CHAUVIN
Paul

I. Best Arm identification :

1). Let $\varepsilon = \bigcup_{i=1}^{m} \bigcup_{t=1}^{\infty} \{|\hat{\mu}_{i,t} - \mu_i| > u(t,\delta')\}$

Then, $P(\varepsilon) \leq \sum_{i=1}^{m} P\left(\bigcup_{t=1}^{\infty} |\hat{\mu}_{i,t} - \mu_i| > u(t,\delta')\right)$

Using the anytime confidence bound, we have :

$P(\varepsilon) \leq \sum_{i=1}^{m} P\left(\bigcup_{t=1}^{\infty} |\hat{\mu}_{i,t} - \mu_i| > u(t,\delta')\right) \leq \sum_{i=1}^{m} \delta'$

$$\leq m\,\delta'$$
$$\leq \delta$$

with $\boxed{\delta' = \dfrac{\delta}{m}}$

2). With the arm elimination condition for the optimal arm :

$$\exists j \;/\; \hat{\mu}_{j,t} - u(t,\delta') \geq \hat{\mu}_t^* + u(t,\delta')$$

In order to delete the optimal arm, $\mu^*$ should be outside the confidence interval : $|\hat{\mu}_t^* - \mu^*| > u(t,\delta')$

$\Rightarrow P(|\hat{\mu}_t^* - \mu^*| > u(t,\delta'), \text{ for one } t) \leq P(\varepsilon)$

$\leq \delta$

$\Rightarrow \forall t, \; P(|\hat{\mu}_t^* - \mu^*| \leq u(t,\delta')) \geq 1-\delta$

## I.3

Let $\hat{\mu}_t^a$ be the estimated reward of the arm with the largest expected reward $\mu^a$

Under $\neg \mathcal{E}$ :  $\hat{\mu}_t^a \geq \mu^a - u(t, \delta')$

$$\hat{\mu}_{i,t} \leq \mu_i + u(t, \delta')$$

and with elimination condition

$$\hat{\mu}_t^a - u(t, \delta') \geq \hat{\mu}_{it} + u(t, \delta')$$

$\Rightarrow$ arm $i$ will be deleted if

$$\mu^a - 2u(t, \delta') \geq \mu_i + 2u(t, \delta')$$

$$\Rightarrow \quad \underbrace{\mu^a - \mu_i}_{D_i} \geq \underbrace{4\, u(t, \delta')}_{C_1}$$

$$u(t, \delta') = \sqrt{\frac{1}{2t} \log\left(4t^2/\delta'\right)}$$

$$u(t, \delta') = \sqrt{\frac{1}{2t} \log\left(\frac{4mt^2}{\delta}\right)} \qquad \text{because } \delta' = \frac{\delta}{m}$$

Hence,

$$D_i^2 \geq \frac{16}{2t} \log\left(\frac{4mt^2}{\delta}\right)$$

$$\Rightarrow \quad \frac{D_i^2}{16} t \geq \log\left(\frac{4mt}{\delta}\right)$$

$$\Rightarrow \quad at \geq \log(bt)$$

with $\boxed{a = \dfrac{D_i^2}{16} \quad ; \quad b = \dfrac{4m}{\delta}}$

$$\Rightarrow \boxed{t \geq \frac{1 + \sqrt{2u}}{a} + u} \quad , \boxed{u = \log\left(\frac{b}{a}\right) - 1}$$

using the footnote.

CHAUVIN I. 4). the arm will be removed after Paul sampling each sub-optimal arm.

The sample complexity is then:

$$O\left(\sum_{i\neq i^*} \frac{\log(b/a)}{a}\right) = O\left(\sum_{i\neq i^*} \frac{\log\left(\frac{M}{\delta D_i^2}\right)}{D_i^2}\right)$$

I. 5).

If multiple best arm exist, the algorithm would never stop as it would not be able to find "bad" arms and $S$ would never be equal to one, so the algorithm would not stop.


## II. Regret Minimization.

II. 1). For fixed $s, a, h, K$ we have:

- Hoeffding's inequality:

$$P(\neg \mathcal{E}_r) = P\left(|\hat{r}_{HK}(s,a) - r_{HK}(s,a)| \geq B_{h,k}^r(s,a)\right)$$

$$\leq \underbrace{2\exp\left(-2N_{h,k}(s,a) B_{h,k}^r(s,a)^2\right)}_{\delta_r}$$

$$\Rightarrow 2N_{h,k}(s,a) B_{hk}^r(s,a)^2 = \log\left(\frac{2}{\delta_r}\right)$$

$$\Rightarrow \boxed{B_{hk}^r(s,a) = \sqrt{\frac{\log\left(\frac{2}{\delta_r}\right)}{2N_{h,k}(s,a)}}}$$

