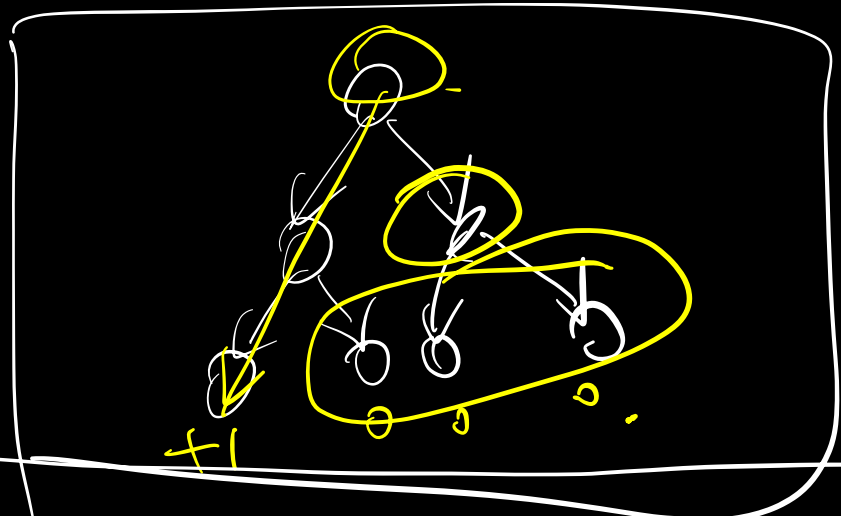


$$E_{\pi} = [ \nabla_{\theta} \gamma^t (s, a) \sum_{t'} r_{t'} ]$$



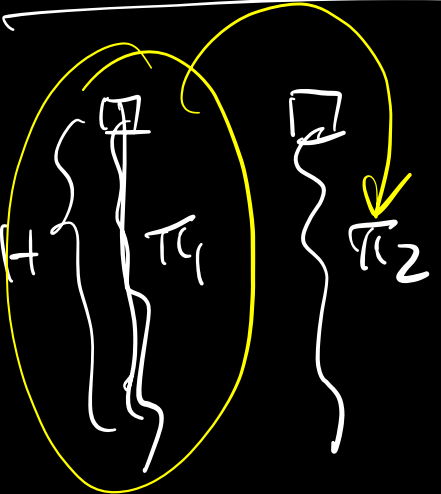
Online RL

For  $t=1, 2, 3, \dots, T$ .

- alg picks  $\pi_t$ .
- run  $\pi_t$  to collect traj.

$$s_H^{(t)}, a_H^{(t)}, r_H^{(t)}, \dots, s_H^{(t)}, a_H^{(t)}, r_H^{(t)}$$

$T$



Sample complexity | Alg: output  $\hat{\pi}$  .  $1/\epsilon^2$ .

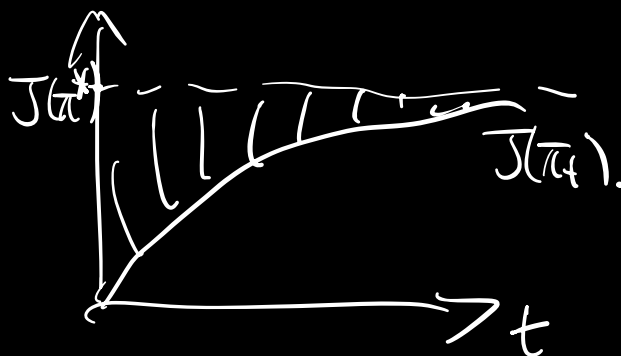
Guarantee: w.p.  $\geq 1 - \delta$ .  $J(\pi^*) - J(\hat{\pi}) \leq \epsilon$ .

$$T = \text{poly}(|S|, |A|, H, 1/\epsilon, \log(1/\delta)).$$

Regret bound. |  $\text{Regret}_T = \sum_{t=1}^T J(\pi^*) - J(\pi_t)$ .

