

# ① 1. Best Arm Identification

• Let's compute  $U(t, \delta)$ :

We want  $U(t, \delta)$  such that:

$$\mathbb{P}\left(\bigcap_{t=1}^{+\infty} \{|\hat{\mu}_{i,t} - \mu_i| \leq U(t, \delta)\}\right) \geq 1 - \delta$$

$$\text{Thus } \mathbb{P}\left(\bigcup_{t=1}^{+\infty} \{|\hat{\mu}_{i,t} - \mu_i| > U(t, \delta)\}\right) \leq \delta$$

Let's note that

$$\begin{aligned} & \mathbb{P}\left(\bigcup_{t=1}^{+\infty} \{|\hat{\mu}_{i,t} - \mu_i| > U(t, \delta)\}\right) \\ & \leq \sum_{t=1}^{+\infty} \mathbb{P}\left(|\hat{\mu}_{i,t} - \mu_i| > U(t, \delta)\right) \text{ by Union bound} \\ & = \sum_{t=1}^{+\infty} \mathbb{P}\left(\left|\sum_{j=1}^t X_{i,j} - t\mu_i\right| > tU(t, \delta)\right) \\ & \leq \sum_{t=1}^{+\infty} 2 \exp\left(-\frac{2U(t, \delta)^2 t^2}{t}\right) \text{ by Hoeffding.} \end{aligned}$$

If we set  $2 \exp\left(-\frac{2U(t, \delta)^2 t^2}{t}\right) = \frac{\delta}{2t^2}$ , we will get:

$$\sum_{t=1}^{+\infty} 2 \exp\left(-\frac{2U(t, \delta)^2 t^2}{t}\right) = \sum_{t=1}^{+\infty} \frac{\delta}{2t^2} = \frac{\delta}{2} \frac{\pi^2}{6} < \frac{\delta}{2} \cdot 2 = \delta.$$

From that we can deduce  $U(t, \delta)$ :

$$2 \exp\left(-\frac{2U(t, \delta)^2 t^2}{t}\right) = \frac{\delta}{2t^2}$$

$$\Leftrightarrow 2U(t, \delta)^2 t = -\ln\left(\frac{\delta}{4t^2}\right)$$

$$\Leftrightarrow U(t, \delta) = \sqrt{\frac{1}{2t} \ln\left(\frac{4t^2}{\delta}\right)}$$

② Let's show that  $P(\mathcal{E}) \leq \delta$  for a certain  $\delta'$ :

$$\begin{aligned} P(\mathcal{E}) &= P\left(\bigcup_{i=1}^k \bigcup_{t=1}^{\infty} \{|\hat{\mu}_{i,t} - \mu_i| > \mathcal{U}(t, \delta')\}\right) \\ &\leq \sum_{i=1}^k P\left(\bigcup_{t=1}^{\infty} \{|\hat{\mu}_{i,t} - \mu_i| > \mathcal{U}(t, \delta')\}\right) \\ &\leq \sum_{i=1}^k \delta' \quad \text{by definition of } \mathcal{U}(t, \delta') \\ &= k \delta' \end{aligned}$$

So, if we set  $\delta' = \frac{\delta}{k}$  we get  $P(\mathcal{E}) \leq \delta$ .

• Let's show that  $P("i^*" \text{ remains in the active set } S") \geq 1 - \delta$ .

$i^*$  is eliminated from  $S$  if:

$$\exists j \in S, \hat{\mu}_{j,t} - \mathcal{U}(t, \delta') \geq \hat{\mu}_{i^*,t} + \mathcal{U}(t, \delta'). \quad (*)$$

If  $\neg \mathcal{E}$  holds,  $\forall j \in S$  we have:

$$\begin{aligned} |\hat{\mu}_{j,t} - \mu_j| &\leq \mathcal{U}(t, \delta'). \\ \neg \mathcal{E} &= \bigcap_{j=1}^k \bigcap_{t=1}^{\infty} \{|\hat{\mu}_{j,t} - \mu_j| \leq \mathcal{U}(t, \delta')\}. \end{aligned}$$

