

Online RL | Tabular :  $M = (S, A, P, R, \gamma, d_0)$ .

For  $t=1, 2, 3, \dots, T$  (assume all traj. from  $s \sim d_0$  terminates in  $H$  steps)

- Learner collects trajectory using  $\pi_t$ .
- $(s_1, a_1, r_1, \dots, s_H, a_H, r_H)$

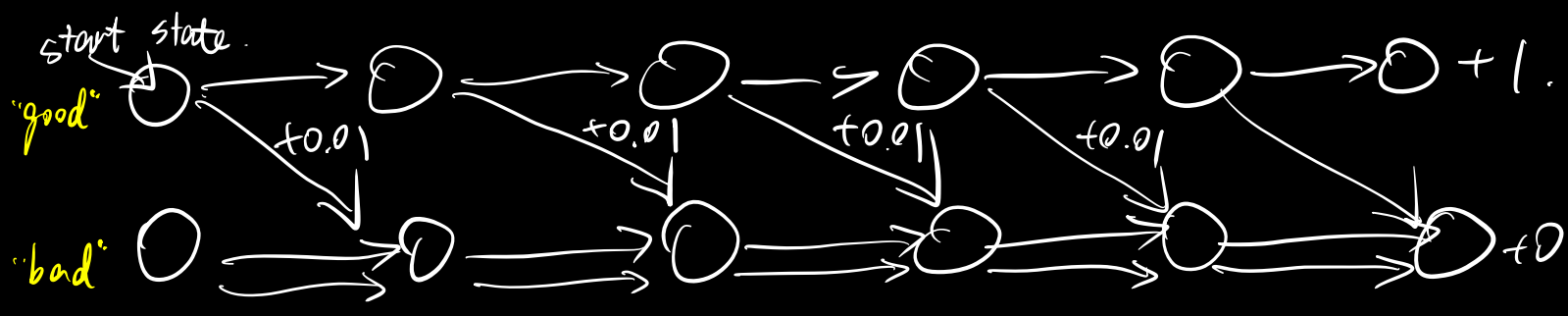
After  $T$  trajectories, output policy  $\hat{\pi}$ .

Guarantee w.p.  $\geq 1 - \delta$ ,  $J(\pi^*) - J(\hat{\pi}) \leq \epsilon \cdot V_{max}$ .

w/ sample complexity  $T = \text{poly}(|S|, |A|, \frac{1}{1-\gamma}, \frac{1}{\epsilon}, \frac{1}{\delta})$ .

Unif exploration is ineffective.

"Combination Lock" (not fundamental hardness).



if explore w/ unif over actions.

(extend to  $\epsilon$ -greedy, softmax).

$R_{max}$  ] At any point during exploration, maintain.

- $n(s,a)$  : #times we see  $(s,a)$  in data
- $n(s,a,s')$  : #times we see  $(s,a,s')$  in data.

\* stop accumulating data after  $n(s,a) = m$ . (artifact for clean concentration analysis)

$R_{max}$  : define  $K := \{ (s,a) : n(s,a) = m \}$ . "Known".

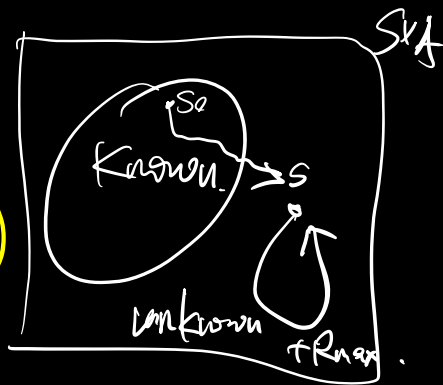
• Build MDP  $\hat{M}_K$  :  $\hookrightarrow$  threshold.

$$\hat{P}_K(s'|s,a) = \begin{cases} n(s,a,s')/n(s,a), & \text{if } (s,a) \in K. \\ \mathbb{I}[s'=s], & \text{if } (s,a) \notin K. \end{cases}$$

$$\hat{R}_K(s,a) = \begin{cases} R(s,a) & \text{if } (s,a) \in K. \\ R_{max}. & \text{if } (s,a) \notin K. \end{cases}$$

• next-exploration policy:  $\pi_{\hat{M}_K}^*$ .

