

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 $S' > 0$:

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

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

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

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

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

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

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

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

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

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

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 $k=1$

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

$$Q_{1,k}^*(s,a) = \hat{r}_{1,k}(s,a)$$

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

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

- Inductive step

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

$$Q_{n,k}(s,a) - Q_n^*(s,a) = \hat{r}_{n,k}(s,a) + b_{n,k}(s,a) + \sum_{s'} \hat{p}_{n,k}(s'|s,a) \min \left\{ 1, \max_a Q_{h+1,k}^*(s',a) \right\} \\ - \hat{r}_n(s,a) - \sum_{s'} p_n(s'|s,a) \max_a Q_{h+1}^*(s',a)$$

$$\geq \hat{r}_{n,k}(s,a) + b_{n,k}(s,a) - \hat{r}_n(s,a) \\ + \sum_{s'} \min \left\{ 1, \max_a Q_{h+1}^*(s',a) \right\} (\hat{p}_{n,k}(s'|s,a) - p_n(s'|s,a))$$

$$\geq \hat{r}_{n,k}(s,a) + b_{n,k}(s,a) - \hat{r}_n(s,a) \\ - \sum_{s'} \min \left\{ 1, \max_a Q_{h+1}^*(s',a) \right\} |\hat{p}_{n,k}(s'|s,a) - p_n(s'|s,a)|$$

$$\geq \hat{r}_{n,k}(s,a) + b_{n,k}(s,a) - \hat{r}_n(s,a) - H \underbrace{\sum_{s'} |\hat{p}_{n,k}(s'|s,a) - p_n(s'|s,a)|}_{||\hat{p}_{n,k}(s'|s,a) - p_n(s'|s,a)||_1}$$

$$\geq \hat{r}_{n,k}(s,a) + b_{n,k}(s,a) - \hat{r}_n(s,a) - H \beta_{n,k}^P(s,a)$$

$$\geq -\beta_{n,k}^P(s,a) + b_{n,k}(s,a) - H \beta_{n,k}^P(s,a) \geq 0$$

$$\Rightarrow b_{n,k}(s,a) \geq \beta_{n,k}^P(s,a) + H \beta_{n,k}^P(s,a)$$

Reinforcement Learning - Homework 3 Niculae's Violante

② 3) 1.

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

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

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

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

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

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

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

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

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

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

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

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

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