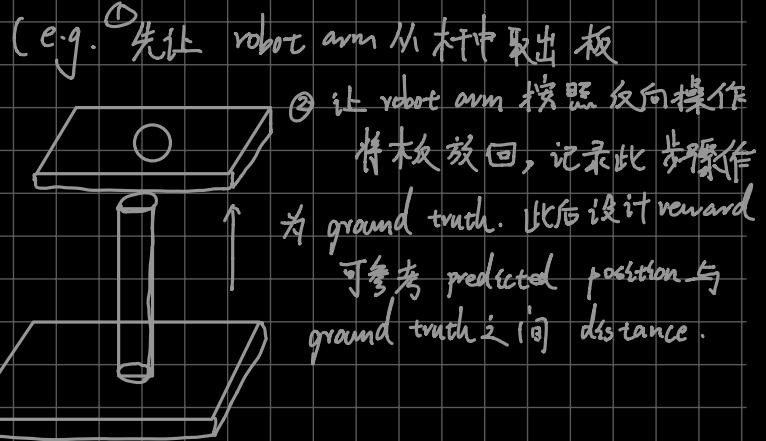


Lecture | RL Overview

Simple Task → Complex Task
↑
Reward design → Curriculum learning

(单式的 reward design 已经不能让 Agent 去完成任务。 e.g.: robot arms 将一堆乱的板子套入直杆中)



Lecture 2 Markov Decision Process

Concepts && Definition:

Framework:

Markov Process \rightarrow Markov Reward Process \rightarrow
Markov Decision Process \rightarrow Optimal Value Function

Markov Property:

- $P[S_{t+1} | S_t] = P[S_{t+1} | S_1, \dots, S_t]$
- The state captures all relevant information from the history.

Markov Process:

$\langle S, P \rangle$, S is a finite set of states (S_1, S_2, \dots), each state has Markov Property.

P is a state transition prob matrix:

$$P_{SS'} = P[S_{t+1} = S' | S_t = S]$$

Random Process:

$X(t, \xi) \Rightarrow$ 对于每个时刻 t , $X(t, \xi)$ 为随机变量

随机过程 $X(t, \xi)$ 为一组依赖于时间的随机变量 $X(t_1, \xi), X(t_2, \xi) \dots X(t_n, \xi)$ 组成。

Goal:

- Goal of RL: Find optimal policy to maximize the total amount of reward.
- Goal of MRP: Maximize the total amount of reward.

MRP:

- $\langle S, P, R, r \rangle \Rightarrow$ Markov Process with values.
- return G_t :

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

\Rightarrow total discount reward from t .

Why introduce discount r :

(i) Avoid infinite returns

(ii) Represent uncertainty about the future

\Rightarrow empirical design

- value function $V(s)$

$$V(s) = E [G_t | s_t = s] \Rightarrow$$

expected return starting from state s .

- Episodic task:

Special MRP that all sequences terminate.

MDP :

$$\langle \underbrace{S, A}, P, R, \gamma \rangle \Rightarrow \text{MRP with decisions}$$

It's an environment that all states are Markov
has Markov property

- policy $\pi(a|s)$:

$$\pi(a|s) = P[A_t = a | s_t = s] \Rightarrow$$

(i) Policy is a distribution of actions given states.

(ii) Policy fully defined the behaviour of an agent

(iii) Policy is time-independent.

- state-value function $V_{\pi}(s)$:

$$V_{\pi}(s) = E_{\pi} [G_t | s_t = s] \Rightarrow$$

value function following policy π .

