


RECALL Hoeffding:

$$\mathbb{P}_S \left[|S_n - \mathbb{E}(S_n)| \geq t \right] \leq 2 \exp\left(-\frac{2t^2}{n}\right)$$

$$\mathbb{P}_S \left[\left| \frac{S_n}{n} - \frac{\mathbb{E}(S_n)}{n} \right| \geq \frac{t}{n} \right] \leq 2 \exp\left(-\frac{2t^2}{n}\right)$$

$$\mathbb{P}_S \left[\left| \hat{P}(s|s,a) - P(s|s,a) \right| \geq \frac{t}{n} \right] \leq 2 \exp\left(-\frac{2t^2}{n}\right)$$

$$\text{Probability} \left(\exists s, a \mid \left| \hat{P}(\cdot|s,a) - P(\cdot|s,a) \right| \geq \frac{t}{n} \right) \leq 2|S||A| \exp\left(-\frac{2t^2}{n}\right)$$

$$\text{We want : } 2|S||A| \exp\left(-\frac{2t^2}{n}\right) \leq \delta$$

$$\text{If } \forall \text{ policies } \pi \text{ we want } \|\delta^{\pi} - \delta\|_{\infty} \leq \frac{\epsilon}{2}$$

then we want

$$\max_{s,a} \left\| \hat{P}(\cdot|s,a) - P(\cdot|s,a) \right\| \leq \frac{(\frac{\epsilon}{2})^2}{2}$$

$$\therefore \frac{t}{n} = \frac{(1-\alpha)^2 \varepsilon}{2} \quad \therefore t = \frac{(1-\alpha)^2 \varepsilon n}{2}$$

$$\therefore 2 \|s\|_A^2 \exp\left(-\frac{2(1-\alpha)^4 \varepsilon^2 n}{4}\right) \leq \delta$$

∴

$$\exp\left(\frac{(1-\alpha)^4 \varepsilon^2 n}{2}\right) \geq \frac{2 \|s\|_A^2}{\delta}$$

$$n \geq \frac{2}{(1-\alpha)^4 \varepsilon^2} \log\left(\frac{2 \|s\|_A^2}{\delta}\right)$$

$$\therefore \underline{\text{\# samples}} = \frac{2 \|s\|_A^2}{(1-\alpha)^4 \varepsilon^2} \log\left(\frac{2 \|s\|_A^2}{\delta}\right)$$

$$= \tilde{O}\left(\frac{\|s\|_A^2}{(1-\alpha)^4 \varepsilon^2}\right)$$

Can one improve this?

Turns out we can:

We do not need $\|g^\pi - g^{\circ\pi}\|_\infty$ for all

polices. Need only for g^* !!

More refined:

Lemma: let $\delta \geq 0$, with $\delta \leq 1 - \delta$

$$\|g^* - \hat{g}^{\pi^*}\|_{\infty} \leq \frac{\gamma}{1-\delta} \sqrt{\frac{2 \log(2|S||A|/\delta)}{N}}$$

Pf.

$$\begin{aligned} \|g^* - \hat{g}^{\pi^*}\|_{\infty} &= \gamma \left\| P^{\pi^*} g^* - \hat{P}^{\pi^*} \hat{g}^{\pi^*} \right\|_{\infty} \\ &\leq \gamma \left\| P^{\pi^*} g^* - \hat{P}^{\pi^*} g^* \right\|_{\infty} + \gamma \left\| \hat{P}^{\pi^*} g^* - \hat{P}^{\pi^*} \hat{g}^{\pi^*} \right\|_{\infty} \\ &= \gamma \left\| P V^* - \hat{P} V^* \right\|_{\infty} + \gamma \left\| \hat{P}^{\pi^*} (g^* - \hat{g}^{\pi^*}) \right\|_{\infty} \\ &\leq \gamma \left\| P V^* - \hat{P} V^* \right\|_{\infty} + \gamma \left\| g^* - \hat{g}^{\pi^*} \right\|_{\infty} \\ \therefore \left\| g^* - \hat{g}^{\pi^*} \right\|_{\infty} &\leq \frac{\gamma}{1-\delta} \left\| P V^* - \hat{P} V^* \right\|_{\infty}. \end{aligned}$$

$$\|P V^* - \hat{P} V^*\|_{\infty} = \max_{s,a} \left| \mathbb{E}_{s' \sim P(\cdot|s,a)} [V^*(s')] - \mathbb{E}_{s' \sim \hat{P}(\cdot|s,a)} [V^*(s')] \right|$$

If this is $\leq \sqrt{\frac{2 \log(2|S||A|/\delta)}{N}}$ we are done!

Use Hoeffding's inequality.

