that

$$\delta_{1k}(s_{1,k}) \leq \sum_{h=1}^H Q_{hk}(s_{hk}, a_{hk}) - r(s_{hk}, a_{hk}) - \mathbb{E}_{Y \sim p(\cdot|s_{hk}, a_{hk})} [V_{h+1,k}(Y)] + m_{hk} \quad (54)$$

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  $V_h^{\pi_k}(s_{hk}) = r(s_{hk}, a_{hk}) + \mathbb{E}_p[V_{h+1,k}(s')] - \delta_{h+1,k}(s_{h+1,k}) - m_{h,k}$
2. Show that  $V_{h,k}(s_{hk}) \leq Q_{h,k}(s_{hk}, a_{hk})$ .
3. Putting everything together prove Eq. 54.

===== **Answer:** =====

1. We have

$$\begin{aligned}
 r(s_{hk}, a_{hk}) + \mathbb{E}_p[V_{h+1,k}(s')] - \delta_{h+1,k}(s_{h+1,k}) - m_{h,k} \\
 &= r(s_{hk}, a_{hk}) + \mathbb{E}_p[V_{h+1,k}(s')] - \mathbb{E}_{Y \sim p(\cdot | s_{hk}, a_{hk})}[\delta_{h+1,k}(Y)] \\
 &= r(s_{hk}, a_{hk}) + \mathbb{E}_p[V_{h+1,k}(s')] + \mathbb{E}_{Y \sim p(\cdot | s_{hk}, a_{hk})}[V_{h+1}^{\pi_k}(Y) - V_{h+1,k}(Y)] \\
 &= r(s_{hk}, a_{hk}) + \mathbb{E}_{Y \sim p(\cdot | s_{hk}, a_{hk})}[V_{h+1}^{\pi_k}(Y)] \\
 &= V_h^{\pi_k}(s_{hk})
 \end{aligned}$$

2. From the algorithm we see that

$$V_{h,k}(s_{hk}) = \min\{H, \max_a Q_{h,k}(s_{hk}, a)\} \leq \max_a Q_{h,k}(s_{hk}, a) = Q_{h,k}(s_{hk}, \arg \max_a Q_{h,k}(s_{hk}, a)) \quad (55)$$

From the definition of  $a_{hk}$  we then have

$$V_{h,k}(s_{hk}) \leq Q_{h,k}(s_{hk}, a_{hk}) \quad (56)$$

3. We have

$$\delta_{1k}(s_{1k}) = V_{1k}(s_{1k}) - V_1^{\pi_k}(s_{1k}) \quad (57)$$

Now using our answer to question 1 we write

$$\delta_{1k}(s_{1k}) = V_{1k}(s_{1k}) - r(s_{1k}, a_{1k}) - \mathbb{E}_p[V_{2,k}(s')] + \delta_{2,k}(s_{2,k}) + m_{1,k} \quad (58)$$

now using Question 1 again to repeatedly expand  $\delta_{2,k}(s_{2,k}), \delta_{3,k}(s_{3,k}), \dots$  we get

$$\delta_{1,k}(s_{1,k}) = \delta_{H+1,k}(s_{H+1,k}) + \sum_{h=1}^H V_{hk}(s_{hk}) - r(s_{hk}, a_{hk}) - \mathbb{E}_p[V_{h+1,k}(s')] + m_{hk} \quad (59)$$

Now by definition we have set

$$\delta_{H+1,k}(s_{H+1,k}) = 0 \quad (60)$$

and using our answer to Question 2 we can re-write the previous equation as

$$\delta_{1,k}(s_{1,k}) \leq \sum_{h=1}^H Q_{hk}(s_{hk}, a_{hk}) - r(s_{hk}, a_{hk}) - \mathbb{E}_p[V_{h+1,k}(s')] + m_{hk} \quad (61)$$

which is the desired form.