So we have:

$$\begin{aligned} |\hat{\mu}_{j,t} - \mu_j| \leq \mathcal{U}(t, \delta') &\Leftrightarrow \hat{\mu}_{j,t} - \mathcal{U}(t, \delta') \leq \mu_j \leq \hat{\mu}_{j,t} + \mathcal{U}(t, \delta'). \\ &\quad \forall j \in S. \end{aligned}$$

From that and (\*), we have.

$$\mu_j \geq \mu_{i^*}.$$



③

By definition of  $i^*$ ,  $\mu_j \geq \mu_{i^*}$  is impossible  
hence  $i^*$  can't be eliminated from  $S$  when  $\neg E$  holds.  
Moreover,  $P(\neg E) \geq 1 - \delta$  and thus  
 $P(i^* \text{ remains in the active set } S'') \geq 1 - \delta$ .

• Let's show that under event  $\neg E$ , an arm  $j \neq i^*$  will be removed from the active set when  
 $\Delta_j \geq C_1 u(t, s')$  for  $C_1 = 4$ .

Indeed, we have:

$$\Delta_j \geq 4 u(t, s') \Leftrightarrow \mu_{i^*} - 2 u(t, s') \geq \mu_j + 2 u(t, s')$$

Moreover, under  $\neg E$ , we have:

$$|\hat{\mu}_{i, \neg E} - \mu_i| \leq u(t, s') \quad \forall i$$

$$\Leftrightarrow \mu_i - u(t, s') \leq \hat{\mu}_{i, \neg E} \leq \mu_i + u(t, s') \quad \forall i.$$

From that we can deduce that:

$$\mu_{i^*} - 2 u(t, s') \geq \mu_j + 2 u(t, s')$$

$$\Leftrightarrow \hat{\mu}_{i^*, \neg E} - u(t, s') \geq \hat{\mu}_{j, \neg E} + u(t, s')$$

And thus  $j$  is eliminated from the active set.

## ④ 2. Regret Minimization in RL

• We want:

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

$$\Leftrightarrow P(\exists k, h, s, a: |\hat{r}_{kh}(s, a) - r_h(s, a)| > \beta_{kh}^r(s, a) \vee \|\hat{p}_{kh}(\cdot | s, a) - p_h(\cdot | s, a)\|_1 > \beta_{kh}^p(s, a)) \leq \frac{\delta}{2}.$$

Let's note that:

$$P(\exists k, h, s, a: |\hat{r}_{kh}(s, a) - r_h(s, a)| > \beta_{kh}^r(s, a) \vee \|\hat{p}_{kh}(\cdot | s, a) - p_h(\cdot | s, a)\|_1 > \beta_{kh}^p(s, a))$$

$$\leq \sum_{k, h, s, a} P(|\hat{r}_{kh}(s, a) - r_h(s, a)| > \beta_{kh}^r(s, a) + P(\|\hat{p}_{kh}(\cdot | s, a) - p_h(\cdot | s, a)\|_1 > \beta_{kh}^p(s, a))$$

$$\leq \sum_{k, h, s, a} 2 \exp\left(-2 \frac{\beta_{kh}^r(s, a)^2 N_{kh}(s, a)}{N_{kh}(s, a)}\right) + (2^S - 2) \exp\left(-\frac{1}{2} N_{kh}(s, a) \beta_{kh}^p(s, a)^2\right)$$

So we want:

$$2 \exp\left(-2 \frac{\beta_{kh}^r(s, a)^2 N_{kh}(s, a)}{N_{kh}(s, a)}\right) = \frac{\delta}{4HKSA}$$

$$\text{and } (2^S - 2) \exp\left(-\frac{1}{2} N_{kh}(s, a) \beta_{kh}^p(s, a)^2\right) = \frac{\delta}{4HKSA}.$$



$$(5) \quad 2 \exp(-2 \beta_{hh}^R(s, a) N_{hh}(s, a)) = \frac{\delta}{4HKSA}$$

$$\Leftrightarrow \beta_{hh}^R(s, a) = \frac{1}{2N_{hh}(s, a)} \ln \left( \frac{8HKSA}{\delta} \right)$$

