

## Part VI Markov Decision Process

- MDP is theoretic foundation of reinforcement learning
- RL can be viewed as an approximate solution to MDP.
- MDP is a generalization of contextual bandit in which the next step depends on current action & state.

### Definitions

- Agent : player, learner, ...
- Environment : Thing that interacts with agent.
- State  $S_t$ : Representation of situation at time t
- State space  $S$  : collection of possible states
- Action  $A_t$ : Action taken by agent at time t
- Action space  $A(s)$ : collection of possible actions at state  $s \in S$ .
- Reward  $R_{t+1}$ : Reward received at the beginning of  $t+1$
- Transition probabilities :  $\forall s, s' \in S, a \in A(s), r \in R$

$$P(s', r | s, a) = \Pr(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a)$$

$$\text{satisfying } \sum_{s'} \sum_r P(s', r | s, a) = 1, \quad \forall s \in S, a \in A(s)$$

• Policy  $\pi$ : How agent acts on the basis of St. Specifically,  
 $\pi(a|s)$  defines the prob. that agent picks action  $a \in A(s)$   
given  $S_t = s$ .

MDP formulates a sequential decision making problem.

MDP and agent's policy together give rise a trajectory

$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, \dots$

Question: What is Markov? Why Markov?

Markov property means that given the present, the past  
and future are independent.

How I reach present state doesn't affect my future,

the only thing matters is that I'm at this state.

Mathematically, it says

$$\Pr(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a, S_i^{t-1}, A_i^{t-1})$$

↙ history doesn't matter

$$= \Pr(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a)$$

Is this practical? Yes if we carefully define states.

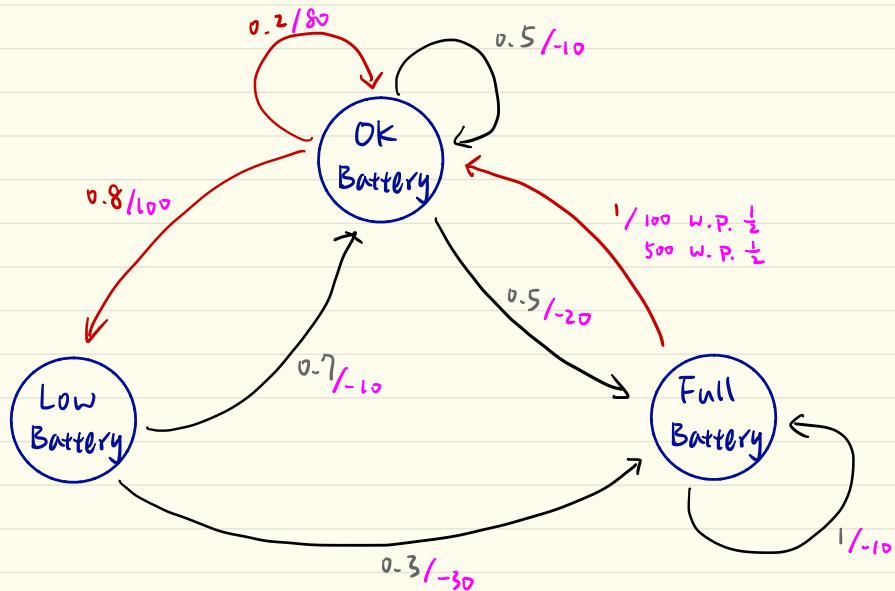
Ex: UberEats with GoGoRo

Red: Take order

Grey : Transition probability

Black: Charge

purple : Reward



Three states:  $S = \{L, O, F\}$

Three action sets:  $A(L) = \{C\}$ ,  $A(O) = \{C, T\}$ ,  $A(F) = \{C, T\}$

Transition probabilities:

$$P(O, -10 | L, C) = 0.7$$

$$P(O, -10 | O, C) = 0.5$$

$$P(F, -10 | F, C) = 1$$

$$P(F, -30 | L, C) = 0.3$$

$$P(F, -20 | O, C) = 0.5$$

$$P(O, 100 | F, T) = \frac{1}{2}$$

$$P(O, 80 | O, T) = 0.2$$

$$P(O, 500 | F, T) = \frac{1}{2}$$

$$P(L, 100 | O, T) \approx 0.8$$

Agent's Goal: The agent's goal in MDP is to maximize not immediate reward but cumulative reward in the long run.

Suppose the trajectory is  $S_0, A_0, R_1, S_1, A_1, R_2, \dots$

the cumulative discounted reward from  $t$  onwards is

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

where  $\gamma \in [0, 1]$  is the discount factor.

Note the recursive form

$$\begin{aligned} G_t &= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

Why discount?

- 1) Mathematically, it makes sure the convergence
- 2) Practically, tradeoff between immediate & future rewards

$\gamma \rightarrow 0$  myopic     $\gamma \rightarrow 1$  farsighted.

Note that  $G_t$  is a random variable evaluated at a outcome trajectory, which is not something we have control over. But we can control the expectation of this RV.

