# Intro to RL + MDPs

- HW 2 should be out tonight.
- Haifeng is around.

What is a "finite horizon $H$ episodic

Markov Decision Process"?

$H=1 \Longleftrightarrow$ roughly contextual bandit.

What is such an MDP?   $H = $ horizon

$$S = \{ \text{set of states} \}$$

$$\Delta(S) = \{ \text{set of probability distributions over } S \}$$

$$A = \{ \text{set of actions} \}$$

Time steps   $h = 1$ to $H$.

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

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

$d_1 \in \Delta(S)$ distribution over initial states.

How do we "play" a single episode of the MDP?

User Selects a policy

$$\pi = (\pi_1, \dots, \pi_{|H|})$$

$$\pi_h : S \longrightarrow \Delta(A)$$

$\uparrow$ from 1 to H.

Then: $s_1 \sim d_1$

We select action $a_1 \sim \pi_1(s_1)$

Observe $r_1 \sim R_1(s_1, a_1)$

$s_2 \sim P_1(s_1, a_1)$

for h=2 to H:

$a_h \sim \pi_h(s_h)$

$r_h \sim R_h(s_h, a_h)$

$s_{h+1} \sim P_h(s_h, a_h)$

Total reward in one episode is $\sum_{h=1}^{H} r_h$.

Play episodes $t=1, \ldots, T$.

For $t=1$ to $T$:

- User selects a policy $\pi^t$ based off past experience.
- Play policy $\pi^t$ in the MDP.
- Gain reward $\sum_{h=1}^{H} r_h^t$ $\longleftarrow$ episode $t$.

Total Reward $= \sum_{t=1}^{T} \sum_{h=1}^{H} r_h^t$ (goal: maximize this.)

Regret $=$ Our total reward $- \mathbb{E}[$ reward of $\pi^* ]$

best policy

What is the "best" $\pi^*$?

$$f(\pi) = \mathbb{E} \text{ reward under policy } \pi. = \mathbb{E}_\pi \sum_{h=1}^{H} r_h$$
(in one episode)

$$\begin{bmatrix} s_1 \sim d_1 \\ \text{for } h=1 \text{ to } H: \\ \quad a_h \sim \pi_h(s_h) \\ \quad r_h \sim R_h(s_h, a_h) \\ \quad s_{h+1} \sim P_h(s_h, a_h). \end{bmatrix}$$
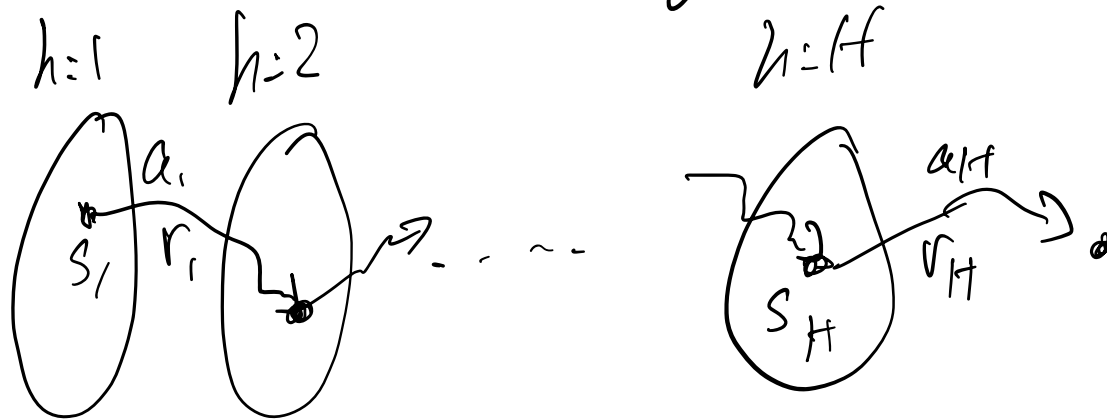
$$\pi^* = \arg\max_\pi f(\pi).$$

<u>Remark:</u> Getting $o(T)$ regret          "Online RL"

$\uparrow$

Finding a nearly optimal policy.          "PAC-RL"