- 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

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



=====Answer:=====

From the definition we have

$$R(T) = \sum_{k=1}^K V_1^*(s_{1,k}) - V_1^{\pi_k}(s_{1,k}) = \sum_{k=1}^K \delta_{1k}(s_{1,k}) \quad (62)$$

Using our previous result we get

$$R(T) \leq \sum_{k,h} m_{h,k} + \sum_{k,h} Q_{hk}(s_{hk}, a_{hk}) - r(s_{hk}, a_{hk}) - \mathbb{E}_p[V_{h+1,k}(s')] \quad (63)$$

$$\leq \sum_{k,h} m_{h,k} + \sum_{k,h} (\hat{r}_{hk}(s_{hk}, a_{hk}) - r(s_{hk}, a_{hk})) + \mathbb{E}_{\hat{p}-p}[V_{h+1,k}(s')] + b_{hk}(s_{hk}, a_{hk}) \quad (64)$$

from Question 2 we know that with probability  $\mathbb{P}(\mathcal{E}) = 1 - \delta/2$  then

$$(\hat{r}_{hk}(s_{hk}, a_{hk}) - r(s_{hk}, a_{hk})) + \mathbb{E}_{\hat{p}-p}[V_{h+1,k}(s')] \leq b_{hk}(s_{hk}, a_{hk}) \quad (65)$$

then from Azuma-Hoeffding with probability at least  $1 - \delta/2$  we get

$$\sum_{k,h} m_{hk} \leq 2H\sqrt{KH \log(2/\delta)} \quad (66)$$

hence putting the two inequalities together we have that with probability at least  $1 - \delta$  we have that

$$R(T) \leq 2H\sqrt{KH \log(2/\delta)} + 2 \sum_{kh} b_{hk}(s_{hk}, a_{hk}) \quad (67)$$

- 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  $R(T) \lesssim H^2 S \sqrt{AK}$

=====Answer:=====

We have that

$$\sum_{h=1}^H \sum_{s,a} \sqrt{N_{hK}(s, a)} \leq \sum_{h=1}^H \sqrt{SAN_{hK}(s, a)} \leq \sum_{h=1}^H \sqrt{SAK} = H\sqrt{SAK} \quad (68)$$

Now we want to bound the bonus. We have that

$$H\beta_{hk}^p(s, a) = H\sqrt{2 \log \left( \frac{4(2^S - 2)AHKS}{\delta} \right)} \sqrt{\frac{1}{N_{hk}(s, a)}} \quad \text{and} \quad \beta_{hk}^r = \sqrt{\frac{H}{2} \log \left( \frac{8AHKS}{\delta} \right)} \quad (69)$$

hence applying the inequalities given above we get

$$\sum_{kh} b_{kh}(s_{hk}, a_{hk}) \leq \sqrt{\frac{1}{2} \log \left( \frac{8AHKS}{\delta} \right)} HK\sqrt{H} + H\sqrt{2 \log \left( \frac{4(2^S - 2)AHKS}{\delta} \right)} (H^2 S^2 A + 2H\sqrt{SAK}) \quad (70)$$

$$\leq HK\sqrt{H} \sqrt{\log \frac{AHKS}{\delta} + \mathcal{O}(1)} + H\sqrt{S} \sqrt{\log \frac{AHKS}{\delta} + \mathcal{O}(1)} (H^2 S^2 A + 2H\sqrt{SAK}) \quad (71)$$

$$= H\sqrt{\log \frac{AHKS}{\delta} + \mathcal{O}(1)} \left( K\sqrt{H} + H^2 S^{5/2} A + 2HS\sqrt{AK} \right) \quad (72)$$

```

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

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

  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

The leading order of  $R(T)$  then seems to be  $H^3 S^{5/2} A$  which is not the desired behavior, however if there was no  $H^2 S^2 A$  term in the originally stated inequality we would get

$$R(T) \lesssim H^2 S \sqrt{AK} \quad (73)$$

which is the desired behavior.

## A Weissmain inequality

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

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

## 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 M  nard, Matteo Pirodda, 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.