

纲要

- 1 强化学习简介
- 2 马尔科夫决策
- 3 动态规划方法
- 4 蒙特卡罗方法
- 5 时序差分学习
- 6 策略梯度方法
- 7 应用实践

马尔科夫决策过程

$\mathbb{P}(s_{t+1}|s_t)$ 表示 t 时刻的状态 s_t 转移到 $t+1$ 时刻的状态 s_{t+1} 的条件概率， $\mathbb{P}(s_{t+1}|s_t, \dots, s_0)$ 表示在历史状态 s_t, \dots, s_0 条件下 $t+1$ 时刻转移到状态 s_{t+1} 的条件概率。两者相等说明转移概率不受历史状态信息影响，即与历史状态 s_{t-1}, \dots, s_0 无关，只与当前状态 s_t 有关。

随机过程满足马尔科夫性，离散状态之间转移概率矩阵 P 为：

$$\begin{matrix} & s_1 & s_2 & s_3 & \cdots & s_{n-1} & s_n \\ \begin{matrix} s_1 \\ s_2 \\ s_3 \\ \vdots \end{matrix} & \begin{pmatrix} p_{11} & p_{12} & p_{13} & \cdots & p_{1,n-1} & p_{1n} \\ p_{21} & p_{22} & p_{23} & \cdots & p_{2,n-1} & p_{2n} \\ p_{31} & p_{32} & p_{33} & \cdots & p_{3,n-1} & p_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \end{pmatrix} \end{matrix}$$

策略函数

一般来说，复杂环境模型包括状态转移函数和回报函数。在问题求解之前，我们需要知道状态转移函数和回报函数，然后通过算法求解最优策略函数 π 。策略函数 π 可分成两种类型，一种是**随机性策略**，表示为：

$$\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1] \quad (5)$$

随机性策略函数 π 输出状态 s 下选择动作 a 的概率。智能体基于动作概率分布进行随机采样，得到最终动作 a 。另一种是**确定性策略**，表示为：

$$\pi : \mathcal{S} \rightarrow \mathcal{A}. \quad (6)$$

确定性策略函数 π 直接输出状态 s 下的动作 a 。从另一个角度而言，策略函数 π 也可以分成连续型策略和离散型策略。

状态转移函数

强化学习的目标是学习智能体的策略 π ， π 可以建模成一个函数，将随机过程的状态空间映射到动作空间，表示为 $\pi: \mathcal{S} \rightarrow \mathcal{A}$ 。动作影响状态转移和奖励回报。

在复杂环境模型已确定的情况下，智能体的策略输出的动作直接影响奖励回报，以随机策略函数举例分析，将状态转移函数重写为：

$$P_{ss'}^\pi = \sum_{a \in A} \pi(a|s) P_{ss'}^a, \quad (7)$$

其中， $\pi(a|s)$ 表示状态 s 下执行动作行为 a 的概率。

奖励函数

奖励函数或回报函数 R 决定了智能体在环境状态 s 下执行动作 a 后得到的奖励值 R_s^a ，可以表示为：

$$R_s^a = \mathbb{E}[R_t | S_t = s, A_t = a]. \quad (8)$$

我们结合策略函数 π ，可以计算智能体在当前时刻状态 s 下选择不同动作后获得的期望奖励回报：

$$R_s^\pi = \sum_{a \in A} \pi(a|s) R_s^a. \quad (9)$$

累积回报

智能体从当前时刻状态 s 开始直至终止状态所获得的**累积奖励回报**可定义为：

$$G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \cdots + \gamma^T R_{t+T} = \sum_{k=0}^T \gamma^k R_{t+k}, \quad (10)$$

其中, γ 是折扣系数, 且 $\gamma \in [0, 1)$ 。强化学习的折扣系数与金融分析中的折扣因子类似。因为 $\gamma < 1$, G_t 不会出现无穷大。智能体在无限长时间的累积收益情况的具体分析如下:

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k} < R_{\max} \sum_{k=0}^{\infty} \gamma^k = R_{\max} \frac{1}{1-\gamma}, \quad (11)$$