Next: Bellman equations/ Dynamic Programming

$h=1$    $h=2$                    $h=H$



Def   Value function under policy $\pi$

$$V_h^\pi(s) = \mathbb{E}_\pi\left[\sum_{\ell=h}^{H} r_\ell \mid s_h = s\right]$$

Q-function under policy $\pi$

$$Q_h^\pi(s,a) = \mathbb{E}_\pi\left[\sum_{\ell=h}^{H} r_\ell \mid s_h = s, a_h = a\right]$$

Def: $V_h^*(s) = V_h^{\bar{\pi}^*}(s)$

$Q_h^*(s,a) = Q_h^{\bar{\pi}^*}(s,a).$

$Q_H^*(s,a) = \mathbb{E}[r_H \mid s_H = s, a_H = a]$

$V_H^*(s) = \mathbb{E}_{\pi^*}[r_H \mid s_H = s]$

$= \max_a \mathbb{E}[r_H \mid s_H = s, a_H = a]$

$= \max_a Q_H^*(s,a).$

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

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

$$Q_h^*(s,a) = \mathbb{E}_{\pi^*}\left[\sum_{\ell=h}^{H} r_\ell \mid s_h=s, a_h=a\right]$$

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

"Bellman Equations"

Dynamic Programming algorithm to compute $\pi^*, V^*,$ and $Q^*$:

($\leftarrow$ "Value iteration")

① $\quad Q_H^*(s,a) = \mathbb{E}\left[V_{H}^* \mid s_H = s, a_H = a\right]$

$V_H^*(s) = \max_a Q_H^*(s,a)$

$\pi_H^*(s) = \operatorname*{argmax}_a Q_H^*(s,a)$

② Use Bellman equations to define
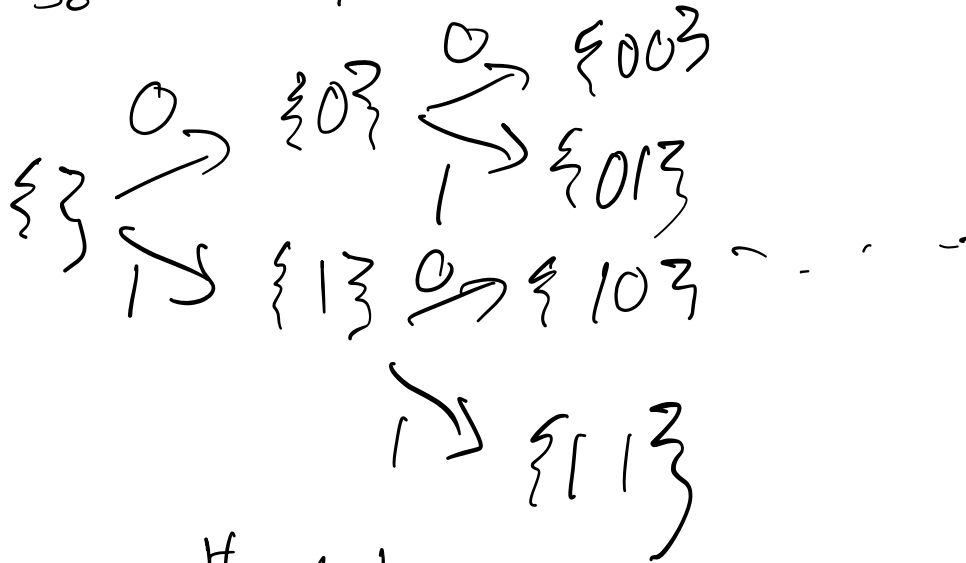
$Q_{H-1}^*, V_{H-1}^*, \pi_{H-1}^*$

Then $\quad Q_{H-2}^*, V_{H-2}^*, \pi_{H-2}^*, \ldots \quad Q_1^*, V_1^*, \pi_1^*.$

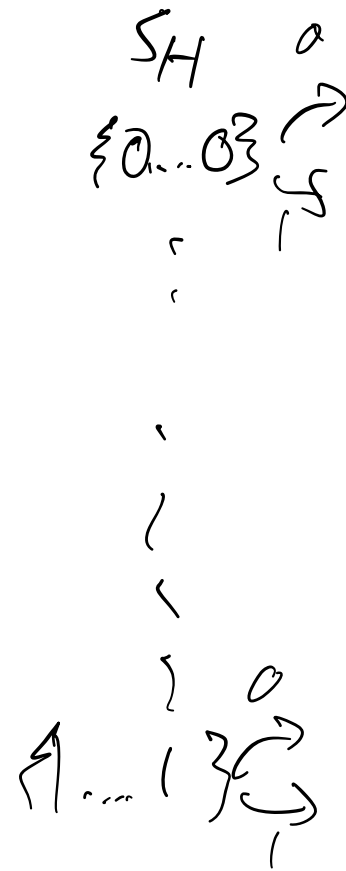① Failure of naive exploration.

