

Exploration | RL settings according to how to interact w/ env.

→ batch: no interactions. given data (collected by someone else)

→ planning: have access to a "generative model".  
i.e. can query  $s' \sim P(\cdot | s, a) \forall s, a$ .  
"sample complexity" → "query complexity".

MCTS.  
tabular.

→ online: can interact w/ env. in terms of trajectories.  
weaker than "gen model".

online setting: learner can interact w/ env.

wit. learn  $\epsilon$ -optimal policy: how many trajectories need to be drawn?

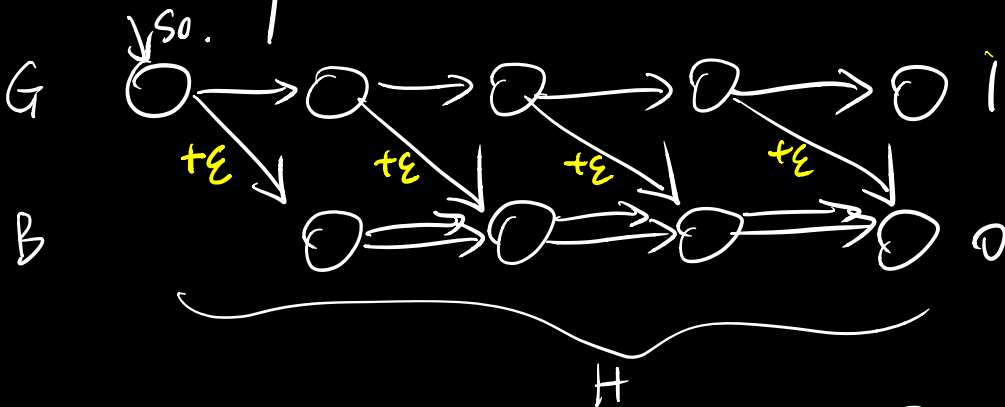
Typical exploration strategy:

1. learning component
2. exploration component

DQN +  $\epsilon$ -greedy exploration  
 $\uparrow$   
 "learning component"  $\left\{ \begin{array}{l} \text{w.p. } 1-\epsilon, \pi_Q \\ \text{w.p. } \epsilon, \text{unif}(A) \end{array} \right.$

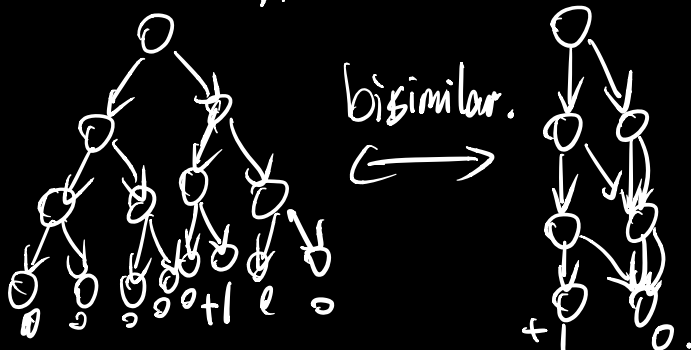
Unif exploration has exponential sample complexity:

Counterexample: "Combination lock".



$\mathcal{P}(1/2^4)$  prob.  
 + greedy doesn't.  
 { Q if Q is init as 0, then no update.  
 ② "anti-shaping".

Fun fact:



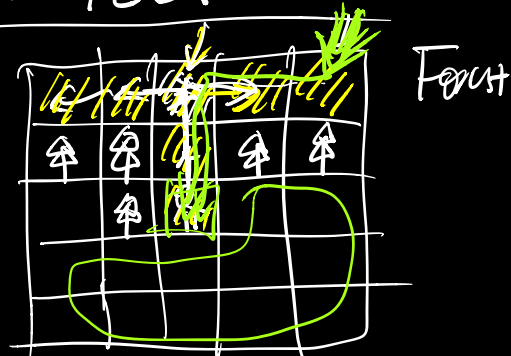
Q: Can we provably explore w/ poly sample complexity.  
at least in tabular RL? A: YES!