Consider the vector

$$v = \boxed{V(s_1) \mid V(s_2) \mid \dots \mid \dots}$$

Pick $s' \sim P(\cdot \mid s, a)$ and consider the random variable $x = V(s')$.

$$\mathbb{E}(x) = \sum_{s' \sim P(\cdot \mid s, a)} P(s' \mid s, a) V(s')$$

Instead pick $s_1 \sim P(\cdot \mid s, a)$, $s_2 \sim P(\cdot \mid s, a)$, \dots

and consider $x_1 = V(s_1)$, $x_2 = V(s_2)$, \dots

$$\mathbb{E}(x_1) = \mathbb{E}(x_2) = \dots = \sum_{s' \sim P(\cdot \mid s, a)} P(s' \mid s, a) V(s')$$

$$\text{Take } x = \frac{x_1 + \dots + x_N}{N} \quad \mathbb{E}(x) = \sum_{s' \sim P(\cdot \mid s, a)} P(s' \mid s, a) V(s')$$

Clearly each $x_i \in [0, 1]$

$$\therefore \Pr \left[|x - \mathbb{E}(x)| > \frac{t}{N} \right] \leq 2 \exp \left(\frac{-2t^2}{N} \right)$$

Probs that for some (S, A) the above fails
 \leq at most

$$2|S||A| \exp\left(-\frac{2t^2}{N}\right)$$

want this to be $\leq \delta$.

$$\therefore \exp\left(-\frac{2t^2}{N}\right) \geq \frac{2|S||A|}{\delta}$$

$$t^2 \geq \frac{N}{2} \ln\left(\frac{2|S||A|}{\delta}\right)$$

$$\therefore t \geq \sqrt{\frac{N \ln\left(\frac{2|S||A|}{\delta}\right)}{2}}$$

$$\|(\hat{P} - \hat{P})\sqrt{V}\|_{\infty} \leq \frac{t}{2} = \sqrt{\frac{\ln\left(\frac{2|S||A|}{\delta}\right)}{2N}} \quad \text{with}$$

prob $(1-\delta)$

Now set:

$$\frac{\delta}{1-\delta} \sqrt{\frac{\log(2|S||A|/\delta)}{2N}} \leq \epsilon$$

$$N \geq \frac{\sigma^2}{(1-\gamma)^2 \epsilon^2} \frac{1}{2} \log \left(\frac{2 \|s\| |A|}{\delta} \right)$$

The factor is $\frac{1}{(1-\gamma)^2 \epsilon^2}$

$$\therefore \# \text{ samples: } \frac{\|s\| |A|}{(1-\gamma)^2 \epsilon^2} 2 \log \left(\frac{2 \|s\| |A|}{\delta} \right)$$

But this requires $\hat{V}(s')$ - we have no access to that

Then: (Azar, Munro, Kappen)

With prob $(1-\delta)$, an error of $O \left(\frac{N \log(N/\delta)}{(1-\gamma)^2 \epsilon^2} \right)$

Sample suffice to find ϵ -optimal estimate of

action value, and to find an ϵ -optimal policy.

- There is a (almost) matching upper bound

Thm:

$$\|\hat{Q}^* - Q^*\|_{\infty} \leq \gamma \sqrt{\frac{c \log(c/s|A|/\delta)}{N}} + \frac{c\gamma}{(1-\gamma)^2} \frac{\log(c/s|A|/\delta)}{N}$$

= Strategic Exploration:

- Don't have access to transitions at each state.
- Can execute trajectories in the MDP.
- Agent has to engage in exploration - to reach new states where not enough samples are

sem!

The algorithm maintains an estimate of transition prob $P(s'|s,a) \forall s', \text{nbrs } s$.

It also estimates reward.

- If s is visited enough, declare it as known.
- Learning is complete when all states are known.
- In a known state - take the optimal action.

Input: parameter $m; \epsilon;$

1. K (known states) = \emptyset . $\forall s, a, s', n(s, a) = n(s, a, s') = 0$
 $R(s, a) = 0$
2. Observe initial state s_0 , π_0 be an initial policy.
- 3 **for** rounds $t = 0, 1, \dots$
- 4 **if** a state has become known (i.e.) $n(s, a) \geq m \forall a$ **then**
- 5 update $K = K \cup \{s\}$.
- 6 let \hat{M}_t have $\hat{P}(s' | s, a) = \frac{n(s, a, s')}{n(s, a)}$, $\hat{r}(s, a) = \frac{R(s, a)}{n(s, a)}$
- 7 let \hat{M}_t be induced MDP.
 $\hat{\pi}_t = \hat{\pi}^*(\hat{M}_t)$ - optimal policy in \hat{M}_t
- 8 **else**
- 9 if $t \geq 1$ $\pi_t = \pi_{t-1}$
- 10 **endif**
- 11 if $s_t \in K$, $a_t = \pi_t(s_t)$, else $a_t = \operatorname{argmin}_a n(s, a)$
- 12 Get reward r_t & observe s_{t+1}
- 13 **if** $s_t \notin K$ **then**
- 14 update $n(s_t, a_t) + 1$; $R(s_t, a_t) + r_t$;
 $n(s_t, a_t, s_{t+1}) + 1$;
- 15 **endif**
- 16 **endfor**

INDUCED MDP: M is given;

Assumes $K \subseteq S$;

$$\forall s \in K, P_{M_K}(s'|s, a) = P_M(s'|s, a) \quad \&$$

$$r_{M_K}(s, a) = r_M(s, a)$$

$$\forall s \notin K, P_{M_K}(s'|s, a) = 1 \quad (s' = s) \quad \text{and}$$

$$r_{M_K}(z|s, a) = 1 \quad (z = 1)$$