

# 二. MDPs

< Markov process / Markov chain >

时间, 状态都离散的过程 → 马尔可夫链.

马尔可夫过程(链)由元组  $(S, P)$  组成

- $S$  为有限状态的集合.
- $P$  为状态转移矩阵  $P_{ss'} = P[S_{t+1}=s' | S_t=s]$

(大写的  $S_t$  为随机状态, 小写的  $s$  为确定状态), 转移矩阵  $P$  ( $|S| \times |S|$  阶方阵)

< Episodes (幕) >

从初始状态  $S_1 = c_1$  开始, 从马尔可夫链中采样子序列, 并终止于终止状态, 那组子序列就叫做 episodes.

< 马尔可夫奖励过程 >

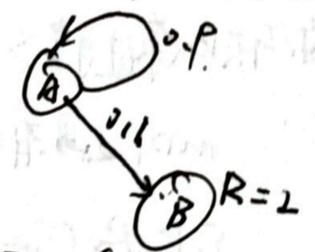
Markov Reward Process (MRP) 由元组  $(S, P, R, \gamma)$  组成.

$S, P$  在 Markov chain 中已有, 额外的

- $R_s$  为奖励函数,  $R_s = E[R_{t+1} | S_t=s]$  (除了奖励的随机性)
- $\gamma$  为折扣因子,  $\gamma \in [0, 1]$

< 奖励函数 >  $R_s = E[R_{t+1} | S_t=s]$  (一般认为到状态  $s$  之后, 系统才给出 Reward, 所以用  $(t+1)$  的下标, 如果在某  $s$ , 奖励没有随机性, 如

那么  $R_s = R_{t+1} = 2$ , 这里方便记为  $R=2$ .



< 回报, return >

从  $t$  时刻的  $S_t$  开始, 直至终止状态时, 所有奖励的折现之和  $C_t$  称为 return.

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

PPT 中称之为 accumulative discounted reward.  $\sum_{t=0}^{\infty} \gamma^t r_t$  (在  $t$  时有值之前  $s$  的 reward 也要计算).

Return 针对某一条幕说明/计算.

< 价值函数, Value Function > (单一条 episodes 的 return 高并不能说明这个状态好).

价值函数  $V(s)$  给出状态  $s$  的长期价值 (long-term value).

input: 状态 output: 价值.

在 MRP 中, 一个  $s$  的 期望回报 就是这个  $s$  的价值函数.

$$V(s) = E[C_t | S_t=s]$$

去除了 Reward 随机性  
去除了状态转移的随机性.

< 贝尔曼方程, Bellman Equation >



如何计算?

把子状态 (进行了) 也期望.

$$V(s) = E[C_t | S_t=s]$$

$$= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t=s]$$

$$= E[R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots) | S_t=s]$$

$$= E[R_{t+1} + \gamma V(S_{t+1}) | S_t=s]$$

$$\Rightarrow V(s) = E[R_{t+1} + \gamma V(S_{t+1}) | S_t=s]$$

前置去除  $\Rightarrow E[R_{t+1} + \gamma V(S_{t+1}) | S_t=s]$

递归.

$$v(s) = E[R_{t+1} + \gamma v(S_{t+1}) | S_t = s]$$

↑ 即时奖励
↑ 下一状态的价值

The other form of Bellman Equation:

$$v(s) = R_s + \gamma \sum_{s' \in S_{t+1}} P_{ss'} v(s')$$

↑ 当前状态奖励
↑ 折扣因子
↑ 转移函数
↑ 下一状态的价值函数

<贝尔曼方程的矩阵形式> ↓ 可将方程化为线性方程

$$\vec{v} = R + \gamma P \vec{v}$$

$$\begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix} = \begin{bmatrix} R_1 \\ \vdots \\ R_n \end{bmatrix} + \gamma \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix} \begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix}$$

$$\vec{v} = (I - \gamma P)^{-1} R \quad \text{计算复杂度为 } O(n^3)$$

直接求解适用于小型子马尔可夫奖励过程

<马尔可夫决策过程> (MDP)