Formal setup:  $M = (S, A, P, R, \gamma, d_0)$ .

episodic: all traj. terminates in  $H$  steps.

find  $\hat{\pi}$ , s.t.  $J(\pi^*) - J(\hat{\pi}) \leq \epsilon \cdot V_{\max}$ .

$$J(\pi) := \mathbb{E}_{s \sim d_0} [V^\pi(s)].$$



$$V_{\max} = \frac{R_{\max}}{1-\gamma}$$

$R_{\max}$ -alg: alg maintains:  $\forall s, a, s', \begin{cases} n(s, a) & \text{(initially: 0)} \\ n(s, a, s') \end{cases}$

hyperparam:  $m$ . ("threshold").

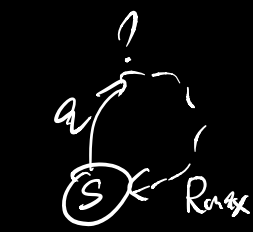
$$K := \{ (s, a) \in S \times A : \underbrace{n(s, a)} = m \}$$

"optimism in face of uncertainty" (OFU)

1. Build MDP  $\hat{M}_K: \forall s, a, s'$

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

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



Stop criterion?

2. Collect a traj.  $s_1, a_1, r_1, \dots, s_H, a_H, r_H$  using  $\pi_{\hat{M}_K}^*$

3.  $\forall h$  s.t.  $n(s_h, a_h) < m, \Rightarrow n(s_h, a_h)++ , n(s_h, a_h, s_{h+1})++$

Analysis.

$$d_M^\pi(s, a) := (1-\gamma) \sum_{h=1}^{\infty} \gamma^{h-1} \mathbb{P}_M [s_h = s, a_h = a | \pi].$$

Define  $M_K$  as "expected ver" of  $\hat{M}_K$ .

$$P_K(s' | s, a) = \begin{cases} P(s' | s, a), & \text{if } (s, a) \in K \\ \mathbb{I}(s' = s) & \text{o.w.} \end{cases} \quad R_K(s, a) = \begin{cases} R(s, a), & \text{if } (s, a) \in K \\ R_{\max}. & \text{o.w.} \end{cases}$$

Def:  $\Delta(M_1, M_2) := \max_{s,a} \|P_1(s,a) - P_2(s,a)\|_1$ .

	M	M <sub>K</sub>	$\hat{M}_K$
K	=M	=M	$\approx M$
unknown	=M	loop	loop

Claim:  $\Delta(M_K, \hat{M}_K)$  is "small" (as a function of m).

Lemma: fix (s,a) where  $n(s,a) = m$ , w.p  $\geq 1 - \delta$ .

$$\| \hat{P}_K(\cdot|s,a) - P_K(\cdot|s,a) \|_1 \leq \sqrt{\frac{2}{m} \log \frac{2(2^{|\mathcal{A}|})}{\delta}}$$

$$\frac{M_K \approx \hat{M}_K}{\Delta}$$

Coro: w.p  $\geq 1 - \delta$ ,  $\Delta(M_K, \hat{M}_K) \leq \sqrt{\frac{2}{m} \log \frac{2 \cdot 2^{|\mathcal{A}|} \cdot |\mathcal{S}| |\mathcal{A}|}{\delta}}$

Lemma (optimism):  $\forall \pi: S \rightarrow A, J_{M_K}(\pi) \geq J_M(\pi)$ .

Further implications of  $\hat{M}_K \approx M_K$ .