其中, R_{\max} 表示最大即时奖励值。

状态值函数

我们在累积回报 G_t 基础上，可以定义 **状态值函数** $V_\pi(s)$ 。 $V_\pi(s)$ 表示从状态 s 出发，智能体基于当前策略函数 π 获得的期望回报，具体数学表示为：

$$V_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]. \quad (12)$$

状态值函数 $V_\pi(s)$ 是智能体在状态 s 获得累积回报 G_t 的期望。对状态值函数 $V_\pi(s)$ 进行简单推导，可以得到

$$\begin{aligned}
 V_\pi(s) &= \mathbb{E}_\pi[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \gamma^3 R_{t+3} + \cdots | S_t = s] \\
 &= \mathbb{E}_\pi[R_t + \gamma(R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots) | S_t = s] \\
 &= \mathbb{E}_\pi[R_t + \gamma G_{t+1} | S_t = s] \\
 &= \mathbb{E}_\pi[R_t + \gamma V_\pi(S_{t+1}) | S_t = s].
 \end{aligned} \quad (13)$$

S_{t+1} 为随机变量， G_{t+1} 的期望值用状态值函数 $V_\pi(S_{t+1})$ 替换。

状态-动作值函数

同样地定义**状态-动作值函数** $Q_{\pi}(s, a)$ 。 $Q_{\pi}(s, a)$ 表示智能体从状态 s 出发，基于当前策略函数 π 执行动作 a 能够获得的期望累积回报，衡量了状态 s 下动作 a 的价值，具体数学表达式为：

$$Q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a]. \quad (14)$$

对状态-动作值函数 $Q_{\pi}(s, a)$ 进行简单推导，可以得到：

$$\begin{aligned}
 Q_{\pi}(s) &= \mathbb{E}_{\pi}[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \gamma^3 R_{t+3} + \cdots | S_t = s, A_t = a] \\
 &= \mathbb{E}_{\pi}[R_t + \gamma(R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots) | S_t = s, A_t = a] \\
 &= \mathbb{E}_{\pi}[R_t + \gamma G_{t+1} | S_t = s, A_t = a] \\
 &= \mathbb{E}_{\pi}[R_t + \gamma Q_{\pi}(S_{t+1}) | S_t = s, A_t = a].
 \end{aligned} \quad (15)$$

S_{t+1} 为随机变量， G_{t+1} 的期望值用动作-状态值函数 $Q_{\pi}(S_{t+1})$ 替换。

状态-动作值函数与状态值函数的关系

基于状态值函数和状态-动作值函数的定义，可以发现两者之间具有紧密的**联系**：

$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) Q_{\pi}(s, a). \quad (16)$$

上式说明状态 s 的值函数 $V_{\pi}(s)$ 是在策略函数 π 下执行动作 a 获得累积收益回报的期望值。 $\pi(a|s)$ 表示状态 s 下动作 a 的概率，状态-动作值函数 $Q_{\pi}(s, a)$ 表示状态 s 下动作 a 的期望累积收益，因此，智能体在状态 s 下遍历所有动作并累积期望收益 $\pi(a|s)Q_{\pi}(s, a)$ ，得到了状态 s 的价值 $V_{\pi}(s)$ 。

状态-动作值函数与状态值函数的关系

我们也可以将两者之间的紧密联系表示为：

$$Q_{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi}(s'). \quad (17)$$

式 (17) 说明，状态-动作值函数等于动作 a 的即时奖励值加上下一个可能状态 s' 的值函数 $V_{\pi}(s')$ 的加权和 $\sum_{s' \in S} P_{ss'}^a V_{\pi}(s')$ 。由于 $V_{\pi}(s')$ 是下一个时刻状态值，我们需要乘上折扣因子 γ 。将式 (17) 代入式 (16)，可以得到状态值函数另外一种更加复杂的表示形式：

$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) \left(R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi}(s') \right). \quad (18)$$

上式只包含了状态值函数 V_{π} ，无状态-动作值函数 Q_{π} ，此方程是求解状态值函数 $V_{\pi}(s)$ 的关键。

状态-动作值函数与状态值函数的关系

类似地，将式（16）代入式（17），可以得到状态-动作值函数另外一种更加复杂的表示形式：

$$Q_{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \left(\sum_{a' \in A} \pi(a'|s') Q_{\pi}(s', a') \right). \quad (19)$$