MDP是具有决策性质的 MRP, 由元组  $(S, A, P, R, \gamma)$  组成

$S, P, R, \gamma$  在 MRP 中已有提及,  $A$  为有限动作的集合

而且  $P$  和  $R_s$  有所改变

$$P_{ss'}^a = P[S_{t+1} = s' | S_t = s, A_t = a]$$

$$R_s^a = E[R_{t+1} | S_t = s, A_t = a]$$

不只是与状态和动作有关, 还对 action 有关系

<策略, Policy>

Policy 完全定义了智能体的行为

马尔可夫决策过程中, Policy 仅取决于当前状态

input:  $s$       output:  $a$  动作

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

在给定  $s$  下,  $\pi$  与  $a$  成正比

$$\begin{bmatrix} \pi(a_1|s) \\ \vdots \\ \pi(a_n|s) \end{bmatrix}$$

⇒ 策略!

而  $\pi(a_n|s)$  与  $\pi(a_1|s)$  成正比

有了  $\pi(a|s)$ , 即 Policy,  $P$  和  $R$  发生改变, Markov chain, 而 MRP, MDP 都会发生变化: 即 policy 改变

$$P_{ss'}^{\pi} = \sum_{a \in A} \pi(a|s) P_{ss'}^a$$

$$R_s^{\pi} = \sum_{a \in A} \pi(a|s) R_s^a$$

<状态价值函数 (State-Value)> ← 对同一个  $s$  求解的  $v$ .

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

遵从策略  $\pi$ .

$V_{\pi}(s)$  为从  $s$  出发遵循  $\pi$  得到的期望回报.

如何求解?

由于  $\pi$  的加入,  $V_{\pi}(s)$  可以区分  $\pi$  的差别, 但是如果没有  $\pi$ , 那么就没有办法, 由此引入:

<动作价值函数, action-Value> ← 对同一个  $s$  求解的  $Q$ .

$$Q^{\pi}(s, a) = E_{\pi} [G_t | S_t = s, A_t = a]$$

遵从策略  $\pi$ .

$Q^{\pi}(s, a)$  为从  $s$  出发遵循  $\pi$  对当前状态  $s$  执行  $a$  得到的期望回报.

<贝尔曼期望方程, Bellman Expectation Equation> Bellman Equation + actions

$$V_{\pi}(s) = E_{\pi} [R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s]$$

