

What is RL? \rightarrow learning what to do in an environment

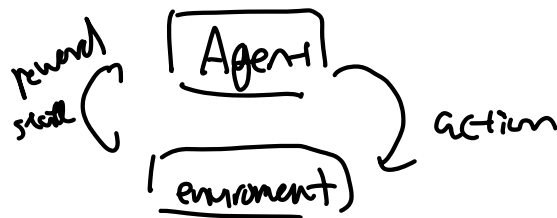
Exploration \checkmark Exploitation

Explore: trying out new things

Exploitation: doing things you already know will work well.

how?

\rightarrow we want to maximize Reward

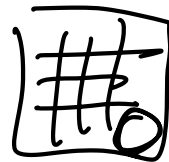


① start at some state $S_t \leftarrow \text{discrete} \rightarrow t = 0, 1, 2, \dots$

② pick Action $A_t \leftarrow$

③ receive S_{t+1}, R_{t+1}

④ repeat



S
 \rightarrow
states

$A(S)$
 \downarrow
actions

$R \subseteq \mathbb{R}$
 \uparrow
rewards

A MDP \rightarrow Markov decision process

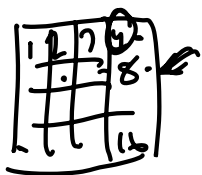
- S
- $A(s) \rightarrow$
- transition probabilities



$S_0 \rightarrow A_0 \rightarrow R_1 \rightarrow S_1 \xrightarrow{A_1} R_2 \rightarrow S_2 \dots \dots$ trajectory

R_t & $S_t \leftarrow$ random variables \rightarrow discrete probability distributions

$p(s', r | s, a)$: given that I'm at state s , if I take action a , what is the probability that I will end up in s' , w/ reward r



$$p(8, 50 | 4, \downarrow) = 1$$

$$\hookrightarrow p(16, 20 | 3, \downarrow) = 0$$

$\sum_{r \in R} p(s', r | s, a)$ = probability of going from state $s \rightarrow s'$ taking action a .

transition probability = $p(s' | s, a)$

Discounted rewards

discount factor $0 \leq \gamma \leq 1$

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots$$

\rightarrow Sum of rewards

maximize $\mathbb{E}[G_t]$

policy & value function

$\pi : \text{state} \rightarrow \text{action}$

$V_\pi(s) = \mathbb{E}[G_t | S_t = s] \rightarrow V_{\pi^*} \geq V_\pi$
state - value function

$q_\pi(s, a) = \mathbb{E}[G_t | S_t = s, A_t = a] \quad q_{\pi^*}(s, a) \geq q_\pi(s, a)$
action-value function

$$V_\pi(s) = \mathbb{E}[G_t | S_t = s]$$



$$\sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_\pi(s')]$$

state-value Bellman equation

iteratively



optimal (maximize)

$$V_{\pi^*} \geq V_\pi \quad \forall \pi$$