上式只包含了状态-动作值函数 Q_{π} ，无状态值函数 V_{π} ，此方程是求解动作-状态值函数 $Q_{\pi}(s, a)$ 的关键。

Bellman方程

在马尔科夫回报过程中，关于状态值函数 $V_{\pi}(s)$ 的Bellman方程可以表示为：

$$V_\pi(s) = \mathbb{E}_\pi[R_t + \gamma V_\pi(S_{t+1}) | S_t = s]. \quad (20)$$

状态值函数 $V_{\pi}(s)$ 的Bellman方程不包含策略函数和动作。

Bellman方程

在马尔科夫决策过程中，关于状态值函数 $V_\pi(s)$ 的Bellman方程可以表示为：

$$\begin{aligned}
 V_\pi(s) &= \sum_{a \in A} \pi(a|s) \left(R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_\pi(s') \right) \\
 &= \sum_{a \in A} \pi(a|s) R_s^a + \gamma \sum_{a \in A} \pi(a|s) \left(\sum_{s' \in S} P_{ss'}^a V_\pi(s') \right). \quad (21)
 \end{aligned}$$

Bellman方程

将式 (9)，即 $R_s^\pi = \sum_{a \in A} \pi(a|s) R_s^a$ ，代入上式，可以得到：

$$\begin{aligned}
 V_\pi(s) &= R_s^\pi + \gamma \sum_{a \in A} \pi(a|s) \left(\sum_{s' \in S} P_{ss'}^a V_\pi(s') \right) \\
 &= R_s^\pi + \gamma \sum_{s' \in S} \left(\sum_{a \in A} \pi(a|s) P_{ss'}^a \right) V_\pi(s'). \tag{22}
 \end{aligned}$$

将式 (7)，即 $P_{ss'}^\pi = \sum_{a \in A} \pi(a|s) P_{ss'}^a$ ，代入上式，可以得到：

$$V_\pi(s) = R_s^\pi + \gamma \sum_{s' \in S} P_{ss'}^\pi V_\pi(s'). \tag{23}$$

Bellman方程

针对马尔科夫决策过程状态空间中的每一个状态 s_1, s_2, \dots ，都可以写出类似的Bellman方程：

$$\begin{aligned}
 V_{\pi}(s_1) &= R_{s_1}^{\pi} + \gamma \sum_{s' \in S} P_{s_1 s'}^{\pi} V_{\pi}(s') \\
 V_{\pi}(s_2) &= R_{s_2}^{\pi} + \gamma \sum_{s' \in S} P_{s_2 s'}^{\pi} V_{\pi}(s') \\
 &\dots \\
 V_{\pi}(s_n) &= R_{s_n}^{\pi} + \gamma \sum_{s' \in S} P_{s_n s'}^{\pi} V_{\pi}(s').
 \end{aligned} \tag{24}$$

Bellman方程

我们可以将上述方程组改写成矩阵形式：

$$\begin{bmatrix} V_{\pi}(s_1) \\ V_{\pi}(s_2) \\ V_{\pi}(s_3) \\ \vdots \\ V_{\pi}(s_n) \end{bmatrix} = \begin{bmatrix} R_{\pi}(s_1) \\ R_{\pi}(s_2) \\ R_{\pi}(s_3) \\ \vdots \\ R_{\pi}(s_n) \end{bmatrix} + \gamma \begin{bmatrix} P_{11}^{\pi} & P_{12}^{\pi} & \cdots & P_{1n}^{\pi} \\ P_{21}^{\pi} & P_{22}^{\pi} & \cdots & P_{2n}^{\pi} \\ P_{31}^{\pi} & P_{32}^{\pi} & \cdots & P_{3n}^{\pi} \\ \vdots & \vdots & \cdots & \vdots \\ P_{n1}^{\pi} & P_{n2}^{\pi} & \cdots & P_{nn}^{\pi} \end{bmatrix} \begin{bmatrix} V_{\pi}(s_1) \\ V_{\pi}(s_2) \\ V_{\pi}(s_3) \\ \vdots \\ V_{\pi}(s_n) \end{bmatrix}. \quad (25)$$