$$\Leftrightarrow \beta_{hh}^R(s, a) = \sqrt{\frac{1}{2N_{hh}(s, a)} \ln \left( \frac{8HKSA}{\delta} \right)}$$

$$\frac{(2^S - 2) \exp\left(-\frac{1}{2} N_{hh}(s, a) \beta_{hh}^P(s, a)^2\right)}{\delta} = \frac{\delta}{4HKSA}$$

$$\Leftrightarrow \beta_{hh}^P(s, a)^2 = \frac{2}{N_{hh}(s, a)} \ln \left( \frac{(2^S - 2) 4HKSA}{\delta} \right)$$

$$\Leftrightarrow \beta_{hh}^P(s, a) = \sqrt{\frac{2}{N_{hh}(s, a)} \ln \left( \frac{(2^S - 2) 4HKSA}{\delta} \right)}$$

So if we set  $\beta_{hh}^R(s, a) = \sqrt{\frac{1}{2N_{hh}(s, a)} \ln \left( \frac{8HKSA}{\delta} \right)}$

and  $\beta_{hh}^P(s, a) = \sqrt{\frac{2}{N_{hh}(s, a)} \ln \left( \frac{(2^S - 2) 4HKSA}{\delta} \right)}$

we get  $\mathbb{P}(\neg \varepsilon) \leq \frac{\delta}{2}$ .

⑥. Let  $b_{h,k}(s,a) := \beta_{h,k}^{\pi}(s,a) + H \beta_{h,k}^P(s,a)$ .

Let's note that  $\hat{\pi}_{H,k}(s,a) + b_{H,k}(s,a) \geq r_{H,k}(s,a)$

under event  $E$ .

Under event  $E$ , we have for  $(s,a) \in S \times A$ :

$$|\hat{\pi}_{H,k}(s,a) - r_{H,k}(s,a)| \leq \beta_{H,k}^{\pi}(s,a).$$

$$\text{Hence } r_{H,k}(s,a) \leq \beta_{H,k}^{\pi}(s,a) + \hat{\pi}_{H,k}(s,a).$$

$$\begin{aligned} \text{So } \hat{\pi}_{H,k}(s,a) + b_{H,k}(s,a) &= \hat{\pi}_{H,k}(s,a) + \beta_{H,k}^{\pi}(s,a) + \underbrace{H \beta_{H,k}^P(s,a)}_{\geq 0} \\ &\geq r_{H,k}(s,a). \end{aligned}$$

$$\text{and thus } Q_{H,k}(s,a) = \hat{\pi}_{H,k}(s,a) + b_{H,k}(s,a)$$

$$\geq r_{H,k}(s,a) = Q_k^*(s,a) \forall s,a.$$

Let's suppose that  $Q_{h,k}(s,a) \geq Q_k^*(s,a) \forall s,a$   
for  $h \leq H-1$

$$\text{Let's show that } Q_{h-1,k}(s,a) \geq Q_{h-1}^*(s,a) \forall s,a.$$

Let  $(s,a) \in S \times A$ .

$$\begin{aligned} Q_{h-1,k}(s,a) &= \hat{\pi}_{h-1,k}(s,a) + b_{h-1,k}(s,a) \\ &\quad + \sum_{s'} \hat{P}_{h-1,k}(s'|s,a) V_{h,k}(s'). \end{aligned}$$



(7)

$$r_{h-1,h}(s,a) + \sum_{s'} p_{h-1,h}(s,a) V_{h,h}(s')$$

$$\text{since } b_{h-1,h}(s,a) = \beta_{h-1,h}^r(s,a) + H \beta_{h-1,h}^p(s,a).$$

and  $\epsilon$  holds so:

$$|\hat{r}_{h-1,h}(s,a) - r_{h-1,h}(s,a)| \leq \beta_{h-1,h}^r(s,a)$$

$$\Leftrightarrow r_{h-1,h}(s,a) \leq \beta_{h-1,h}^r(s,a) + \hat{r}_{h-1,h}(s,a).$$

and

$$\sum_{s'} (\hat{p}_{h-1,h}(s'|s,a) - p_{h-1,h}(s'|s,a)) \underbrace{V_{h,h}(s')}_{\leq H}$$