② Difficulty of large state spaces.

"Good configuration lock"

$S_0$     $S_1$     $S_2$

$\{\} \xrightarrow{0} \{0\} \xrightarrow{0} \{00\}$

$\xrightarrow{1} \{01\}$

$\xrightarrow{1} \{1\} \xrightarrow{0} \{10\}$

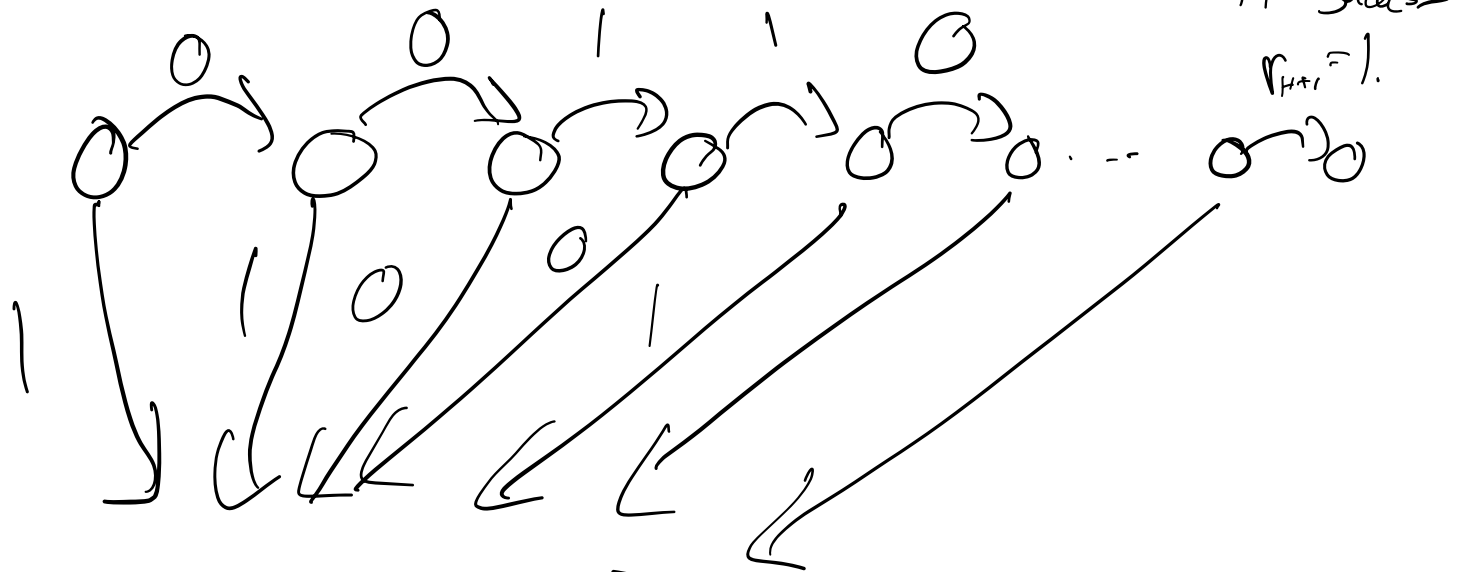$\xrightarrow{1} \{11\}$

Requires $2^H \simeq |S|$
many plays to get reward with good probability.

$r_h = 0 \quad h < H \quad$ password
$r_H = 1$ if $\quad l$
you played 0011001l ....
(or some other fixed str)

$S_H \quad 0$
$\{0...0\} \circlearrowright$

$\{1...1\} \xrightarrow{0} \circlearrowright$

(11) "Bad combination lock"

$s_0$

$s_H$ SUCCESS

$r_{H+1} = 1$.

Failure