因此，我们可以进一步用矩阵符号表示：

$$V_{\pi} = R_{\pi} + \gamma P_{\pi} V_{\pi}, \quad (26)$$

其中， V_{π} 和 R_{π} 为列向量， P_{π} 为状态转移概率矩阵。

Bellman方程

求解状态值函数的Bellman方程组 (26)，可得：

$$\begin{aligned}
 V_{\pi} &= R_{\pi} + \gamma P_{\pi} V_{\pi} \\
 (I - \gamma P_{\pi}) V_{\pi} &= R_{\pi} \\
 V_{\pi} &= (I - \gamma P_{\pi})^{-1} R_{\pi}.
 \end{aligned} \tag{27}$$

因此，状态值函数 V_{π} 可以基于状态转移函数 $P: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ 和回报奖励函数 $R: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathcal{R}$ 直接求解。

策略迭代算法

在马尔科夫决策过程的状态值函数 V_π 解析公式

$$V_\pi = (I - \gamma P_\pi)^{-1} R_\pi$$

中：

- 解析解存在需要矩阵 $(I - \gamma P_\pi)$ 可逆。
- 对于现实问题，复杂随机过程的状态转移概率矩阵不一定满足 $(I - \gamma P_\pi)$ 可逆。
- 即使逆矩阵存在，由于马尔科夫决策过程的状态数量多，状态转移矩阵规模大。
- 矩阵求逆计算复杂度较高。
- 超大规模矩阵逆计算在有限时间和有限内存资源条件下也基本不可能完成。

策略迭代算法

因此，我们可以考虑采用数值方法中的迭代方法求解此类问题，迭代公式表示如下：

$$V_{k+1} = R_{\pi} + \gamma P_{\pi} V_k. \quad (28)$$

迭代公式为了求出状态值函数，用等式右边的状态值函数 V_k 计算出等式左边的状态值函数 V_{k+1} 后，继续将 V_{k+1} 代入等式右边，迭代计算状态值函数 V_{k+2} ，以此类推，循环迭代。
迭代公式的矩阵形式展开后可以得到每一个状态值函数迭代公式为

$$V_{k+1}(s) = \sum_{a \in A} \pi(a|s) \left(R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_k(s') \right). \quad (29)$$

状态值函数迭代过程为策略评估，即给定策略函数 π ，可以估计各个状态值函数。

策略迭代算法

通常地，状态值函数的初始值都设置成0，即 $V_0 = 0$ 。我们通过迭代公式计算状态值函数 V_k ，直至收敛到 V^* ，则 V^* 为迭代公式的不动点，且不动点 V^* 必定满足：

$$V^* = R_{\pi} + \gamma P_{\pi} V^*. \quad (30)$$

将式(28)减去式(30)，可得：

$$V_{k+1} - V^* = \gamma P_{\pi} (V_k - V^*). \quad (31)$$

用 e_k 表示第 k 步数值误差，定义为数值解 V_k 与精确解 V^* 之差：

$$e_k = V_k - V^*, \quad (32)$$

策略迭代算法

同样， e_{k+1} 表示第 $k+1$ 步数值误差，定义为第 $k+1$ 步数值解 V_{k+1} 与精确解 V^* 之差：

$$e_{k+1} = V_{k+1} - V^*, \quad (33)$$

我们将上述两个误差定义公式代入式(31)，可以得到迭代公式：

$$e_{k+1} = \gamma P_{\pi} e_k. \quad (34)$$

进一步迭代计算可得：

$$e_k = \gamma^k (P_{\pi})^k e_0. \quad (35)$$

基于状态转移概率矩阵 P_{π} 的性质和折扣系数 $\gamma < 1$ ，可以得到：

$$\lim_{k \rightarrow \infty} e_k = 0. \quad (36)$$

即当迭代次数足够多时，误差趋于0，此时数值解 V_k 收敛到精确解 V^* 。

策略改进

在定义模型值函数后，智能体通过值函数得到最优策略。策略改进的过程中可以采用贪心策略：

$$\pi_{k+1}(s) = \arg \max_a Q_{\pi_k}(s, a), \quad (37)$$

其中， $Q_{\pi_k}(s, a)$ 表示智能体在状态 s 下动作 a 的价值，状态-动作值函数的下标 π_k 表示当前策略函数：

$$Q_{\pi_k}(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi_k}(s'). \quad (38)$$

