

Reinforcement Learning - HW 3

Elias Masquil

1) Best Arm Identification

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

Any-time bound:

$$P(\epsilon) = P\left(U_{i=1}^k - U_{t=1}^{\infty} \{|\hat{\mu}_{i,t} - \mu_i| > U(t, \delta')\}\right) \leq \delta$$

$$P\{|\hat{\mu}_{i,t} - \mu_i| > U(t, \delta')\} \leq \delta' \quad \text{for a single arm and time}$$

By Hoeffding's inequality

$$\begin{aligned} P(|\hat{\mu}_{i,t} - \mu_i| \geq U(t, \delta')) &\leq 2e^{-2tU(t, \delta')} = \delta' \\ \Rightarrow U(t, \delta') &= \sqrt{\frac{1}{2t} \log \frac{2}{\delta'}} \end{aligned}$$

By Union's bound

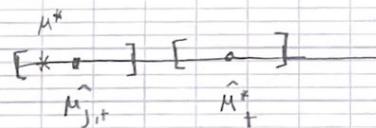
$$P(\epsilon) \leq \sum_{i=1}^k \sum_{t=1}^{\infty} \delta' \quad \left(\sum_i \frac{1}{t^2} = \frac{\pi^2}{6} \right)$$

$$\text{If we consider } \delta' = \frac{6}{\pi^2} \cdot \frac{1}{k} \cdot \frac{1}{t^2} \cdot \delta \Rightarrow P(\epsilon) \leq \delta$$

- ii) Show that with $\rho \geq 1 - \delta$, the optimal arm remains in the active set.

Arm elimination condition (for the optimal arm)

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



For deleting the optimal arm, μ^* should be outside of the confidence interval: $|\hat{\mu}_t^* - \mu^*| > U(t, \delta')$

$$\Rightarrow P(|\hat{\mu}_t^* - \mu^*| > U(t, \delta') \text{ for some } t) \leq P(\epsilon) \leq \delta$$

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

Reinforcement Learning - HW 3

Elias Masquil

i) iii) Let's define $\hat{\mu}_t^*$ as the estimated reward of the arm with the largest expected reward μ^*

Under event \mathcal{E} :

$$\hat{\mu}_t^* \geq \mu^* - U(t, s')$$

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

And the elimination condition is:

$$\hat{\mu}_t^* - U(t, s') \geq \hat{\mu}_{i,t} + U(t, s')$$

Then, the arm i will be dropped if

$$\mu^* - 2U(t, s') \geq \mu_i + 2U(t, s')$$

$$\Delta_i \geq 4U(t, s')$$

$$U(t, s') = \frac{1}{2\pi} \log \left(\frac{2\pi^2 t^2 K}{6s} \right)$$

$$\Delta_i^2 \geq 16 \cdot \frac{1}{2\pi} \log(bt) \quad \text{with } b = \pi \sqrt{\frac{K}{3s}}$$

$$at \geq \log(bt) \quad \text{with } a = \frac{\Delta_i^2}{16}$$

Using the suggestion

$$T \geq \frac{1 + \sqrt{2U + U}}{a} \quad \text{with } U = \log \frac{b}{a} - 1$$

$$(ii) T_i \geq \frac{2 \log \left(\frac{16\pi}{\Delta_i^2} \sqrt{\frac{K}{3s}} \right) - 2 + \log \left(\frac{16\pi}{\Delta_i^2} \sqrt{\frac{K}{3s}} \right)}{\Delta_i^2}, 16$$

iv) After sampling each sub-optimal arm T_i times, that arm will be eliminated.

Then, the sample complexity is just the sum of all T_i $\forall i=1, n$:

$$\mathcal{O} \left(\sum_{i \neq i^*} \frac{\log \left(\frac{16\pi}{\Delta_i^2} \sqrt{\frac{K}{3s}} \right)}{\Delta_i^2} \right)$$

Reinforcement Learning - HW 3

Elias Masquil

1) Best Arm Identification

- v) I think it won't work since it won't converge
After eliminating all sub-optimal arms, the algorithm might keep iterating across the optimal ones. Since $\Delta_i = 0$ for optimal arms, they won't be eliminated and then the time of removal will be ∞ .

2) Regret minimization in RL

- i) For fixed s, a, h, k we have two confidence intervals, for the rewards and for the transitions:

- Using Hoeffding's inequality:

$$P(\gamma_{\epsilon}) = P(|\hat{r}_{hk}(s,a) - r_{hk}(s,a)| \geq B_{hk}^r(s,a)) \leq 2e^{-2N_{hk}(s,a)B_{hk}^r(s,a)^2}$$

$$2e^{-2N_{hk}(s,a)B_{hk}^r(s,a)^2} = \delta_r$$

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

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

- Using Weissman inequality:

$$P(\gamma_{\epsilon_p}) = P(\|\hat{p}_{hk}(\cdot|s,a) - p_{hk}(\cdot|s,a)\|_1 \geq B_{hk}^p(s,a)) \leq (2^s - 2)e^{-\frac{N_{hk}(s,a)B_{hk}^p(s,a)^2}{2}}$$

