# Exploration and Learning in "Tabular" MDPs:

$\Delta(x) = \{$ probability dist. on $x \}$

What is an MDP? (finite horizon $H$, episodic)

$$( S, A, P_h, R_h, d_0 )$$

(for $h = 1 \dots$ to $H$)

$P_h : S \times A \to \Delta(S)$

$R_h : S \times A \to \Delta(\mathbb{R})$

Policy $= \pi_h : S \longrightarrow \Delta(A)$  for $h=1$ to $H$.

$\pi^* =$ optimal policy

Tabular : $|S|$ and $|A|$ are not too big

Given MDP, find $\pi^*$ via Bellman Equations.

$$Q_h^*(s,a)$$

$$||$$

$$\mathbb{E}_{\substack{s_{h+1} \sim P_h(s,a) \\ r_h \sim R_h(s)}} \left[ r_h + V_{h+1}^*(s_{h+1}) \right]$$

$$V_h^*(s) = \max_a Q_h^*(s,a)$$

$$\pi_h^*(s) = \text{argmax}_a Q_h^*(s,a)$$

Problem for Reinforcement Learning:

oftentimes $P_h$, $R_h$ are unknown.

BUT: given $Q_h^*$, can compute $\pi_h^*$.

Q: How do we find $Q_h^*$? (or $\pi_h^*$?)

Today: Tabular MDP where $P_h, R_h$ etc unknown (e.g. $Q_h^*$)

$R_h$: reward function is unknown,
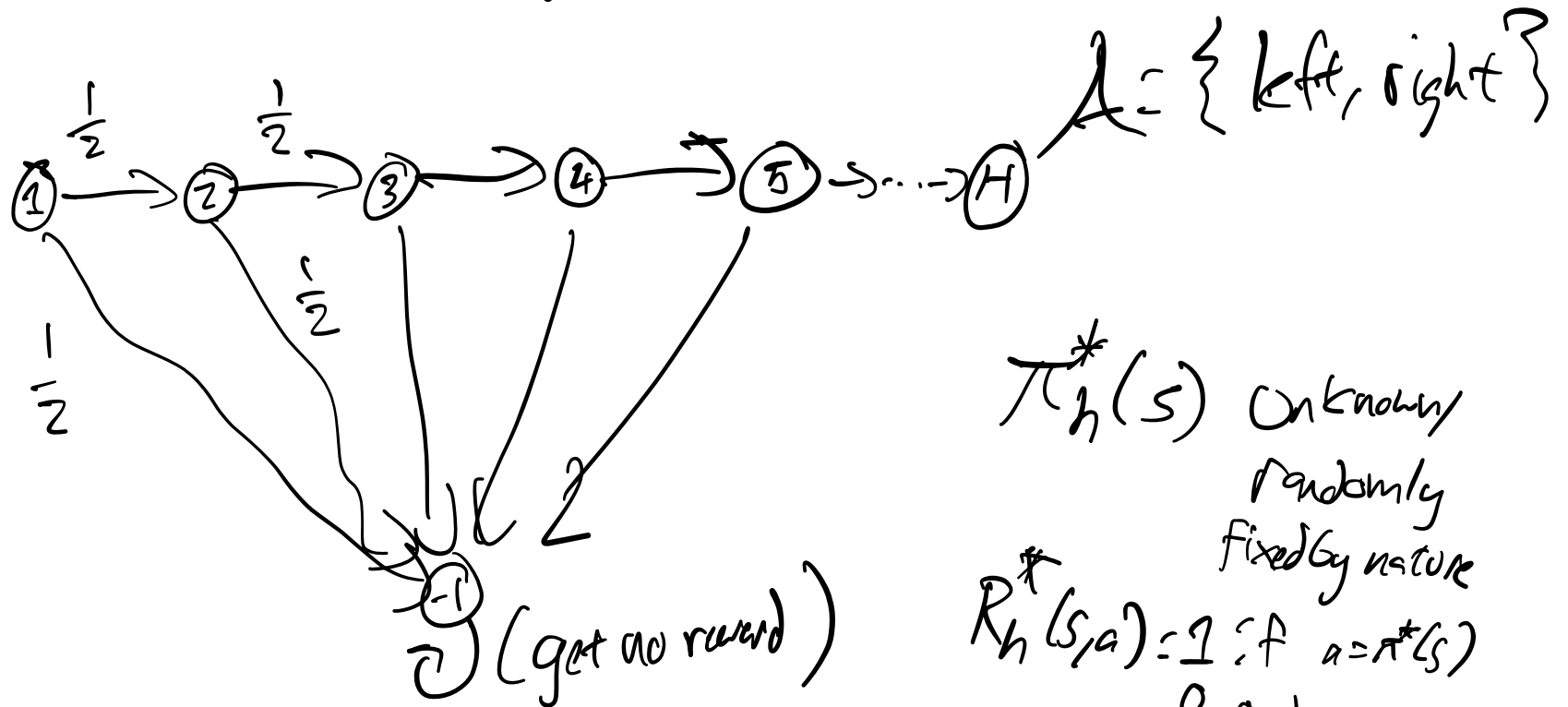   but can usually deal with it using UCB.
      (like bandits)

So assume $R_h$ is known. (and $d_0$)

For Tabular MDP, main issue is $P_h$ is unknown.