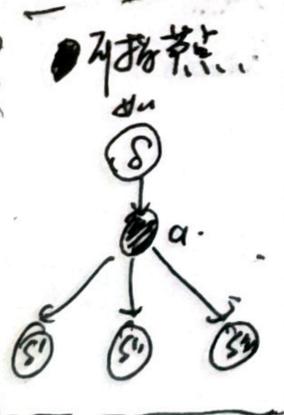
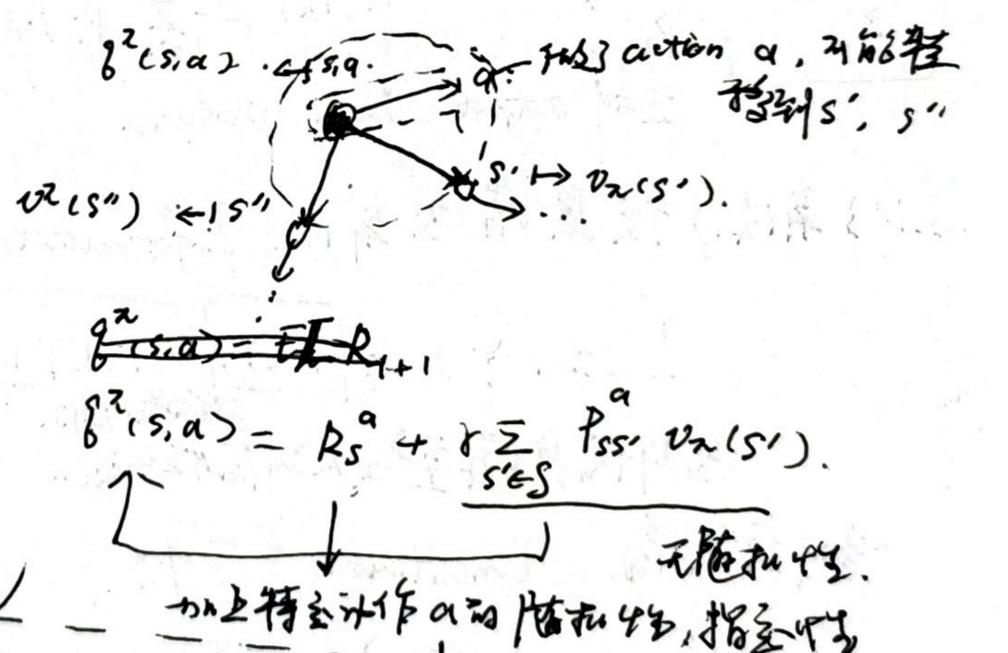
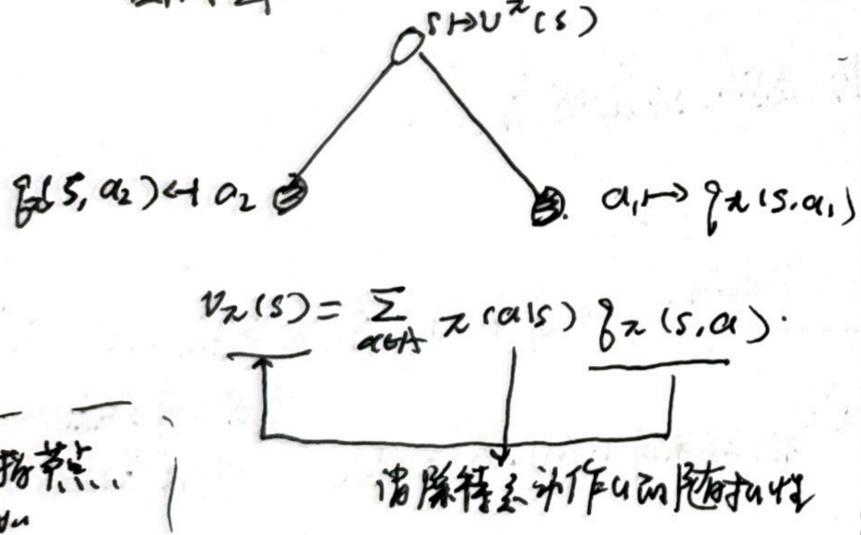
对动作清除随机. 对  $S_t \rightarrow S_{t+1}$  的  $A_t$  转移清除随机.

$$Q^{\pi}(s, a) = E_{\pi} [R_{t+1} + \gamma Q^{\pi}(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$$

$a$  不同, 后续  $A$  也不同.  $a$  不同, 后续  $A$  也不同.  $\downarrow$  随机.

回顾图:

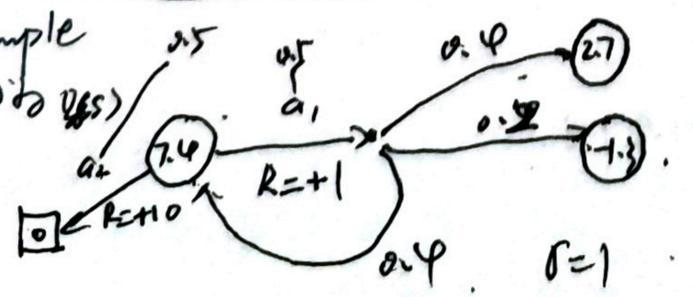
$V_{\pi}(s)$  比  $Q^{\pi}(s, a)$  少了  $A_t$  的随机性.  $\leftarrow$  缺少了  $a$  动作的随机性.



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

$$Q^{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \sum_{a' \in A} \pi(a'|s') Q^{\pi}(s', a')$$

example



$$V_{\pi}(s) = 0.5 \times (10 + 0) + 0.5 \times (1 + 0.5 \times 2.7 + 0.2 \times (1.3) + 0.4 \times 7.4)$$