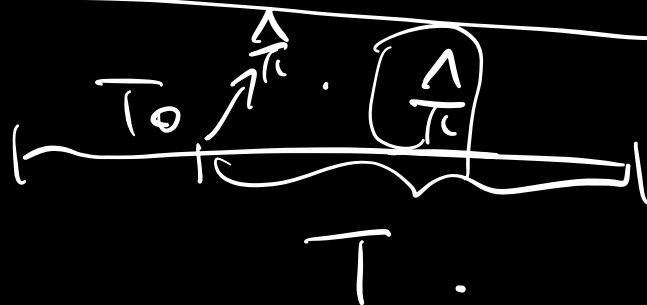
Guarantee: w.p.  $\geq 1 - \delta$ .

$$\text{Regret}_T = \text{poly}(S, A, H, \log(1/\delta)) \underline{\underline{o(T)}}.$$



Sample comp  $\rightarrow$  regret.

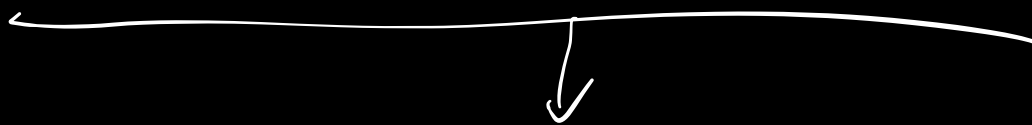
Alg:  $1/\epsilon^2$ .



$$\sum_{t=1}^T J(\pi^*) - J(\pi_t)$$

$$\leq T_0 \cdot V_{\max} + \underbrace{(T - T_0)}_T \underbrace{(J(\pi^*) - J(\hat{\pi}))}_{\epsilon}$$

$$\leq \underline{T_0} V_{\max} + \frac{cT}{\sqrt{T_0}}$$



$$T_0 V_{\max} + \frac{cT}{2\sqrt{T_0}} + \frac{cT}{2\sqrt{T_0}}$$

achieved by

$$T_0 V_{\max} = \frac{cT}{2\sqrt{T_0}} \quad \text{IV}$$

$$\Rightarrow \sqrt[3]{T_0 V_{\max} \cdot \frac{cT}{2\sqrt{T_0}} \cdot \frac{cT}{2\sqrt{T_0}}} = \sqrt[3]{\frac{c^2 T^2 V_{\max}}{4}}$$

$$2\sqrt{T_0} = \frac{cT}{V_{\max}}$$

$$T_0 = \left( \frac{cT}{2V_{\max}} \right)^{2/3}$$

regret  $\rightarrow$  SC. run regret-min for  $T$ .

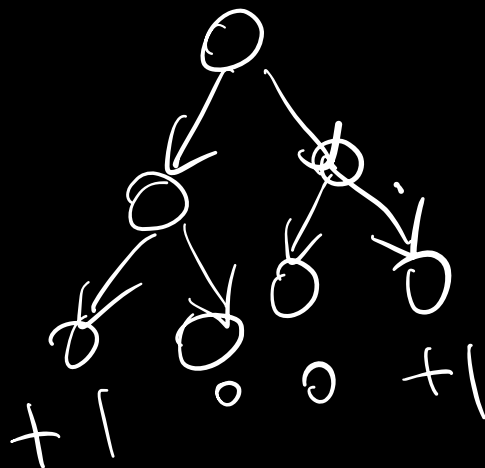
Output:

draw  $i \sim \text{Unif}([T])$ .  
 $\downarrow$   
 run  $\pi_i$ .  $(1, 2, \dots, T)$

$$\begin{aligned} J(a^*) - \frac{\sum_{i=1}^T J(\pi_i)}{T} &= \frac{1}{T} \sum_{i=1}^T J(a^*) - J(\pi_i) \\ &= \frac{1}{T} O(\sqrt{T}). \end{aligned}$$

$$\pi: S \rightarrow \Delta(A).$$

$$= \left( \frac{1}{\sqrt{4}} \right).$$



L.  
R.

$$\pi_1(\cdot | S_1), \quad \pi_2(\cdot | S_2),$$

$$\pi_1 \quad \pi_2 \quad 1/2, \quad 1/2,$$

$$a_1 \sim \frac{\pi_1(\cdot | S_1) + \pi_2(\cdot | S_2)}{2}$$

$$S_1, a_1, S_2 \quad \boxed{a_2}?$$

$$Pr(i | S_1, a_1, S_2),$$