"Optimism in face of uncertainty" (OFU)



Remark : alg needs to determine stopping criteria.  
(ignored for now; will comment on later)

Define  $M_K$ :

$$P_K(s'|s,a) = \begin{cases} P(s'|s,a), & \text{if } (s,a) \in K. \\ \mathbb{I}[s'=s], & \text{o.w.} \end{cases}$$

$$R_K(s,a) = \begin{cases} R(s,a), & \text{if } (s,a) \in K. \\ R_{\max}, & \text{o.w.} \end{cases}$$

	$M$	$M_K$	$\hat{M}_K$
$(s,a) \in K$	$= M$	$= M$	$\approx M$
$(s,a) \notin K$	$= M$	self-loop	self-loop

↖ ↗

Define:  $\text{dis}(M_K, \hat{M}_K) = \max_{s,a} \|P_K(\cdot|s,a) - \hat{P}_K(\cdot|s,a)\|_1$

Lemma 1. For any fixed  $s,a$ , w/  $n(s,a) = m$ .

$$\text{w.p. } \geq 1 - \delta, \quad \|P_K(\cdot|s,a) - \hat{P}_K(\cdot|s,a)\|_1 \leq \sqrt{\frac{2}{m} \log \frac{2^{251}}{\delta}}$$

$$\text{Lemma 2. } \|V_{M_K}^* - V_{\hat{M}_K}^*\|_{\infty} \leq \frac{\text{dis}(M_K, \hat{M}_K) \cdot V_{\max}}{2(1-\delta)}$$

$$\text{Lemma 3. } \forall \pi, \|V_{M_K}^{\pi} - V_{\hat{M}_K}^{\pi}\|_{\infty} \leq (\dots)$$

Lemma 4. Consider  $M_K$  &  $M$ ,  $\forall \pi$ .

$$|J_M(\pi) - J_{M_K}(\pi)| \leq V_{\max} \cdot \mathbb{P}_M[\text{escape}_K(\tau) | \pi]$$

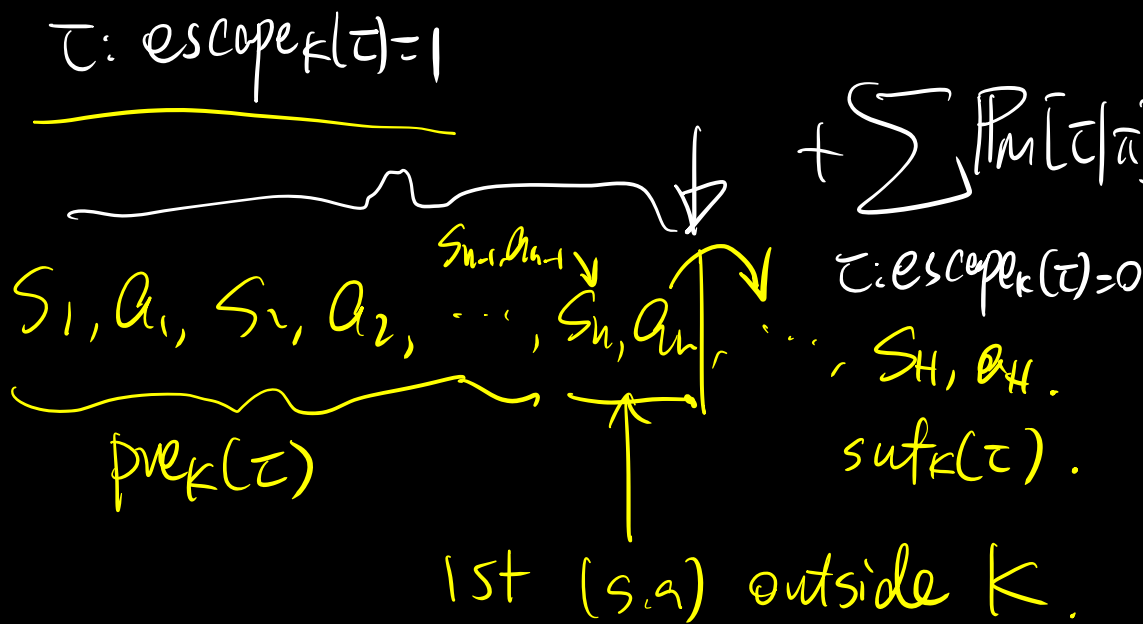
where  $\text{escape}_K(\tau) = \mathbb{I}[\exists h, (s_h, a_h) \notin K]$ .

•  $P_M[-|\pi]$  is the prob. when traj.  $\tau$  is generated, by running  $\pi$  in  $M$ .