## More Definitions

- State-transition prob

$$P(s' | s, a) = \Pr \{ S_t = s' | S_{t-1} = s, A_{t-1} = a\} = \sum_{r \in R} P(s', r | s, a)$$

- Expected reward for  $(s, a)$

$$r(s, a) = \mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a] = \sum_{r \in R} \sum_{s' \in S} P(s', r | s, a)$$

- Expected reward for  $(s, a, s')$

$$r(s, a, s') = \mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a, S_t = s']$$

$$= \sum_{r \in R} r \cdot P(r | s, a, s') = \sum_{r \in R} r \cdot \frac{P(s', r | s, a)}{P(s' | s, a)}$$

To evaluate a policy  $\pi$ , we define

- Value of state  $s$  under  $\pi$

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right]$$

state-value function for  $\pi$

- Value of action  $a$  in state  $s$  under  $\pi$

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right]$$

action-value function for  $\pi$ , or Q value

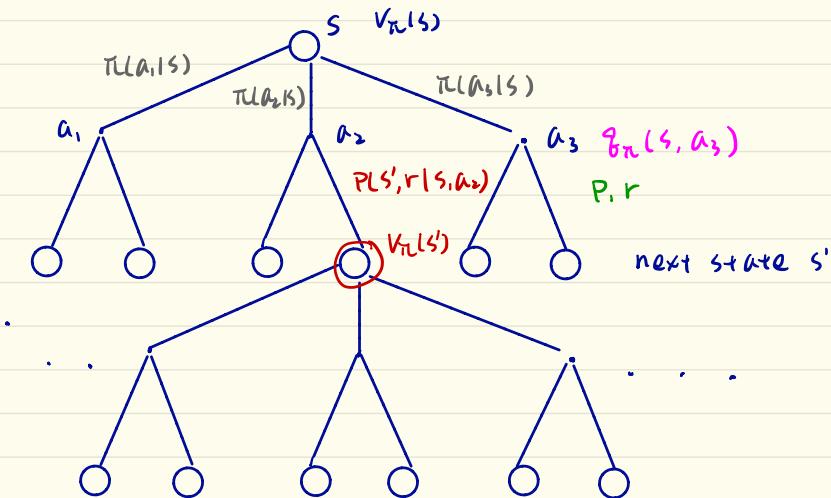
Focus on  $V_{\pi}(s)$  for now:

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

$$= \sum_a \pi(a|s) \cdot \sum_{s'} \sum_r p(s', r | s, a) \cdot [r + \gamma \mathbb{E}_{\pi}[G_{t+1} | S_{t+1} = s']]$$

$$= \sum_a \pi(a|s) \cdot \sum_{s'} \sum_r p(s', r | s, a) \cdot [r + \gamma V_{\pi}(s')] \quad \cdots (4)$$

This is a so called Bellman equation for  $V_{\pi}$



This helps us evaluate a policy. But how to find best policy?

## Bellman Optimality Equation

Optimal state-value function      Optimal action-value function

$$V^*(s) = \max_{\pi} V_{\pi}(s)$$

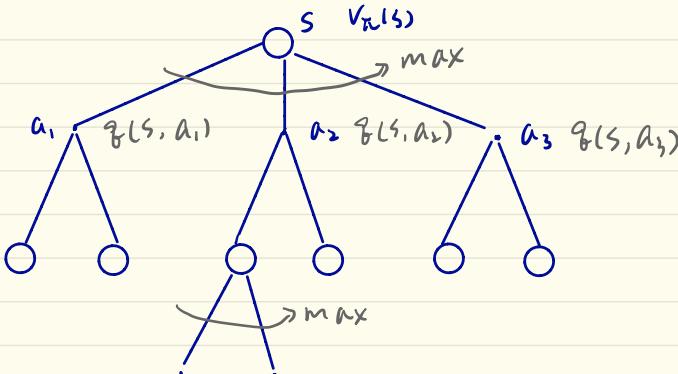
$$q^*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

$$V^*(s) = \max_{a \in A(s)} q^*(s, a) = \max_{a \in A(s)} \mathbb{E}_{\pi^*}[G_t | S_t = s, A_t = a]$$

$$= \max_{a \in A(s)} \mathbb{E}_{\pi^*}[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a]$$

$$= \max_{a \in A(s)} \mathbb{E}[R_{t+1} + \gamma \cdot V^*(S_{t+1}) | S_t = s, A_t = a]$$

$$= \max_{a \in A(s)} \sum_{s', r} P(s', r | s, a) [r + \gamma \cdot V^*(s')]$$



$V^*(s)$  must be equal to the expected return for the best action from that state  $s$ .

Remarks : 1)  $\pi^*$  is not unique but  $V^*$  is unique  
2) It is very difficult (if not impossible) to solve Bellman equation in closed form. There are several ways to approximate it in dynamic programming.

3)  $\max_{\pi \in \Pi(s)} \sum_{s', r} P(s', r | s, a) [r + \gamma \cdot V^*(s')] \quad (\text{**})$  looks very greedy. Why is it optimal? The future expected reward is encoded in (\*\*).