"Simulation Lemma"

①  $\forall \pi, |J_{\hat{M}_K}(\pi) - J_{M_K}(\pi)| \leq \Delta(M_K, \hat{M}_K) \cdot \frac{V_{max}}{2(1-\gamma)}$

②  $\|V_{M_K}^* - V_{\hat{M}_K}^*\|_{\infty} \leq \Delta(M_K, \hat{M}_K) \cdot \frac{V_{max}}{2(1-\gamma)}$

Let  $T_K$  and  $\hat{T}_K$  be the Bellman op in  $M_K, \hat{M}_K$  resp.

$$\|V_{M_K}^* - V_{\hat{M}_K}^*\|_{\infty} = \|V_{M_K}^* - \hat{T}_K V_{M_K}^* + \hat{T}_K V_{M_K}^* - \hat{T}_K V_{\hat{M}_K}^*\|_{\infty} \leq \|V_{M_K}^* - \hat{T}_K V_{M_K}^*\|_{\infty} + \gamma \|V_{M_K}^* - V_{\hat{M}_K}^*\|_{\infty}$$

$$\|V_{M_K}^* - \hat{T}_K V_{M_K}^*\|_{\infty} = \|T_K V_{M_K}^* - \hat{T}_K V_{M_K}^*\|_{\infty}$$

$$= \gamma \max_{s,a} \left| \mathbb{E}_{s' \sim P_K(\cdot|s,a)} [V_{M_K}^*(s')] - \mathbb{E}_{s' \sim \hat{P}_K(\cdot|s,a)} [V_{M_K}^*(s')] \right|$$

$$= \gamma \max_{s,a} \left| \langle P_K(s,a) - \hat{P}_K(s,a), V_{M_K}^* - \frac{V_{max}}{2} \cdot \mathbf{1} \rangle \right|$$

$$\leq \gamma \max_{s, a} \left( \left\| P_K(s, a) - \hat{P}_K(s, a) \right\|_{\mathbb{R}^S} \cdot \frac{V_{\max}}{2} \right)$$

$$= \gamma \cdot \Delta(M_K, \hat{M}_K) \cdot \frac{V_{\max}}{2}$$

$$\|T_1 - T_2\| = \max_i |T_i - T_{i+1}|$$

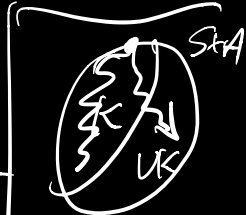
Key Lemma: "Induced Lags"  $\forall \pi: S \rightarrow A$

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

$\mathbb{P}_M[\cdot | \pi]$  considers dist. of rand. traj. generated in  $M$ .

$\tau$  is rand. traj.  $\tau = (s_1, a_1, s_2, a_2, \dots, s_h, a_h)$

$\text{escape}_K(\cdot) = 1$  if  $\exists h, s_h, a_h \in K$ .



Proof:  $J_{M_K}(\pi) = \sum_{\tau: \text{escape}_K(\tau)=1} \mathbb{P}_{M_K}[\tau | \pi] R(\tau)$

$$+ \sum_{\tau: \text{escape}_K(\tau)=0} \mathbb{P}_{M_K}[\tau | \pi] R(\tau)$$

$$R(\tau) = \sum_{h=1}^H \gamma^{h-1} R(s_h, a_h)$$

$$= \sum_{\text{escape}_K(\tau)=1} \mathbb{P}_{M_K}[\tau | \pi] \cdot (R(\text{pre}_K(\tau)) + R(\text{suf}_K(\tau)))$$

$$+ \sum_{\text{escape}_K(\tau)=0} \mathbb{P}_{M_K}[\tau | \pi] R(\tau)$$



$$\leq \sum_{\text{escape}_K(\tau)=1} \mathbb{P}_{M_K}[\tau | \pi] \cdot (R(\text{pre}_K(\tau)) + V_{\max})$$

$$+ \sum_{\text{escape}_K(\tau)=0} \mathbb{P}_{M_K}[\tau | \pi] R(\tau)$$

$$= \sum_{\text{pre}_K(\tau)} \mathbb{P}_{M_K}[\text{pre}_K(\tau) | \pi] (R(\text{pre}_K(\tau)) + V_{\max}) +$$

$$\underline{J_M(\pi)} \equiv \sum_{\text{pre}_k(\tau) \in R(\text{pre}_k(\tau))} \mathbb{P}_M[\text{pre}_k(\tau) | \pi]$$

