What is RL?

- Agent & Environment



$$\mathbb{P}[S_{t+1} \mid S_t]$$
$$= \mathbb{P}[S_{t+1} \mid S_1, \dots S_t]$$

- RL agent components
  - ↳ policy
  - ↳ value function
  - ↳ model $\Big\langle$ ⟶ model - based
    ⟶ model - free
  - ↳ prediction & control

- Markov decision processes (MDP)
  - ↳ $\langle S, P, R, A, \gamma \rangle$ — discount factor $\in [0,1]$
    - states
    - reward $R_s = \mathbb{E}[R_{t+1} \mid S_t = s]$
    - actions
    - state transition probability

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

  - ↳ state - value function

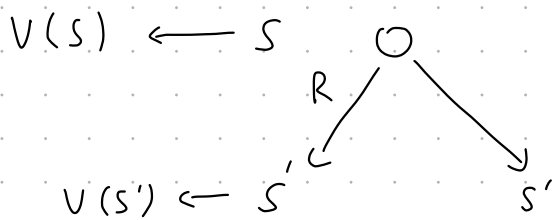    $$V(S) = \mathbb{E}[G_t \mid S_t = s] = \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s\right]$$
    $$= \mathbb{E}[R_{t+1} + G_{t+1}] = \mathbb{E}[R_{t+1} + \gamma V(S_{t+1}) \mid S_t = s]$$

    Bellmann equation

Nice fact: Bellmann equation is recursive & linear

$V(s) \longleftarrow s$



$V(s') \longleftarrow s'$

$$V(s) = R_{t+1} + \gamma \sum_{s' \in S} P_{ss'} V(s')$$

↳ policy $\pi$:

$$\pi(a \mid s) = \mathbb{P}[A_t = a \mid S_t = s]$$

$$V_\pi(s) = \mathbb{E}_\pi[G_t \mid S_t = s] \qquad \text{value - function}$$

$$q_\pi(s,a) = \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a] \qquad \text{action - value function}$$
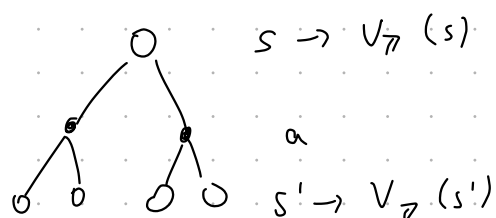
↳ return to Bellmann equation

$$V_\pi(s) = \mathbb{E}_\pi[R_{t+1} + \gamma V_\pi(S_{t+1}) \mid S_t = s]$$

$$q_\pi(s,a) = \mathbb{E}_\pi[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a]$$

$V_\pi(s) \longleftarrow s$



$q_\pi(s,a) \quad a$

$s, a \to q_\pi(s,a)$

$s' \to V_\pi(s')$

$$V_\pi(s) = \sum_{a \in A} \pi(a \mid s) \, q_\pi(s,a) \qquad q_\pi(s,a) = R_{t+1}^a + \gamma \sum_{s'} P_{ss'} V_\pi(s')$$

$$= R_{t+1} + \gamma V_\pi(s')$$

⇓ combine



$s \to V_\pi(s)$

$a$

$s' \to V_\pi(s')$

$$V_\pi(s) = \sum_a \pi(a \mid s) \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_\pi(s') \right)$$

$$ex: \left( R_s^a + \gamma V_\pi(s') \right]$$

2) optimality

$$V_*(s) = \max_{\pi} \left( V_\pi(s) \right)$$

$$q_*(s,a) = \max_{\pi} \left( q_\pi(s,a) \right)$$

$\pi_*$ for any MDP

1) there is always an optimal policy $\pi^*$

2) all optimal policies achieve the $V_*$ and $q_*$

$$\pi_*(a|s) = \begin{cases} 1 & \text{if } a = \text{argmax}_a \; q_*(s,a) \\ 0 & \text{else} \end{cases}$$

---

## Dynamic programming

1) optimal substructure

2) overlapping subproblems

$\Bigg\}$ MDP ✓ because of Bellmann equation

- model - based (update based on estimate)

- boohtraps

## Monte Carlo

value = empirical mean return

- model - free (no knowledge of environment required)

- sampling method

- inherently slow (one episode $\hat{=}$ one sample)

## TD learning

- model free
- bootstraps
- online update

$$V(S_t) \leftarrow V(S_t) + \alpha \underbrace{\left( \underbrace{R_{t+1} + \gamma V(S_{t+1})}_{\text{target}} - V(S_t) \right)}_{\delta_t}$$

learning rate

$\Downarrow$

prediction problem

## Control problem

on-policy          vs          off-policy

(on the job)                    (from observing)

$\hookrightarrow$ learn about policy $\pi$

from experience sampled from $\pi$

Ex: SARSA

- apply TD to $q(s, a)$
- every time step, evaluate & improve policy