$$\leq H \|\hat{p}_{h-1,h}(\cdot|s,a) - p_{h-1,h}(\cdot|s,a)\|_1$$

$$\leq H \beta_{h-1,h}^p(s,a)$$

$$\Leftrightarrow \sum_{s'} p_{h-1,h}(s'|s,a) V_{h,h}(s')$$

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

So we have:

$$Q_{h-1,h}(s,a) \geq r_{h-1,h}(s,a) + \mathbb{E}_{s' \sim p_{h-1,h}} [V_{h,h}(s')]$$

if  $V_{h,h}(s') = H$  for  $s' \in S$ , we have

$$V_{h,h}(s') = H \geq \cancel{Q_h^*} \max_a Q_h^*(s',a).$$

if  $V_{h,h}(s') = \max_a Q_{h,h}(s',a)$ , we have.

$$(8) \quad V_{h,h}(s') = \max_a Q_{h,h}(s',a) \geq \max_a Q_h^*(s',a) \quad \text{by assumption.}$$

by assumption.

$$\text{so } \forall s' \in S, V_{h,h}(s') \geq \max_a Q_h^*(s',a).$$

$$\begin{aligned} \text{Thus } r_{h-1,h}(s,a) + \mathbb{E}_{s' \sim p(\cdot|s,a)} [V_{h,h}(s')] \\ \geq r_{h-1,h}(s,a) + \mathbb{E}_{s' \sim p(\cdot|s,a)} [\max_a Q_h^*(s',a)] \\ = Q_{h-1}^*(s,a) \end{aligned}$$

$$\text{Finally, } Q_{h-1,h}(s,a) \geq Q_{h-1}^*(s,a) \quad \forall (s,a).$$

$$\text{and thus } Q_{h,h}(s,a) \geq Q_h^*(s,a) \quad \forall (s,a) \quad \forall h=1, \dots, H$$

$$\begin{aligned} \bullet 1. \text{ Let's show that } V_h^{\pi_h}(s_{hh}) &= r(s_{hh}, a_{hh}) \\ &+ \mathbb{E}_{s' \sim p(\cdot|s_{hh}, a_{hh})} [V_{h+1,h}(s')] \\ &- \delta_{h+1,h}(s_{hh}, a_{hh}) \\ &- m_{h,h}. \end{aligned}$$

We have:

$$\begin{aligned} V_h^{\pi_h}(s_{hh}) &= Q_h^{\pi_h}(s_{hh}, \pi_{h,h}(s_{hh})) \\ &= Q_h^{\pi_h}(s_{hh}, a_{hh}) = r(s_{hh}, a_{hh}) + \mathbb{E}_{s' \sim p(\cdot|s_{hh}, a_{hh})} [V_{h+1}^{\pi_h}(s')] \\ &= r(s_{hh}, a_{hh}) + \mathbb{E}_{s' \sim p(\cdot|s_{hh}, a_{hh})} [V_{h+1,h}(s') - \delta_{h+1,h}(s')]. \end{aligned}$$



$$\begin{aligned}
& \textcircled{9} = r(s_{hh}, a_{hh}) + \mathbb{E}_{s' \sim p(\cdot | s_{hh}, a_{hh})} [V_{h+1, h}(s')] - \mathbb{E}_{s' \sim p(\cdot | s_{hh}, a_{hh})} [\delta_{h+1, h}(s')] \\
& + \delta_{h+1, h}(s_{h+1, h}) - \delta_{h+1, h}(s_{h+1, h}) \\
& = r(s_{hh}, a_{hh}) + \mathbb{E}_{s' \sim p(\cdot | s_{hh}, a_{hh})} [V_{h+1, h}(s')] - \delta_{h+1, h}(s_{h+1, h}) \\
& - m_{hh}.
\end{aligned}$$