$$+ \sum_{\text{escape}_k(\tau)=0} \mathbb{P}_M[\tau | \pi] \cdot R(\tau)$$

$\text{pre}_k(\tau) = s_1, a_1, \dots, s_n, a_n$   
 s.t.  $\forall t \leq n-1, (s_t, a_t) \in K$   
 &  $(s_n, a_n) \notin K$

$$\mathbb{P}_{M_k}(\text{pre}_k(\tau) | \pi) = d_0(s_1) \cdot \pi(a_1 | s_1) \cdot P_k(s_2 | s_1, a_1) \cdot \pi(a_2 | s_2) \cdot \dots \cdot P_k(s_n | s_{n-1}, a_{n-1}) \cdot \pi(a_n | s_n)$$

Claim:  $\mathbb{P}_{M_k}(\text{pre}_k(\tau) | \pi) = \mathbb{P}_M[\text{pre}_k(\tau) | \pi]$

$\forall \text{escape}_k(\tau)=0, \mathbb{P}_{M_k}[\tau | \pi] = \mathbb{P}_M[\tau | \pi]$

$$J_{M_k}(\pi) - J_M(\tau) \leq \sum_{\text{pre}_k(\tau) \in \Delta} \mathbb{P}_M[\text{pre}_k(\tau) | \pi] \cdot V_{\max}$$

$$= \mathbb{P}_M[\text{escape}(\tau) | \pi]$$

Sample complexity analysis.  $\forall$  episode, either of following happens.

- ①.  $\pi_{M_k}^*$  is  $\epsilon$ -optimal  $\leftarrow \checkmark$  "terminate - or - explore"
- ②.  $\pi_{M_k}^*$  escapes w/ significant prob.  $\Delta$

Proof: w.t.s.  $\neg \text{①} \Rightarrow \text{②}$ . optimism.

$$\begin{aligned} \epsilon \cdot V_{\max} &\leq \underline{J_M(\pi^*)} - J_M(\pi_{M_k}^*) \leq \underline{J_{M_k}(\pi^*)} - J_M(\pi_{M_k}^*) \\ &\leq J_{M_k}^* - J_M(\pi_{M_k}^*) \\ &\leq \underline{J_{M_k}^*} + \underbrace{\Delta(M_k, \hat{M}_k)}_{\Delta} \cdot \frac{V_{\max}}{2(1-\gamma)} - J_M(\pi_{M_k}^*) \\ &\leq \underline{J_{M_k}(\pi_{M_k}^*)} + \Delta(M_k, \hat{M}_k) \cdot \frac{V_{\max}}{1-\gamma} - \underline{J_M(\pi_{M_k}^*)} \end{aligned}$$

$$\leq \underbrace{\Delta(M_k, \hat{M}_k)}_{\Delta} \cdot \frac{V_{\max}}{1-\gamma} + V_{\max} \cdot \underbrace{P_M[\text{escape}_k(\tau) | \pi_{\hat{M}_k}^*]}_{\Delta}$$

$$\max_{s, a \in K} \|P_k(\cdot | s, a) - \hat{P}_k(\cdot | s, a)\|_1 \leq \sqrt{\frac{2}{m} \log(\dots)}$$

Outline of remainder of proof:

Set  $m$  large enough, s.t.  $\Delta(M_k, \hat{M}_k) \cdot \frac{V_{\max}}{1-\gamma}$

$$\Rightarrow P_M[\text{escape} | \pi_{\hat{M}_k}^*] \geq \frac{\epsilon}{2} \quad \left| \leq \frac{\epsilon}{2} \cdot V_{\max} \right.$$

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

crude calc:  $m \cdot |S| \times |A| \cdot \left(\frac{\epsilon}{2}\right)$

Remarks: Q can terminate if  $\underbrace{J_M(\pi_{\hat{M}_k}^*)}_{\Delta} \approx \underbrace{J_{\hat{M}_k}^*}_{\Delta}$

②  $R_{\max}$  doesn't necessarily explore.  $\leftarrow$  can estimate from Monte-Carlo  
 $\rightarrow$  "reward-free exploration"

③ "hard-to-reach" states don't bother us.