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

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

$$\text{Then } P(\gamma_{\epsilon}) \leq \underset{s,a,h,k}{P(\gamma_{\epsilon})} + P(\gamma_{\epsilon_p})$$

$P(\gamma_{\epsilon})$ Probability of having at least one of the true $p(\cdot|s,a)$ or $r_{hk}(s,a)$ outside of the confidence bounds in a fixed s,a,h,k

Reinforcement Learning - HW 3

Elias Masquil

If we set $\delta_r = \delta_a = \frac{\delta^*}{2}$

$$P(\text{TE}_{s,a,h,k}) \leq \delta^*$$

We have a bound for any k, h, s, a .

$$\text{Then } P(\varepsilon) = 1 - P(\text{TE}) = 1 - P\left(\bigcup_{sahk} \text{TE}_{s,a,h,k}\right)$$

By Union bound

$$1 - P\left(\bigcup_{sahk} \text{TE}_{s,a,h,k}\right) \geq 1 - \sum_{s,a,h,k} P(\text{TE}_{s,a,h,k})$$

And we want $P(\varepsilon) \geq 1 - \delta/2$

$$\sum_{s,a,h,k} P(\text{TE}_{s,a,h,k}) \leq \frac{\delta}{2}$$

$$SAHK \delta^* = \frac{\delta}{2} \Rightarrow \delta^* = \frac{\delta}{2SAHK}$$

Now we can correctly define the original confidence bounds, using the δ^* value.

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

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

Reinforcement Learning - HW3

Elias Masquil

2) Regret minimization in RL

ii) $\hat{r}_{H,K}(s,a)$

Base case

$$h = H$$

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

$$Q_H^*(s,a) = \hat{r}_{H,K}(s,a)$$

Since we are under ϵ event, we also know that

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

Then

$$Q_{H,K}(s,a) \geq r_{H,K}(s,a) + b_{H,K}(s,a) - B_{H,K}^r(s,a)$$

If we define a bonus such that

$$b_{H,K}(s,a) \geq B_{H,K}^r(s,a) \quad \text{the base case is true.}$$

Inductive step

Assume that $Q_{h,K}(s,a) \geq Q_h^*(s,a)$, let's prove

the inequalities for $h-1$

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

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

$$Q_{h-1,K}^*(s,a) = \hat{r}_{h-1,K}(s,a) + \sum_{s'} \hat{p}_{h-1,K}(s'|s,a) \max_a Q_{h-1,K}^*(s,a)$$

Taking the difference

$$Q_{h-1,K}(s,a) - Q_{h-1,K}^*(s,a) = \hat{r}_{h-1,K}(s,a) + b_{h-1,K}(s,a) - \hat{r}_{h-1,K}(s,a)$$

$$+ \sum_{s'} \hat{p}_{h-1,K}(s'|s,a) \min \left\{ H, \max_a Q_{h,K}(s,a) \right\}$$

$$- \hat{p}_{h-1,K}(s'|s,a) \max_a Q_{h-1,K}^*(s,a)$$

Reinforcement Learning - HW 3

Elias Masouil

By using the inductive assumption

$$\begin{aligned} Q_{h-1,K}^* - Q_{h-1,K}^* &\geq \hat{r}_{h-1,K}(s,a) + b_{h-1,K}(s,a) - r_{h-1,K}(s,a) \\ &\quad + \sum_{s'} \min\{H, \max_a Q_{h,K}(s',a)\} (\hat{P}_{h-1,K}(s'|s,a) - P_{h-1,K}(s'|s,a)) \geq \\ &\geq \hat{r}_{h-1,K}(s,a) + b_{h-1,K}(s,a) - r_{h-1,K}(s,a) - \sum_{s'} \min\{H, \max_a Q_{h,K}(s',a)\} |\hat{P} - P| \\ &\geq \hat{r}_{h-1,K}(s,a) + b_{h-1,K}(s,a) - r_{h-1,K}(s,a) - H \geq |\hat{P}_{h-1,K}(s'|s,a) - P_{h-1,K}(s'|s,a)| \end{aligned}$$

Since we are under event \mathcal{E} , both confidence intervals hold

$$\begin{aligned} &\geq \hat{r}_{h-1,K}(s,a) + b_{h-1,K}(s,a) - r_{h-1,K}(s,a) - H B_{h-1,K}^P(s,a) \geq \\ &\geq -B_{h-1,K}^R(s,a) - H B_{h-1,K}^P(s,a) + b_{h-1,K}(s,a) \\ &\geq 0 \iff b_{h-1,K}(s,a) \geq B_{h-1,K}^R(s,a) + H B_{h-1,K}^P(s,a) \end{aligned}$$