Proof.  $J_M(\pi) \stackrel{!}{=} J_{M_K}(\pi) + V_{\max} \cdot P_M[\dots]$

$$= \sum_{\tau} P_M[\tau|\pi] \cdot R_M(\tau) \quad \left| \begin{array}{l} R_M(\tau) \\ = \sum_{h=1}^H \gamma^{h-1} R(s_h, a_h) \end{array} \right.$$

$$= \sum_{\tau: \text{escape}_K(\tau)=1} P_M[\tau|\pi] \cdot (R_M(\text{pre}_K(\tau)) + R_M(\text{suf}_K(\tau))) + \sum_{\tau: \text{escape}_K(\tau)=0} P_M[\tau|\pi] \cdot R_M(\tau)$$



$$\leq \sum_{\text{esc}=1} P_M[\tau|\pi] (R_M(\text{pre}_K(\tau)) + V_{\max})$$

$$+ \sum_{\text{esc}=0} P_M[\tau|\pi] R_M(\tau).$$

$$= \sum_{pre_k(\tau)} P_M [pre_k(\tau) | \pi] (R_M(pre_k(\tau)) + V_{max}) + \sum_{esc=0} P_M[\tau | \pi] R_M(\tau).$$

$$= \sum_{pre_k(\tau)} P_{M_k} [pre_k(\tau) | \pi] (R_{M_k}(pre_k(\tau)) + V_{max}) + \sum_{esc=0} P_{M_k}[\tau | \pi] R_{M_k}(\tau)$$

$$P[\tau] = d_0(s_1) \cdot \pi(a_1 | s_1) \cdot \underbrace{P(s_2 | s_1, a_1)} \cdot \dots$$

$$\leq J_{M_k}(\pi) + \sum_{pre_k(\tau)} P_{M_k} [pre_k(\tau) | \pi] \cdot \underbrace{V_{max}}_g$$

$$= J_{M_k}(\pi) + \left( \sum_{pre_k(\tau)} P_M [pre_k(\tau) | \pi] \right) V_{max}$$

$$= J_{M_k}(\pi) + P_M [escape_k(\tau) | \pi] \cdot V_{max} \quad \square$$

Remark: only uses the fact  $M = M_k$  on  $k$ .  
(didn't use the  $R_{max}$  design).

Putting things together:

"optimal or explore": either near-optimal (in  $M$ )  
 or explore (s.a.)  $\notin K$ .

$$\varepsilon \cdot V_{\max} < \underbrace{J_M(\pi_M^*) - J_M(\pi_{\hat{M}_K}^*)}_{\Delta} \leq \underbrace{J_{M_K}(\pi_M^*) - J_M(\pi_{\hat{M}_K}^*)}_{\text{(optimism)}}$$

$$\leq \underbrace{J_{M_K}(\pi_{\hat{M}_K}^*)}_{\Delta} - J_M(\pi_{\hat{M}_K}^*)$$

$$\leq \underbrace{J_{\hat{M}_K}(\pi_{\hat{M}_K}^*)}_{\Delta} + \text{dis}(M_K, \hat{M}_K) \cdot \frac{V_{\max}}{2(1-\gamma)} \quad (\text{Lemma 2})$$

$$\leq \underbrace{J_{M_K}(\pi_{\hat{M}_K}^*)}_{\Delta} + \text{dis}(M_K, \hat{M}_K) \cdot \frac{V_{\max}}{1-\gamma} - \underbrace{J_M(\pi_{\hat{M}_K}^*)}_{\Delta}$$

$$\leq \mathbb{P}_M[\text{escape} \mid \pi_{\hat{M}_K}^*] \cdot V_{\max} + \text{dis}(M_K, \hat{M}_K) \cdot \frac{V_{\max}}{1-\gamma}$$

$\Downarrow$

$$\varepsilon < \mathbb{P}_M[\text{escape} \mid \pi_{\hat{M}_K}^*] + \text{dis}(M_K, \hat{M}_K) / (1-\gamma)$$

$\downarrow$  exploration policy       $\downarrow$  controlled by  $m$

then set  $m = \tilde{O}\left(\frac{|S|}{\epsilon^2(1-\gamma)^2} \log \frac{1}{\delta}\right)$  &

$\Rightarrow \text{dis}(M_k, \hat{M}_k) \leq \epsilon(1-\gamma)/2$  &

$\Rightarrow \epsilon < \mathbb{P}_M[\text{escape} | \dots] + \epsilon/2$

$\Rightarrow \mathbb{P}_M[\text{escape} | \dots] > \frac{\epsilon}{2}$

---

$$\frac{|S \times A| \cdot m}{\epsilon/2} \Rightarrow \tilde{O}\left(\frac{|S|^2 |A|}{\epsilon^3 (1-\gamma)^2} \log \frac{1}{\delta}\right)$$

Remark: Stopping criterion: *optimistic value*

$$J_M(\pi_M^*) - J_M(\pi_{\hat{M}_k}^*) \lesssim J_{\hat{M}_k}(\pi_{\hat{M}_k}^*) - J_M(\pi_{\hat{M}_k}^*)$$

if this is  $\leq O(\epsilon)$ .

can stop.

$\uparrow$

known

$\triangleleft$

estimated  
by MC.