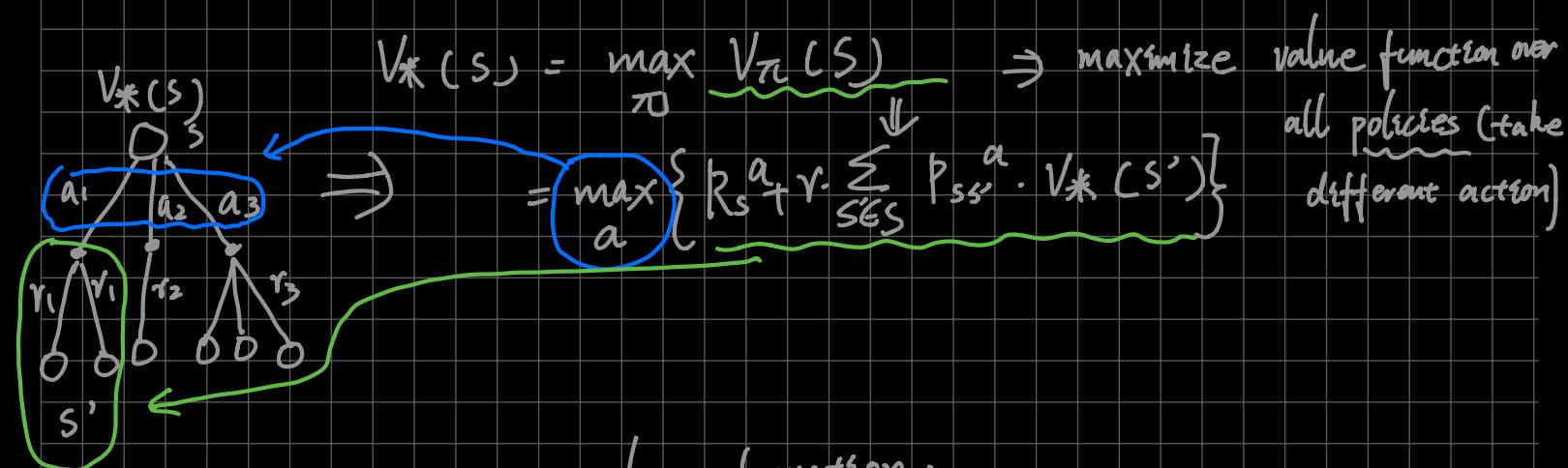
- action-value function $q_\pi(s, a)$:

$$q_\pi(s, a) = \mathbb{E}_\pi [G_t \mid s_t = s, a_t = a] \Rightarrow$$

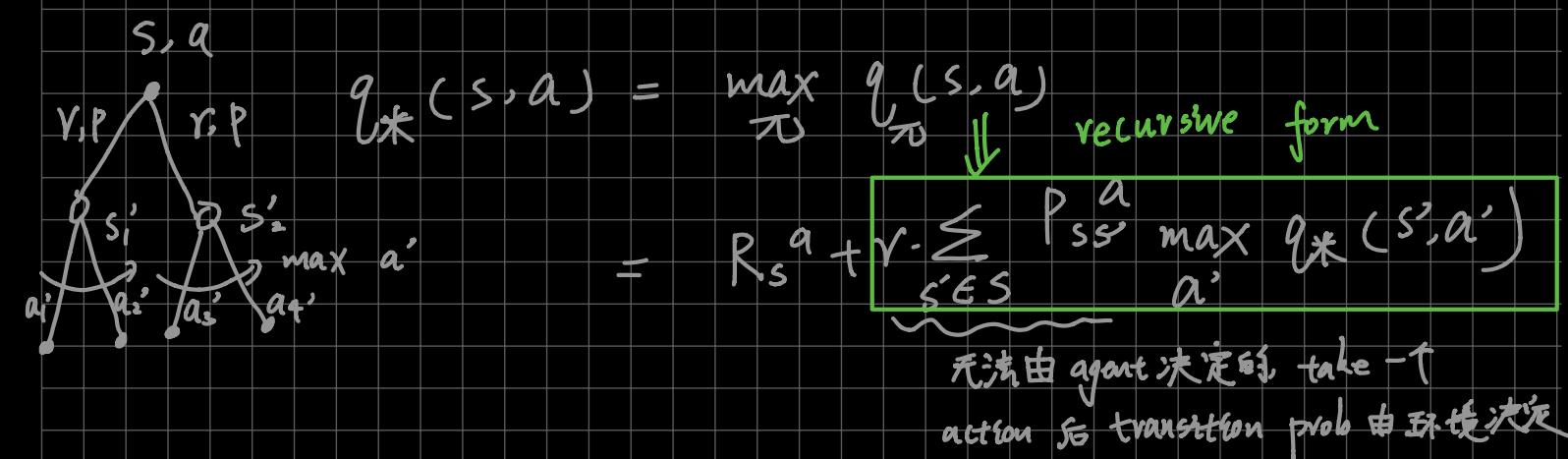
expected return starting from s, a and follows policy π .