将其代入式（37）后，可得：

$$\pi_{k+1}(s) = \arg \max_a \left(R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi_k}(s') \right). \quad (39)$$

此过程叫做策略改进，基于给定的值函数 $V_{\pi_k}(s)$ 改进策略函数。

策略迭代算法的具体迭代过程伪代码

Algorithm 8: 策略迭代算法伪代码

```

Input: 马尔科夫决策过程五元组  $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ 
Output: 最优策略  $\pi^*$ 
1 初始化状态值函数  $V(s) = 0$ , 初始化策略函数  $\pi$  为随机策略
2 for  $k = 0, 1, 2, 3, \dots$  do
3     % 策略评估
4     for  $l = 0, 1, 2, 3, \dots$  do
5         for  $s \in \mathcal{S}$  do
6              $V'(s) = \sum_{a \in \mathcal{A}} \pi(a|s) (R_s^a + \gamma \sum_{s' \in \mathcal{S}} P_{sa'}^a V_k(s'))$ 
7             if  $V' == V$  then
8                 | 停止迭代;
9             else
10                |  $V = V'$ 
11    % 策略改进
12    for  $s \in \mathcal{S}$  do
13        |  $\pi'(s) = \arg \max_{a \in \mathcal{A}} (R_s^a + \gamma \sum_{s' \in \mathcal{S}} P_{sa'}^a V_k(s'))$ 
14    % 策略迭代终止判断
15    if  $\pi' == \pi$  then
16        | 停止迭代
17    else
18        |  $\pi = \pi'$ 
19  $\pi^* = \pi'$ 
    
```

值函数迭代算法

值函数迭代算法直接求解马尔科夫决策过程中的值函数，智能体基于值函数选择最优动作。迭代算法直接进行值函数迭代，值函数迭代算法无策略改进过程，值函数迭代更新公式如下：

$$V_{k+1}(s) = \max_{a \in A} \left(R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_k(s') \right), \quad (40)$$

公式中的 $\max_{a \in A}$ 操作是值函数迭代算法的关键。直接使用当前值函数的最大值来更新值函数，此过程类似如下操作：

$$V_{k+1}(s) = \max_{a \in A} Q_k(s, a), \quad (41)$$

其中，

$$Q_k(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_k(s'). \quad (42)$$

值函数迭代算法

换言之，值函数迭代算法将策略改进过程融入了值函数迭代过程，其关键操作为值函数更新中的 $\max_{a \in A}$ 操作。迭代过程收敛后得到最终值函数，进而求得最优策略函数：

$$\pi(s) = \arg \max_{a \in A} \left(R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_k(s') \right). \quad (43)$$

因此，值函数迭代算法中的关键步骤为值函数更新中的 $\max_{a \in A}$ 操作， $\max_{a \in A}$ 操作需要遍历动作空间中的所有动作 a ，并选择期望累积收益最大的动作。值函数更新是基于期望收益最大的动作更新，具体迭代过程如值函数迭代算法伪代码：

纲要

- 1 强化学习简介
- 2 马尔科夫决策
- 3 动态规划方法
- 4 蒙特卡罗方法**
- 5 时序差分学习
- 6 策略梯度方法
- 7 应用实践

蒙特卡罗估计

如果马尔科夫决策过程**不包含模型状态转移函数或不**存在显式的**状态转移函数等**，那么动态规划方法（策略迭代和值函数迭代算法）不具有可行性，只能基于状态值函数或者状态-动作值函数的初始定义分析。假设 T 步动作的状态值函数定义如下：

$$V_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s] = \mathbb{E}_{\pi}[R_t + \gamma R_{t+1} + \cdots + \gamma^{T-1} R_{t+T-1} | S_t = s]. \quad (44)$$

一般而言，无穷步动作的状态值函数可定义为：

$$\begin{aligned}
 V_{\pi}(s) &= \mathbb{E}_{\pi}[G_t | S_t = s] \\
 &= \mathbb{E}_{\pi}[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \gamma^3 R_{t+3} + \cdots | S_t = s] \\
 &= \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k} | S_t = s\right].
 \end{aligned} \quad (45)$$