2.47476

7.4

< 最优价值函数, Optimal Value Function >

最优目标

最优状态价值函数

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

最优动作价值函数

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

最优价值函数已知, MDP 动态求解

都是对不同策略取最优

Bellman equation

~~$$V_{\pi}^*(s)$$~~

$$Q_{\pi}^*(s, a) = E [R_{t+1} + \gamma \max_{a'} Q^*(s', a') | s, a]$$

< 最优策略, Optimal Policy >

∴ < 偏序 > 对于任意两个状态  $s$  都有  $V_{\pi}(s) \geq V_{\pi'}(s), \pi \geq \pi'$

在 MDP 中, 至少存在一个策略不低于其他策略, 即  $\pi_{\pi} \geq \pi, \forall \pi, \pi_{\pi} \rightarrow \text{optimal policy}$

$$\pi_{\pi}(a|s) = \begin{cases} 1 & \text{当 } a = \arg \max_{a \in A} Q_{\pi}(s, a) \\ 0 & \end{cases}$$

arg  $\Rightarrow$  argument, ~~自变量~~

$\Rightarrow$  只留最优的动作, 舍弃其他

< 贝尔曼最优方程 > (Bellman Optimality Equation) 非线性, 不可通过矩阵运算

$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) Q_{\pi}(s, a) = \max_a Q_{\pi}(s, a) \rightarrow V_{\pi}^*(s)$$

$$Q_{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi}^*(s')$$

other form

$$V_{\pi}^*(s) = \max_a (R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi}^*(s'))$$

$$Q_{\pi}^*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \max_{a'} Q_{\pi}^*(s', a')$$



三. From Q-learning to Actor Critic.

Q(s, A) 的求解:  $\rightarrow$  Sarsa 算法更新 Q(s, A), 策略评估, 用贪婪优化.



$$Q(s, A) = E_{\pi} [G_t | S_t = s, A_t = A]$$

$$= E_{\pi} [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) | S_t = s, A_t = A]$$

$$Q(s, A) \leftarrow Q(s, A) + \alpha (R_t + \gamma Q(S_{t+1}, A_{t+1}) - Q(s, A))$$

$G_t$  用  $\epsilon$ -greedy 策略

在策略中的 Sarsa 依赖于 Bellman 期望方程.

- 1: 随机初始化  $Q(s, a), \forall s \in S, a \in A(s)$ . 令  $Q(\text{terminal-state}, \cdot) = 0$ .
- 2: REPEAT (对于每个 episode):
- 3:     初始化  $S$ .
- 4:     使用从  $Q$  得到的策略 (eg.  $\epsilon$ -greedy) 在状态  $S$  时选择动作  $A$ .
- 5:     REPEAT (对于 episode 中的每一步):
- 6:         执行动作  $A$ , 观察得到  $R, S'$
- 7:         根据从  $Q$  中得到的策略 (eg.  $\epsilon$ -greedy) 在状态  $S'$  时选择动作  $A'$ .
- 8:          $Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma Q(s', A') - Q(s, A)]$
- 9:          $S \leftarrow S'; A \leftarrow A'$
- 10:     直到  $S$  为终止状态.

< On-policy, off-policy, 在线, 离线 >

< Off-policy >

目标策略: 用来学习的策略. 行为策略: 生成训练样本的策略.

评估目标策略  $v(a|a)$  以估计  $V_{\pi}(s)$  或  $Q_{\pi}(s, a)$ .

同时运行两个策略  $\mu(a|s)$   $\{S_t, A_t, R_t, \dots, R_T\} \sim \mu$

< Q-Learning > = off policy. 依赖于 Bellman 最优方程  $G_t$ .  $Q$  更新的核心. 偷别人的东西.

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$

行为策略和评估策略都随  $Q$  的更新而更新.

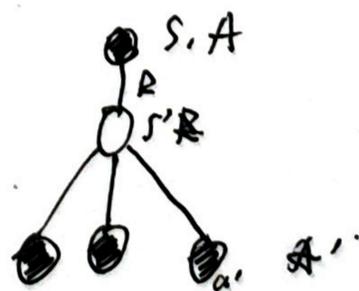
目标策略是贪婪的, 行为策略是  $\epsilon$ -greedy 的.