Want to learn $P_h$ by playing many episodes.

Q: Can we learn $P_h(s, a)$ for all $(s, a) \in S \times A$?

# P: Some states may be hard to reach.



$A = \{ \text{left, right} \}$

$\pi_h^*(s)$ unknown
randomly
fixed by nature

$R_h^*(s,a) = 1$ if $a = \pi^*(s)$
$\phantom{R_h^*(s,a) =} 0$ o.w.

for every $s \geq 1$, you pull either left or right arm
and if you pull "correct one"

For any policy
$$Pr(s_H = H) = \left(\frac{1}{2}\right)^{H-1}.$$

# RMAX algorithm:

"Optimism under uncertainty"

- Maintain $K \subseteq S \times A$

$$K = \{(s,a): n(s,a) \geq N\} \xleftarrow{\text{parameter}}$$

$r_h \in [0,1]$ almost surely

$\gg$ #states $= |S|$

- If $(s,a) \notin K$, (we assume it is "maximally good" (optimism)

it pretend $Q_h^*(s,a) = "H - h + 1"$.

Rmax algorithm: For episode $t=1$ to $T$:

- Compute $Q_t^*$ under the (optimism) assumption for $(s,a) \notin K$.

For $(s,a) \in K$:
Estimate $\hat{P}_h(s,a)$ from past experience.

- Play optimal strategy according to $Q_t^*$.
- $n(s,a) = n(s,a)+1$ for any $(s,a)$ tried in this episode.

$$Q_{t,h}^*(s,a) = \mathbb{E}\left[ r_h + V_{t,h+1}^*(s_{h+1}) \right]$$

$r_h \sim R_h(s,a)$

$s_h \sim \hat{P}_h(s,a)$

$\hat{P}_h(s_h, a_a) = $ Observed distribution over $S$ from past.

FOR $(s,a) \in K$.

$$V_{t,h}^*(s) = \max_a Q_{t,h}^*(s,a)$$

(otherwise, $Q_{t,h}^*$ is optimistic).

Play $\pi_t^* \Longrightarrow$ observe $(s_1, a_1), (s_2, a_2) \dots (s_H, a_H)$

Then $n(s_i, a_i) = n(s_i, a_i) + 1$ for all $i \in [H]$

$m(s_i, a_i, s_{i+1}) = m(s_i, a_i, s_{i+1}) + 1$

$n(s,a) = 0 \quad \forall S \times A$ initially.

$m(s, a, s') = 0$ initially

Really $K \subseteq \{ (s,a,h) : n(s,a,h) \geq N \}$

$n(s_h, a_h, h) = n(s_h, a_h, h) + 1$

$m(s_h, a_h, h, s_{h+1}) = m(s_h, a_h, h, s_{h+1}) + 1$

$$\hat{P}_h(s,a)(s_{h+1}) = \frac{m(s,a,h,s_{h+1})}{\sum_{s'} m(s,a,h,s')}$$

__Thm:__ With probability at least 99%,
after $T = poly(|S|, |A|, H, 1/\epsilon)$ plays
$\pi_T^*$ is $\epsilon$-optimal.
$\left( f(\pi_T^*) \geq f(\pi^*) - \epsilon \right)$.

# Pf sketch of Thm:

At each time $t$, either

(A) $\pi_t^*$ has $\gg \frac{\varepsilon}{H}$ probability of

escaping $\mathcal{K}$

or

(B) $\pi_t^*$ is $\varepsilon$-optimal.

(B) $\leadsto$ Done. (by optimism + $\mathcal{K}$ being well-observed)

(A) $\leadsto$ after $\gtrsim \frac{H}{\varepsilon}$ rounds, some $n(s,a,h)$ some

increases by $1$ for $(s,a,h) \notin \mathcal{K}$.

SO $\quad n(s_a,h) < N \iff (s,a,h) \in K$
$$N = \frac{H^2 |S| \log |S||A|H}{\varepsilon^2}$$

case Ⓐ can only occur roughly

$\simeq \frac{H}{\varepsilon} |S||A|H$ many times $\qquad d_{TV}(\hat{P}_h(s,a),$
$$P_h(s,a)) \leq \frac{\varepsilon}{H}$$

(after this many times, $n(s_a,h) \geq N$

for all $(s,a,h)$,

i.e $K = \phi$)

Rigorous pf idea: ① $Q^*_t = Q^{*, M_t}$ is the $Q^*$ function

for "optimistic MDP" $M_t$.

② If not in case Ⓐ then w.p. $\simeq 1-\varepsilon$ optimistic MDP and true MDP behave same.

Example of soluble non-Tabular MDPs?

$$S = \mathbb{R}^d$$

action $u_h$ instead of an

Linear dynamics.

$$S_{h+1} = A S_h + B u_h + \text{noise}$$

$u_h =$ "control input" given by policy

Objective:

$$R_h(S_h) = \frac{1}{2} \langle S_h, M_h, S_h \rangle$$

objective

Fact: $u_h = K_h S_h$.    LQR.