2. Let's show that  $V_{h, h}(s_{hh}) \leq Q_{h, h}(s_{hh}, a_{hh})$

We have:

$$\begin{aligned}
V_{h, h}(s_{hh}) &= \min \{ H, \max_a Q_{h, h}(s_{hh}, a_{hh}) \} \\
&\leq \max_a Q_{h, h}(s_{hh}, a) \\
&= Q_{h, h}(s_{hh}, a_{hh}) \text{ because } a_{hh} = \arg \max_a Q_{h, h}(s_{hh}, a)
\end{aligned}$$

We do have  $V_{h, h}(s_{hh}) \leq Q_{h, h}(s_{hh}, a_{hh})$

3. Let's prove eq. 1.

Let's note that for  $h=1, \dots, H$ , we have:

$$\begin{aligned}
\delta_{hh}(s) &= V_{hh}(s) - V_h^{\pi^h}(s) \leq Q_{hh}(s_{hh}, a_{hh}) \\
&\quad - r(s_{hh}, a_{hh}) - \mathbb{E}_{s' \sim p(\cdot | s_{hh}, a_{hh})} [V_{h+1, h}(s')] \\
&\quad + \delta_{h+1, h}(s_{h+1, h}) + m_{hh}.
\end{aligned}$$

by q. 1 and q. 2.

⑩

Hence:

$$\begin{aligned} \delta_{1,h}(s_{1,h}) &\leq Q_{1,h}(s_{1,h}, a_{1,h}) - r(s_{1,h}, a_{1,h}) \\ &\quad - \mathbb{E}_{s' \sim p(\cdot | s_{1,h}, a_{1,h})} [V_{2,h}(s')] + \delta_{2,h}(s_{2,h}) \\ &\quad + m_{1,h} \end{aligned}$$

$$\leq Q_{1,h}(s_{1,h}, a_{1,h}) - r(s_{1,h}, a_{1,h})$$

$$- \mathbb{E}_{s' \sim p(\cdot | s_{1,h}, a_{1,h})} [V_{2,h}(s')] + m_{1,h}$$

$$+ Q_{2,h}(s_{2,h}, a_{2,h}) - r(s_{2,h}, a_{2,h})$$

$$- \mathbb{E}_{s' \sim p(\cdot | s_{2,h}, a_{2,h})} [V_{3,h}(s')] + \delta_{3,h}(s_{3,h}) + m_{2,h}$$

$\ll \dots$

$$\leq Q_{1,h}(s_{1,h}, a_{1,h}) - r(s_{1,h}, a_{1,h})$$

$$- \mathbb{E}[V_{2,h}(s')] + m_{1,h} + Q_{2,h}(s_{2,h}, a_{2,h})$$

$$- r(s_{2,h}, a_{2,h}) - \mathbb{E}_{s' \sim p(\cdot | s_{2,h}, a_{2,h})} [V_{3,h}(s')] + m_{2,h}$$

$$+ \dots + Q_{H,h}(s_{H,h}, a_{H,h}) - r(s_{H,h}, a_{H,h})$$

$$- \mathbb{E}_{s' \sim p(\cdot | s_{H,h}, a_{H,h})} [V_{H+1,h}(s')] + \underbrace{\delta_{H+1,h}(s_{H+1,h}) + m_{H,h}}_{=0}$$

$$= V_{H+1,h}(s_{H+1,h}) - V_{H+1,h}^{\pi,h}(s_{H+1,h}) = 0$$

$$= \sum_{h=1}^H Q_{h,h}(s_{h,h}, a_{h,h}) - r(s_{h,h}, a_{h,h})$$

$$- \mathbb{E}_{s' \sim p(\cdot | s_{h,h}, a_{h,h})} [V_{h+1,h}(s')] + m_{h,h}$$



(11) Finally,

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

• Let's suppose that  $\epsilon$  and  $\left\{ \sum_{k,h} m_{1k,h} \leq 2H\sqrt{KH \log\left(\frac{2}{\delta}\right)} \right\}$  hold.