$$\begin{aligned} \text{iii) 1)} V_h^{T_K}(s_{h,K}) &= r(s_{h,K}, a_{h,K}) + \sum_{s' \in s_{h+1,K}} p(s'|s,a) V_{h+1}^{T_K}(s') \\ &= r(s_{h,K}, a_{h,K}) + \sum_{s' \in s_{h+1,K}} p(s'|s,a) (V_{h+1}(s') - \delta_{h+1,K}(s')) \\ &= r(s_{h,K}, a_{h,K}) + \mathbb{E}_P[V_{h+1}(s')] - \mathbb{E}_P[\delta_{h+1,K}(s')] \\ &= r(s_{h,K}, a_{h,K}) + \mathbb{E}_P[V_{h+1}(s')] - \delta_{h+1,K}(s_{h+1,K}) - m_{h,K} \end{aligned}$$

$$\begin{aligned} \text{2)} V_{h,K}(s_{h,K}) &= \min_a \{H, \max_{a_{h,K}} Q_{h,K}(s,a)\} \\ &\leq Q_{h,K}(s_{h,K}, a_{h,K}) \quad \text{since } a_{h,K} \text{ was} \\ &\quad \text{the greedy action } (a_{h,K} = \arg \max_a Q_{h,K}(s_{h,K}, a)) \end{aligned}$$

Reinforcement Learning - HW 3

Elias Masquil

2) Regret minimization in RL

iii) 3)

$$\delta_{1K}(s_{1K}) = V_{1K}(s) - V_1^{\pi_K}(s) \leq$$

$$\leq Q_{1K}(s_{1K}, a_{1K}) - r(s_{1K}, a_{1K}) - \mathbb{E}_P[V_{2,K}(s')] \quad \text{↑ develop this}$$

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

$$\leq \dots$$

$$\leq \sum_{h=1}^H Q_{hK}(s_{hK}, a_{hK}) - r(s_{hK}, a_{hK}) - \mathbb{E}_P[V_{h+1,K}(s)] + m_{hK}$$

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

Since our estimates are optimistic

$$\leq \sum_{k=1}^K V_1(s_{1,k}) - V_1^{\pi_K}(s_{1,k}) = \sum_{k=1}^K \delta_{1K}(s_{1,k})$$

$$\leq \sum_{k=1}^K \sum_{h=1}^H Q_{hK}(s_{hK}, a_{hK}) - r(s_{hK}, a_{hK}) - \mathbb{E}_{Y \sim P}[V_{h+1,K}(Y)] + m_{hK}$$

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

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

$$\leq \sum_{k,h} |\hat{r}_{hk}(s_{hk}, a_{hk}) - r(s_{hk}, a_{hk})| + H \sum_{s'} |\hat{p}_{hk}(s'|s_{hk}, a_{hk}) - p(s'|s_{hk}, a_{hk})|$$

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

→ Note that by Azuma-Hoeffding:

$$\sum_{k,h} m_{hk} \leq 2H\sqrt{K H \log(2/\delta)} \text{ with proba } \geq 1 - \delta/2$$

→ The other terms are bounded by the confidence intervals

$$\sum_{k,h} |\hat{r}_{hk}(s_{hk}, a_{hk}) - r(s_{hk}, a_{hk})| + H \sum_{s'} |\hat{p}_{hk}(s'|s_{hk}, a_{hk}) - p(s'|s_{hk}, a_{hk})|$$

$$\leq \sum_{k,h} B_{hk}^r(s, a) + H B_{hk}^p(s, a) \leq \sum_{k,h} b_{hk}(s_{hk}, a_{hk})$$

with proba $\geq 1 - \delta/2$

Reinforcement Learning - HW 3

Elias Masquel

Finally, combining the two previous bounds:

$$R(T) \leq \sum_{k=1}^H b_{k,n}(s_{kn}, a_{kn}) + \sqrt{2HKH \log(2/\delta)}$$

with probability $1 - \delta$

$$\text{v) } \sum_{h=1}^H \sum_{s,a} \sqrt{N_{h,K}(s,a)} = HSA \sum_h \sum_{s,a} \frac{\sqrt{N_{h,K}(s,a)}}{HSA}$$

Using Jensen inequality for concave functions

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

Since $\sum_{s,a} N_{h,K}(s,a) \leq K$

$$\Rightarrow \leq \sqrt{HSA} \sqrt{\sum_h K} = H\sqrt{SAK}$$

Now making the inequalities between the bonus and the confidence bounds, hold with equality. we have

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

$$\leq \frac{2}{\sqrt{2}} \sqrt{\log\left(\frac{8SAHK}{\delta}\right)} \sum_{h,K} \frac{1}{\sqrt{N_{h,K}(s,a)}} + H \sqrt{2 \log\left(\frac{(2^{h-2})4SAHK}{\delta}\right)} \sum_{h,K} \frac{1}{\sqrt{N_{h,K}(s,a)}} + 2H\sqrt{KH \log(2/\delta)}$$

$$\approx H^2 S^2 A + H\sqrt{SAK} + H\sqrt{S(H^2 S^2 A + H\sqrt{SAK})} + 2H\sqrt{KH}$$

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