蒙特卡罗估计

同时，状态-动作值函数定义为：

$$\begin{aligned}
 Q_{\pi}(s, a) &= \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a] \\
 &= \mathbb{E}_{\pi}[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \gamma^3 R_{t+3} + \cdots | S_t = s, A_t = a] \\
 &= \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k} | S_t = s, A_t = a\right].
 \end{aligned}
 \tag{46}$$

上述公式基于数学期望定义了状态值函数和状态-动作值函数。在统计分析过程中，期望可以通过随机采样样本的均值进行估计和近似。在实际计算中，我们需要基于策略 π 采样得到一些样本，代入公式中估计出均值大小，近似（估计）值函数的期望值，此为蒙特卡罗方法的重要思想。

蒙特卡罗估计

我们从状态 s_0 开始随机采样了 n 条完整的轨迹数据，运用蒙特卡罗方法估计状态 s_0 的状态值函数。经验轨迹 i 中状态 s_0 对应的累积奖励值为：

$$G_i = \sum_{t=0}^T r_{it}, \quad (48)$$

其中 $i \in \{1, 2, \dots, n\}$ 。我们计算 n 条完整轨迹的平均累积回报值作为状态 s_0 的状态值函数的估计：

$$V_{\pi}(s_0) = \frac{1}{n} \sum_{i=1}^n G_i. \quad (49)$$

上述 $V_{\pi}(s)$ 的估计过程简化了对轨迹样本的统计分析。在实际计算过程中，可以进一步考虑初访和每访的区别。

蒙特卡罗估计

因此，我们将公式中 $\frac{1}{n}$ 替换成 α ，状态 s 在第 n 次采样后，更新状态值函数：

$$V_n(s) = V_{n-1}(s) + \alpha(G_n - V_{n-1}(s)), \quad (51)$$

其中 α 为机器学习中的学习率。

蒙特卡罗估计

运用蒙特卡罗方法估计状态 s_0 时选择动作 a_0 的状态-动作值函数，计算随机采样的 n 条完整^{完整}的轨迹数据的平均累积回报：

$$Q_{\pi}(s_0, a_0) = \frac{1}{n} \sum_{i=1}^n G_i. \quad (53)$$

同样，蒙特卡罗算法采用增量值更新状态-动作值函数：

$$Q_n(s, a) = Q_{n-1}(s, a) + \alpha(G_n - Q_{n-1}(s, a)), \quad (54)$$

其中， α 为学习率。获得状态-动作值函数 $Q(s, a)$ 后，可以计算最优策略函数：

$$\pi(s) = \arg \max_a Q(s, a). \quad (55)$$

蒙特卡罗估计

蒙特卡罗强化学习算法也是一个迭代算法，每次采样都使用最新的状态-动作值函数 $Q(s, a)$ 来构建策略函数。在实际采样过程中，为了增加随机采样完整轨迹的多样性，一般并不完全按照最新的状态-动作值函数 $Q(s, a)$ 进行采样，而是采用 ϵ -贪心算法，具体采样策略为：

$$\pi(s, a) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|\mathcal{A}|}, & a = \arg \max_a Q(s, a) \\ \frac{\epsilon}{|\mathcal{A}|}, & a \neq \arg \max_a Q(s, a) \end{cases}, \quad (56)$$