Let's show that  $R(T) \leq 2 \sum_{k,h} b_{1k,h}(s_{1k}, a_{1k}) + 2H\sqrt{KH \log\left(\frac{2}{\delta}\right)}$ ,

with probability  $1 - \delta$ :

$$R(T) = \sum_{k=1}^K V_1^*(s_{1k}) - V_1^{\pi_k}(s_{1k}) \\ = \sum_{k=1}^K V_1^*(s_{1k}) - V_{1k}(s_{1k}) + \underbrace{V_{1k}(s_{1k}) - V_1^{\pi_k}(s_{1k})}_{=\delta_{1k}(s_{1k})}.$$

Let's note that  $V_{1k}(s_{1k}) \geq \max_a Q_{1k}^*(s_{1k}, a)$ .

If  $V_{1k}(s_{1k}) = H$ , then  $V_{1k}(s_{1k}) = H \geq \max_a Q_{1k}^*(s_{1k}, a)$ .

If  $V_{1k}(s_{1k}) = \max_a Q_{1k}(s_{1k}, a) \geq \max_a Q_{1k}^*(s_{1k}, a)$

because  $Q_{1k}$  is optimistic under  $\epsilon$ .

Thus,

$$R(T) \leq \sum_{k=1}^K \underbrace{V_1^*(s_{1k}) - \max_a Q_{1k}^*(s_{1k}, a)}_{=V_{1k}^*(s_{1k})} + \delta_{1k}(s_{1k}) \\ = 0.$$

$$\textcircled{12} \leq \sum_{k,h} Q_{kh}(s_{kh}, a_{kh}) - r(s_{kh}, a_{kh}) - \mathbb{E}_{Y \sim P(\cdot | s_{kh}, a_{kh})} [V_{k+1,h}(Y)] + m_{kh} \quad \text{by equation (1).}$$

$$= \sum_{k,h} \hat{\pi}_{kh}(s_{kh}, a_{kh}) + b_{kh}(s_{kh}, a_{kh}) + \left( \sum_{s'} \hat{p}_{kh}(s' | s_{kh}, a_{kh}) \cdot V_{k+1,h}(s') \right) - \mathbb{E}_{Y \sim P(\cdot | s_{kh}, a_{kh})} [V_{k+1,h}(Y)] + m_{kh} - r(s_{kh}, a_{kh})$$

$$= \sum_{k,h} \hat{\pi}_{kh}(s_{kh}, a_{kh}) - r(s_{kh}, a_{kh}) + \sum_{s'} (\hat{p}_{kh}(s' | s_{kh}, a_{kh}) - p(s' | s_{kh}, a_{kh})) V_{k+1,h}(s') + b_{kh}(s_{kh}, a_{kh}) + \sum_{k,h} m_{kh}$$

$$\leq \sum_{k,h} \tilde{\beta}_{kh}(s_{kh}, a_{kh}) + H \beta_{kh}^p(s_{kh}, a_{kh})$$

$$+ b_{kh}(s_{kh}, a_{kh}) + 2H \sqrt{KH \log\left(\frac{2}{\delta}\right)}$$

because we are under  $\varepsilon$  and

$$\left\{ \sum_{k,h} m_{kh} \leq 2H \sqrt{KH \log\left(\frac{2}{\delta}\right)} \right\}$$

$$= 2 \sum_{k,h} b_{kh}(s_{kh}, a_{kh}) + 2H \sqrt{KH \log\left(\frac{2}{\delta}\right)}$$

by definition of the bonus.

$$\text{Finally, } R(T) \leq 2 \sum_{k,h} b_{kh}(s_{kh}, a_{kh}) + 2H \sqrt{KH \log\left(\frac{2}{\delta}\right)}$$

$$\text{Moreover, } P(\neg \varepsilon \cap \left\{ \sum_{k,h} m_{kh} > 2H \sqrt{KH \log\left(\frac{2}{\delta}\right)} \right\}) \leq \frac{\delta}{2} + \frac{\delta}{2} = \delta$$

$$\text{Thus } P\left(\varepsilon \cap \left\{ \sum_{k,h} m_{kh} \leq 2H \sqrt{KH \log\left(\frac{2}{\delta}\right)} \right\}\right) \geq 1 - \delta$$



(13)

$$R(T) \leq 2 \sum_{h,k} b_{hk} (s_{hk}, a_{hk}) + 2H \sqrt{KH \log\left(\frac{2}{\delta}\right)}$$

with probability at least  $1 - \delta$ .

.