Optimal Value Function: (specifies the best possible performance)

- state-value function:

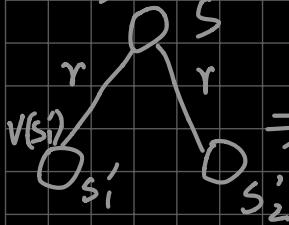


- action-value function:



Bellman Equation: (or Bellman Expectation Equation)

$$V(s)$$



• MRP:

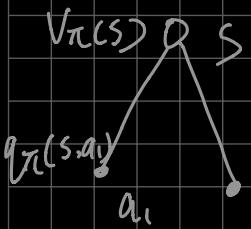
$$\Rightarrow V(s) = E[G_t | s_t = s] = E[\underbrace{R_{t+1} + r \cdot V(s_{t+1})}_{\text{recursive form}} | s_t = s]$$

• MDP:

$$V^\pi$$

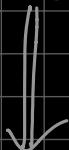
$$\textcircled{1} \quad V_\pi(s) = E[G_t | s_t = s] = E_\pi[R_{t+1} + r \cdot V_\pi(s_{t+1}) | s_t = s]$$

regressive definition



$$\Rightarrow \textcircled{2} \quad V_\pi(s) = \sum_{a \in A} \pi(a|s) \cdot q_\pi(s, a)$$

layer definition



$$V_\pi(s)$$

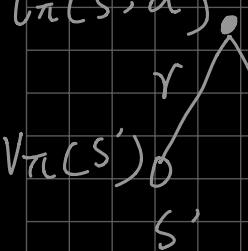
$$\Rightarrow \textcircled{3} \quad V_\pi(s) = \sum_{a \in A} \pi(a|s) \underbrace{(R_s^a + r \cdot \sum_{s' \in S} P_{ss'}^a \cdot V_\pi(s'))}_{2 \text{ layer definition}}$$

Q π :

$$\textcircled{1} \quad q_\pi(s, a) = E_\pi[G_t | s_t = s, a_t = a] = E_\pi[R_{t+1} + r \cdot q_\pi(s_{t+1}, a_{t+1}) | s_t = s, a_t = a]$$

regressive definition

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



$$\Rightarrow \textcircled{2} \quad q_\pi(s, a) = R_s^a + r \cdot \sum_{s' \in S} P_{ss'}^a \cdot V_\pi(s')$$

• Optimality Equation:

It is non-linear, so can't directly be solved. Need to use iterative solution methods. (Q-learning, SARSA)

Matrix Form: (Can be directly solved in this form)

- MRP:

$$V = R + \gamma \cdot P \cdot V \xrightarrow{\text{per state form}}$$

following policy π

$$\begin{bmatrix} V(1) \\ \vdots \\ V(n) \end{bmatrix} = \begin{bmatrix} R_1 \\ \vdots \\ R_n \end{bmatrix} + \gamma \cdot \underbrace{\begin{bmatrix} P_{11} & \cdots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \cdots & P_{nn} \end{bmatrix}}_{\text{state}} \begin{bmatrix} V(1) \\ \vdots \\ V(n) \end{bmatrix}$$

Solve: $V = (I - \gamma \cdot P)^{-1} \cdot R, O(n^3) \Rightarrow$

only for small MRPs, large MRPs costs too high. Use iterative methods (DP, MC-TD) instead.

- MDP:

$$V_\pi = R^\pi + \gamma \cdot P^\pi \cdot V_\pi$$

Solve: $V_\pi = (I - \gamma \cdot P^\pi)^{-1} \cdot R^\pi$