4) Historically, MDP was heavily studied for stochastic optimal control.

In what follows, we introduce two ways to solve Bellman equation by techniques from dynamic programming

## Solving Bellman

Two approaches:

- Policy iteration = Policy evaluation + Policy improvement

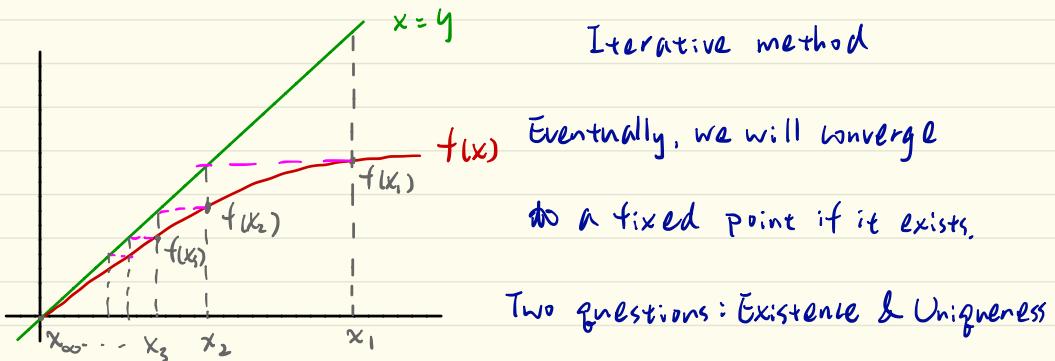
To evaluate a policy  $\pi$ , we need to solve

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

same  $V_\pi(\cdot)$

$$= f(V_\pi)$$

Consider solving  $x = f(x)$ , i.e., a fixed point of  $f(\cdot)$



Two questions: Existence & Uniqueness

Let  $V_k(\cdot)$  be the value that we obtain in the  $k$ -th iteration

under  $\pi$ . Iterative policy evaluation update rule

$$V_{k+1}(s) = \sum_a \pi(a|s) \cdot \sum_{s'} \sum_r p(s', r | s, a) \cdot [r + \gamma V_k(s')]$$

Let  $\pi_k$  be the policy adopted at  $k$ -th iteration

Find an improved policy based on current estimated  $V_{\pi_k}$

Policy improvement:

$$\pi_{k+1} = \arg \max_a \sum_{s', r} P(s', r | s, a) [r + \gamma V_{\pi_k}(s')]$$

Policy iteration:  $\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \rightarrow \dots \xrightarrow{I} \pi^* \xrightarrow{E} V_{\pi^*}$

For finite MDP, since there are only finite number of policies  
this sequence always converges to optimal policy & value.

Drawback: Complexity! Need to run policy iteration many times.

Value iteration:

Directly turning Bellman Optimality equation into update rule.

$$V_{k+1}(s) = \max_a \sum_{s'} \sum_r P(s', r | s, a) \cdot [r + \gamma V_k(s')]$$

We stop once value's improvement become small.

Remark: Solving MDP requires the full knowledge of dynamics.

Even if we have, it is still super complex. This is where RL can help.

# Reinforcement Learning

Idea: MDP is awesome, but requires knowledge of underlying distributions. Besides, it is very difficult to solve.

$$\text{RL} = \text{Monte Carlo Method} + \text{MDP}$$

MC method: Use experience! In MC, we use repeated simulations to estimate.

- Model-based RL: Estimate underlying transition probabilities  
 $P(s', r | s, a) \quad \forall s \in S, s' \in S, a \in A(s), r \in R.$
- Model-free RL: Why bother to learn distributions if all we care is value? Directly learn value! (much easier!)

Now, since we don't know whether our current knowledge is accurate, we again face exploration vs exploitation!

The problem is a generalization of contexture bandit but now your current action affects next state.

There are many RL algorithm, we will talk about only one.

## Q-learning (Watkins (1989))

At state  $s_t$ , choose an action  $A_t \in A(s_t)$  according to Q function

Observe reward  $R_t$  and use this reward to update Q function

• Action selection:

At state  $s_t$ , select  $A_t = \arg\max_a Q(s_t, a)$  w.p.  $1 - \epsilon$

Randomly pick  $A_t \in A(s_t)$  w.p.  $\epsilon$  i.e.,  $\epsilon$ -greedy

• Update:

$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, A_t)]$$

$$\text{New } \leftarrow \text{old} + \text{step size} [\text{target} - \text{old}]$$

where  $\alpha$  is the step size introduced in Bandit problem and

$\gamma$  is the discounting factor introduced in MDP.

Remark: 1)  $\alpha \in [0, 1]$ : To deal with non-stationary problem, we give more recent observations higher weight.

$\gamma \in [0, 1]$ : To trade off between immediate and future rewards.

2) Deep reinforcement learning: Q-learning + NN-based nonlinear function approximation. A.K.A DQN.