

Online RL

$$M = (S, A, P, R, H, \underline{d_0}).$$

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

• Learner chooses π_t to collect a traj.

• $(s_t, a_t, r_t, \dots, s_H, a_H, r_H)$. "pure exploration"

After T rounds. learner output $\hat{\pi}$.

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

"simple regret". $\approx f(\frac{1}{T})$.

to achieve ϵ . how large $T = g(\frac{1}{\epsilon})$ needs to be.

"sample complexity"

Regret: ("cumulative") regret = $\sqrt{V_{\max} \cdot T}$

$$\text{Regret} := \sum_{t=1}^T (J(\pi^*) - J(\pi_t)) = \underline{\underline{o(T)}}.$$

"exploration - exploitation".

$$\sqrt{T}$$

$$\text{tabular MDP: } \sqrt{HSA T}$$

Sample comp \Rightarrow regret.

$$\hat{\pi} \quad J(\pi^*) - J(\hat{\pi}) = O\left(\frac{1}{\sqrt{T}}\right).$$

$\pi_1, \pi_2, \dots, \pi_T.$

(1) Run SC alg for $T_1 \leq T$ episodes

(2) $\pi_{T_1:T} = \hat{\pi}.$

$$\text{Regret} = T_1 + (T - T_1) \cdot \sqrt{\frac{T}{T_1}}.$$

$$\leq T_1 + T \sqrt{\frac{T}{T_1}} = O(T^{2/3}).$$

$$x + T \sqrt{\frac{T}{x}} = x + \frac{T \sqrt{T}}{2 \sqrt{x}} + \frac{T \sqrt{T}}{2 \sqrt{x}}$$

$$\geq \sqrt[3]{x \cdot \frac{T \sqrt{T}}{2 \sqrt{x}} \cdot \frac{T \sqrt{T}}{2 \sqrt{x}}}$$

$$= \sqrt[3]{\frac{T^2}{4}} = O(T^{2/3}).$$

$$T_1 = O(T^{2/3}).$$

Regret \Rightarrow SC. $\pi_1, \pi_2, \dots, \pi_T$

$$\sum_{t=1}^T \left(\underline{J}(\pi^*) - \underline{J}(\pi_t) \right) = O(\sqrt{T}).$$

$\hat{\pi} =$ unif mixture of $\pi_1, \pi_2, \dots, \pi_T$.

i.e. draw $t \in \text{unif}\{1, \dots, T\}$,
then run π_t .

$$J(\hat{\pi}) = \frac{1}{T} \sum_{t=1}^T J(\pi_t)$$

$$\begin{aligned} J(\pi^*) - J(\hat{\pi}) &= \frac{1}{T} \left(\underbrace{\sum (J(\pi^*) - J(\pi_t))}_{O(\sqrt{T})} \right) \\ &= \frac{1}{T} \sqrt{T} = \frac{1}{\sqrt{T}}. \end{aligned}$$