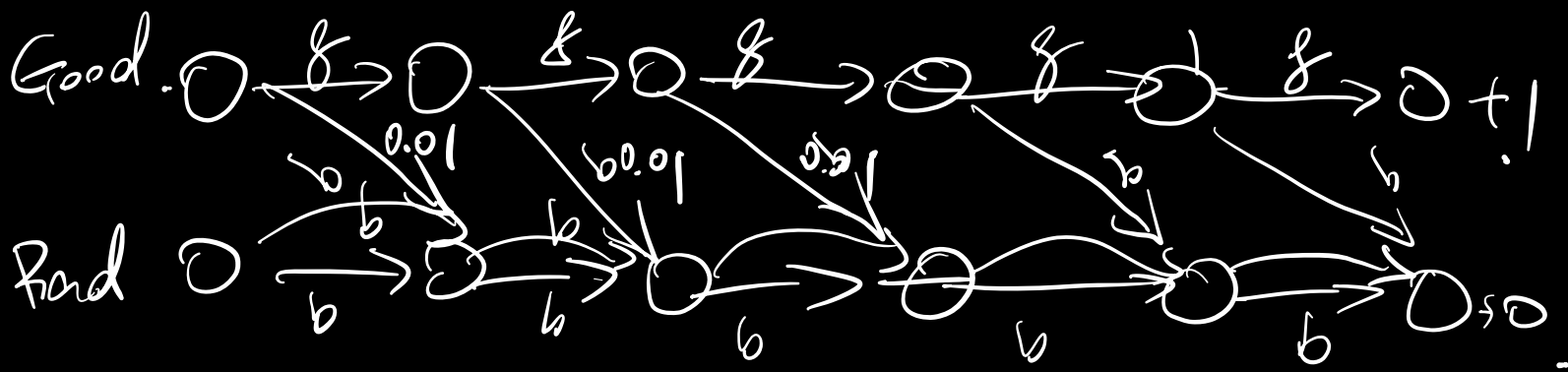


Rmax Exploration "Combination Lock".



$$\nabla J(\pi) \frac{1}{1-\gamma} \mathbb{E}_{d^\pi} \left[\sum_{u'} \gamma^u \left(\sum_{u'} \gamma^u \right) \right]$$

Want: sample complexity $\text{poly}(|S|, |A|, \frac{1}{\epsilon}, \frac{1}{\delta})$
 output $\hat{\pi}$ s.t. $J_{\pi^*} - J_{\hat{\pi}} \leq \epsilon \cdot V_{\max}$

Rmax. At any point of exploration.
 Maintain.

- $n(s, a)$: # times see (s, a)
- $n(s, a, s')$: # times (s, a, s') parameter.

Define $K := \{ (s, a) : n(s, a) = m \}$.

Build MDP \hat{M}_K :

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

- $\hat{R}_K(s, a) = \sum R(s, a) \text{ if } (s, a) \in K$

$P_K(s', a) = \begin{cases} P(s', a) & \text{if } (s, a) \in K \\ R_{\max} & \text{otherwise} \end{cases}$

Collect next episode w/ $\pi_{M_K}^*$.

Optimism in face of uncertainty

Analysis: Define M_K similar to \hat{M}_K .

| | M | M_K | \hat{M}_K |
|----------|-------|---------------------|---------------------|
| Known(K) | $= M$ | $= M$ | $\approx M$ |
| Unknown | $= M$ | R_{\max} -loop | R_{\max} -loop |

except $P_K(s'|s, a) = P(s'|s, a) \forall (s, a) \in K$.

1. Optimism: $\forall \pi: S \rightarrow \Delta(A), V_{M_K}^\pi(s) \geq V_M^\pi(s)$.

2. (Induced Ineq)

Given M_1, M_2 , agree on $K \subseteq (S \times A)$.

Define $escape_K(\tau) = \begin{cases} 1 & \text{if } \tau \text{ visits } (s, a) \notin K \\ 0 & \end{cases}$

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

$|J_{M_1}(\pi) - J_{M_2}(\pi)| \leq V_{\max} \cdot P_{M_1}^\pi[escape(\tau)]$

$\mathbb{P}_M^{\pi_M^*}[\text{escape}(\tau)] \geq \underline{O(\varepsilon)}$ \rightarrow suboptimality.

MSA
 ε

$$\begin{aligned} \frac{\varepsilon V_{\max}}{\Delta} &< J_M(\pi_M^*) - J_M(\pi_{\hat{M}_K}^*) \\ &\leq J_{M_K}(\pi_M^*) - J_M(\pi_{\hat{M}_K}^*) \\ &\leq J_{M_K}(\pi_{\hat{M}_K}^*) - J_M(\pi_{\hat{M}_K}^*) \\ &\leq J_{M_K}(\pi_{\hat{M}_K}^*) - J_M(\pi_{\hat{M}_K}^*) + \underbrace{O\left(\frac{1}{m}\right)}_{\text{small}} \\ &\leq V_{\max} \cdot \underbrace{\mathbb{P}_M^{\pi_{\hat{M}_K}^*}}_{\Delta} [\text{escape}_K(\tau)] + \underbrace{\text{small}}_{\Delta} \\ &\geq \frac{\varepsilon V_{\max}}{2} \end{aligned}$$

$\Rightarrow \mathbb{P}[\text{escape}] \geq \frac{\varepsilon}{2}$

$\frac{\varepsilon V_{\max}}{2}$

Proof of induced norm.

$$\frac{|\bar{J}_{M_1}(\pi) - \bar{J}_{M_2}(\pi)|}{1-\gamma} = \frac{f(s_0) - \bar{J}_{M_1}(\pi)}{1-\gamma} = \frac{\mathbb{E}_{d_{M_1}^\pi} [f - \bar{T}^\pi f]}{1-\gamma}$$

$$= |\bar{J}_{M_1}(\pi) - Q_{M_2}^\pi(s_0, \pi)|.$$

$$= \frac{1}{1-\gamma} \left| \mathbb{E}_{d_{M_1}^\pi} [Q_{M_2}^\pi - \bar{T}_{M_1}^\pi Q_{M_2}^\pi] \right|.$$

$$\sum_{(s,a) \in K} d_{M_1}^\pi(s,a) \cdot (Q_{M_2}^\pi - \bar{T}_{M_2}^\pi Q_{M_2}^\pi)$$

$$\sum_{(s,a) \in K} d_{M_1}^\pi(s,a) \cdot O(V_{max})$$