- Weissman inequality:

$$P(\neg \mathcal{E}_p) = P\left(\| \hat{P}^{1}_{hk}(\cdot | s,a) - P_h(\cdot | s,a)\|_1 \geq B^P_{hk}(s,a)\right)$$

$$\leq (2^s - 2) \exp\left(- \frac{N_{hk}(s,a)\, B^P_{hk}(s,a)^2}{2}\right)$$

$$N_{hk}(s,a)\, B^P_{hk}(s,a)^2 = 2 \log\left(\frac{2^s - 2}{\delta_p}\right)$$

$$\Rightarrow \quad B^P_{hk}(s,a)^2 = \sqrt{\frac{2}{N_{hk}(s,a)} \log\left(\frac{2^s - 2}{\delta_p}\right)}$$

- Both inequalities give:

$$P(\neg \mathcal{E}_{s,a,k,h}) \leq P(\neg \mathcal{E}_r) + P(\neg \mathcal{E}_p)$$

we set $\delta_r = \delta_a = \frac{\delta'}{2}$

$$P(\neg \mathcal{E}_{s,a,k,h}) \leq \delta'$$

the bound for any $s,a,k,h$.

Then, $P(\mathcal{E}) = 1 - P(\neg \mathcal{E}) = 1 - P\left(\bigcup_{s,a,k,k} \neg \mathcal{E}_{s,a,k,k}\right)$

By Union bound:

$$1 - P\left(\bigcup_{s,a,k,k} \neg \mathcal{E}_{s,a,k,k}\right) \geq 1 - \sum_{s,a,k,k} P(\neg \mathcal{E}_{s,a,k,k})$$

we want $P(\mathcal{E}) \geq 1 - \delta/2$

$$\sum_{s,a,k,k} P(\neg \mathcal{E}_{s,a,k,k}) \leq \frac{\delta}{2}$$

SAHK $\delta' = \frac{\delta}{2} \Rightarrow \delta' = \frac{\delta}{2 SAHK}$

4/9

Hence, the confidence bounds are:

$$B_{hk}^{r}(s,a) = \sqrt{\frac{\log\left(\frac{8SAHK}{\delta}\right)}{2 N_{hk}(s,a)}}$$

$$B_{hk}^{P}(s,a) = \sqrt{\frac{2}{N_{hk}(s,a)} \log\left(\frac{(2^{S}-2)\,4SAHK}{\delta}\right)}$$

## II.2)

### Base case : $h = H$.

$$Q_{H,k}(s,a) = \hat{r}_{H,k}(s,a) + b_{H,k}(s,a)$$

$$Q_{H}^{\star}(s,a) = r_{H,k}(s,a)$$

we are under $\mathcal{E}$ event, hence :

$$\hat{r}_{H,k}(s,a) \geqslant r_{H,k}(s,a) - B_{H,k}^{r}(s,a)$$

then : $Q_{H,k}(s,a) \geqslant r_{Hk}(s,a) + b_{H,k}(s,a) - B_{H,k}^{r}(s,a)$

with bonus : $b_{H,k}(s,a) \geqslant B_{H,k}^{r}(s,a)$ , base case

is true.

### Inductive step.

Assume $\quad Q_{h,k}(s,a) \geqslant Q_{h}^{\star}(s,a)$

let's prove $\quad Q_{h-1,k}(s,a) \geqslant Q_{h-1}^{\star}(s,a)$

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

$$= \hat{r}_{h-1,k}(s,a) + b_{h-1,k}(s,a) + \sum_{s'} \hat{P}_{h-1,k}(s'|s,a) \min\left(H, \max_{a} Q_{h,k}(s,a)\right)$$

5/9

$$Q^{R}_{h-1,k}(s,a) = r_{h-1,k}(s,a) + \sum_{s'} P_{h-1,k}(s'|s,a) \max_a Q^*_{h,k}(s,a)$$

$$\Rightarrow Q_{h-1,k}(s,a) - Q^{R}_{h-1,k}(s,a) = \hat{r}_{h-1,k}(s,a) + b_{h-1,k}(s,a)$$

$$- r_{h-1,k}(s,a) + \sum_{s'} \hat{P}_{h-1,k}(s'|s,a) \min\left(H, \max_a Q(s,a)\right)$$

$$- P_{h-1,k}(s'|s,a) \max_a Q^{R}_{h,k}(s,a)$$

$$\Rightarrow Q_{h-1,k} - Q^{R}_{h-1,k} \geqslant \hat{r}_{h-1,k}(s,a) + b_{h-1,k}(s,a) - r_{h-1,k}(s,a)$$

$$+ \sum_{s'} \min\left(H, \max_a Q\right) \left(\hat{P}_{h-1,k} - P_{h-1,k}\right)$$

$$\geqslant \hat{r}_{h-1,k} + b_{h-1,k} - r_{h-1,k}$$

$$- \sum_{s'} \min\left(H, \max_a Q\right) |\hat{P} - P|$$

$$\geqslant \hat{r}_{h-1,k} + b_{h-1,k} - r_{h-1,k}$$

$$- H \sum_{s'} |\hat{P} - P|$$

$$\geqslant \hat{r}_{h-1,k} + b_{h-1,k} - r_{h-1,k} - H B^{P}_{h-1,k}(s,a)$$

because confidence
intervals holds as we're
under event $\Sigma$.