目标策略简化:

$$R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$$

$$= R_{t+1} + \gamma Q(S_{t+1}, \arg \max_{a'} Q(S_{t+1}, a'))$$

$$= R_{t+1} + \max_{a'} \gamma Q(S_{t+1}, a')$$



Q-Learning

$$Q(s, A) \leftarrow Q(s, A) + \alpha (R + \gamma \max_{a'} Q(s', a') - Q(s, A))$$

# Sarsamax (Q-learning)

- 1: 随机初始化  $Q(s, a)$ ,  $\forall s \in S, a \in A(s)$ , 将  $Q(\text{terminal-state}, \cdot) = 0$ .
- 2: REPEAT (对于每一个 episode):
- 3:     初始化  $s$ .
- 4:     REPEAT (对于 episode 中每一步):
- 5:         根据从  $Q$  学到的策略 (eg.  $\epsilon$ -greedy) 在  $S$  中选择  $A$ .
- 6:         执行  $A$ , 观察到  $R, s'$ .
- 7:         更新  $Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma \max_{a'} Q(s', a') - Q(s, A)]$
- 8:          $s \leftarrow s'$
- 9:     直到  $s$  为终止状态.

我更新了  
之后跟操  
作者共同  
优化  $Q$ .

不需要自己选择动作.

未建模, 因为操作者是人.

Sarsa 走正路, 采取点收益最高的路径.

Q-Learning 走正路, 点收益不如 Sarsa, 所谓“高悬空中”!

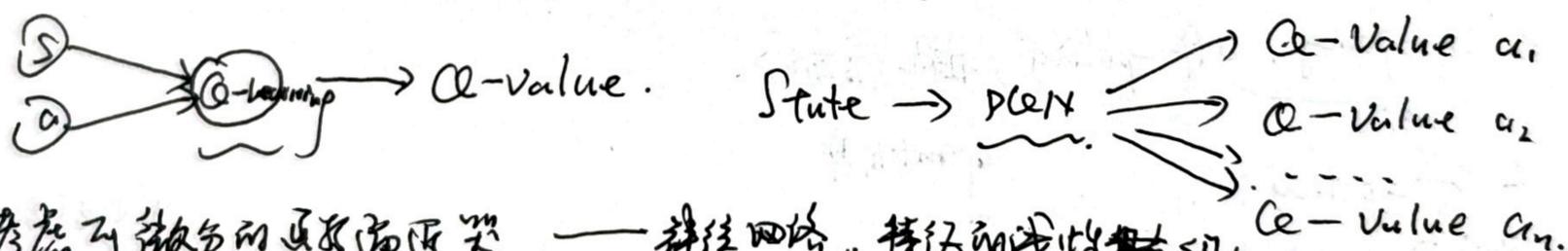
## < DQN >

如果  $S$  有限, 上述方法可用线性或其他优化方法, 但是有无限状态呢?

以前的过程会建立一个  $Q$ -table, 学习的过程就是填表的过程, 但是太多的数据, 无法存储  
解决方法: (用函数逼近价值函数,  $\text{input} \rightarrow s, \text{output } Q$ )

$$\hat{V}(s; w) \approx V_w(s) \quad \text{or} \quad \hat{Q}(s, a; w) \approx Q_w(s, a)$$

input  
↓  
system  
↓  
output



考虑用微分的函数逼近器。——神经网络, 神经网络线性组合。  
此外, 用于非平稳非独立分布的数据训练方法。QL 中的数据特征

## < 线性最小二乘 >

线性最小二乘 = 矩阵求逆

非线性最小二乘 = < Gradient Descent >

While True.  $\rightarrow$  样本量大, 使用SGD优化  $L(w)$

$$\nabla_w L(w) = \nabla_w \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, w), y_i)$$

$$w \leftarrow w - \alpha \nabla_w L(w) \quad \Delta w = \alpha (y_i - f(x_i, w)) \frac{\partial f}{\partial w}$$

GD, SGD, MBGD, GD with Momentum =  $-\frac{1}{2} \alpha \nabla_w L(w)$

$$\Rightarrow w \leftarrow w + \Delta w$$

$$\Delta w = -\frac{1}{2} \alpha \nabla_w J(w) = \alpha E_w [(V_w(s) - \hat{V}(s, w)) \nabla_w \hat{V}(s, w)]$$

$$\pm \text{SGD 求导 } \Delta w = \alpha [(V_w(s) - \hat{V}(s, w)) \nabla_w \hat{V}(s, w)]$$

状态  $s$  可以用特征向量表示.