$$\geqslant - B^{r}_{h-1,k}(s,a) - H B^{P}_{h-1,k}(s,a) + b_{h-1,k}(s,a)$$

$$\geqslant 0$$

$$\Longleftrightarrow \quad b_{h-1,k}(s,a) \geqslant B^{r}_{h-1}(s,a) + H B^{P}_{h-1,k}$$

## II. 3.1)

$$V_h^{T_k}(s_{hk}) = r(s_{hh}, a_{hh}) + \sum_{s'} p(s' | (s, a))\, V_{h+1,k}^{T_k}(s^i)$$

$$= r(s_{hh}, a_{hh}) + \sum_{s'} p(s'|s,a)\left(V_{h+1}(s') - \delta_{h+1,k}(s^i)\right)$$

$$= r(s_{hh}, a_{hh}) + E_p\left[V_{h+1,k}(s')\right] - E_p\left[\delta_{h+1,k}(s')\right]$$

$$= r(s_{hh}, a_{hh}) + E_p\left[V_{h+1,k}(s')\right] - \delta_{h+1,k}(s_{h+1,k}) - m_{h,k}$$

## II. 3.2)

$$V_{hk}(s_{hh}) = \min\left\{ H,\ \max_a Q_{h,k}(s, a) \right.$$

$$\leq Q_{h,k}(s_{hh}, a_{hh})$$

## III. 3.3)

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

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

$$+ \delta_{2k}(s_{2k}) + m_{1k}$$

$$\leq \dots \qquad \text{(we replace } \delta_{2k}(s_{2k}) \text{ in the same way}$$

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

$$+ m_{hk}$$

# II. 4).

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

$$\leq \sum_{k=1}^{K} \hat{V}_1(s_{1,k}) - V_1^{\pi_k}(s_{1,k}) = \sum_{k=1}^{K} \delta_{1k}(s_{1k})$$

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

$$- E_{y \sim p}\left[V_{h+1,k}(y)\right] + m_{hk}$$

$$= \sum_{h,k} \hat{r}_{hk}(s_{hk}, a_{hk}) - r(s_{hk}, a_{hk})$$

$$+ \sum_{s'}\left(\hat{P}_{hk}(s' \mid s_{hk}, a_{hk}) - p(s' \mid s_{hk}, a_{hk})\right)\left(V_{h+1,k}(s')\right)$$

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

$$\leq \sum_{k,h} |\hat{r}_{hk} - r| + H \sum_{s'} |\hat{p} - p| + b_{hk} + m_{hk}$$

- we have $\sum_{k,h} m_{hk} \leq 2H\sqrt{KH \log\left(\frac{2}{\delta}\right)}$ with

probability $\geq 1 - \frac{\delta}{2}$ thanks to Azuma-Hoeffding

- Other terms are bounded by confidence intervals.

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

with probability $1 - \delta$.

**II.5).**

$$\sum_{h=1}^{H} \sum_{s,a} \sqrt{N_{h,k}(s,a)} = HSA \sum_{h} \sum_{s,a} \frac{\sqrt{N_{hk}(s,a)}}{HSA}$$

$$\leq HSA \sqrt{\sum_{h} \sum_{s,a} \frac{N_{hk}(s,a)}{HSA}}$$

$$= \sqrt{HSA \sum_{h} \sum_{SA} N_{hk}(s,a)}$$

$$\leq \sqrt{HSA} \sqrt{\sum_{h} K} \qquad \text{because } \sum_{s,a} N_{k,h} \leq K$$

**Hence :**

$$R(T) \leq 2 \sum_{h,k} \sqrt{\frac{\log\left(\frac{8SAHK}{\delta}\right)}{2 N_{hk}(s,a)}} + H \sqrt{\frac{2}{N_{hk}}\left(\log\frac{(2^3-2)4SAHK}{\delta}\right)}$$

$$+ 2H \ KH \log(2/\delta)$$

$$\leq H^2 S^2 A + H\sqrt{SAK} + H\sqrt{S}\left(H^2 S^2 A + H\sqrt{SAK}\right) + 2H\sqrt{KH}$$

$$\leq H^2 S\sqrt{AK}$$