其中, $|A|$ 表示动作空间 A 的大小, 即动作数量; ϵ 决定了智能体的探索能力, ϵ 越大, 智能体的行为随机性越大, 探索能力越大。

$$\pi^*(s) = \arg \max_a Q(s, a)$$

时序差分算法

我们简单比较**蒙特卡罗强化学习算法**和动态规划（Dynamic Programming, DP）算法的差异和实用场景。

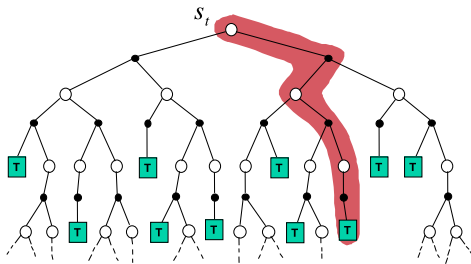


图 1: 图片来自David Silver 个人网站¹

¹<https://www.davidsilver.uk/teaching/>

时序差分算法

我们简单比较蒙特卡罗强化学习算法和动态规划（Dynamic Programming, DP）算法的差异和实用场景。

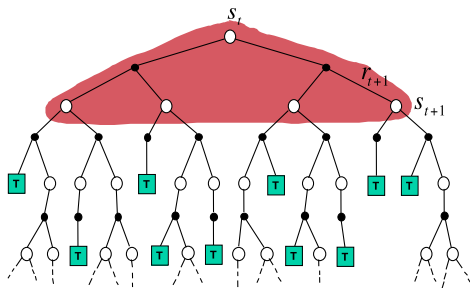


图 2: 图片来自David Silver 个人网站²

²<https://www.davidsilver.uk/teaching/>

时序差分算法

时序差分学习（Temporal-Difference Learning, TD Learning）算法是真正意义上的强化学习基础算法，而动态规划和蒙特卡罗算法都是经典的解决马尔科夫决策过程问题的方法。面对动态规划和蒙特卡罗方法的诸多问题，我们将介绍经典的强化学习的基础算法，由强化学习之父 Richard Sutton 教授提出的时序差分算法：

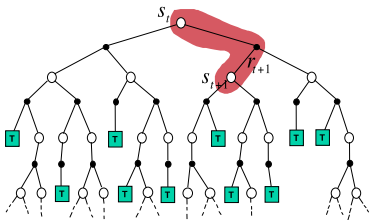


图 3: 图片来自David Silver 个人网站³

时序差分算法与动态规划和蒙特卡罗算法比较

我们简单比较时序差分算法、动态规划和蒙特卡罗算法这三个方法的状态值函数估计过程。在**动态规划方法**中，我们有：

$$V_\pi(S_t) = \mathbb{E}_\pi[R_t + \gamma V_\pi(S_{t+1}) | S_t = s]. \quad (57)$$

在蒙特卡罗方法中，我们有：

$$V_{\pi}(S_t) = \mathbb{E}_{\pi}[G_t | S_t = s]. \quad (58)$$

在时序差分方法中，我们有：

$$V_\pi(S_t) \approx [R_t + \gamma V_\pi(S_{t+1}) | S_t = s]. \quad (59)$$

Q-learning

强化学习算法Q-learning的状态-动作值函数更新公式为:

$$Q(s, a) = Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a)). \quad (60)$$

其中, r 表示智能体在状态 s 下动作 a 的即时回报, $\max_{a'} Q(s', a')$ 表示智能体跳转至下一个状态 s' 能获得的最大累积回报, 此过程需要遍历所有的动作。状态-动作值函数的更新公式 $r + \gamma \max_{a'} Q(s', a')$ 被称作**时序差分 (TD) 目标值**。式 (60) 与采用增量值更新状态-动作值函数的蒙特卡罗算法类似:

$$Q_n(s, a) = Q_{n-1}(s, a) + \alpha(G_n - Q_{n-1}(s, a)). \quad (61)$$

一条完整轨迹的累积回报 G_n 替换成 $r + \gamma \max_{a'} Q(s', a')$ 。

SARSA算法

Q-learning算法是强化学习中最负盛名的算法之一，但Q-learning算法也存在很多问题和不足，如过估计（Overestimation）问题。**时序差分SARSA算法**是另一个经典强化学习算法。SARSA算法与Q-learning算法的主要区别是状态-动作值函数更新公式：

SARSA : $Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma Q(s', a') - Q(s, a))$,

$$\text{Q-learning : } Q(s, a) \leftarrow Q(s, a) + \alpha \left(r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right). \quad (62)$$

策略梯度方法

智能体与环境交互得到的所有轨迹 τ 的期望收益为:

$$\bar{R}_\theta = \sum_{\tau} R(\tau) p_\theta(\tau). \quad (66)$$

基于策略梯度的**强化学习**的目标是找到最优化的策略函数参数 θ ，即最优策略 π_θ ，使得期望累积收益 \bar{R}_θ 最大，用数学语言描述为：

$$\pi_\theta = \arg \max_{\pi_\theta} \bar{R}_\theta = \arg \max_{\pi_\theta} \sum_{\tau} R(\tau) p_\theta(\tau). \quad (67)$$

显然, $\bar{R}_\theta = \sum_{\tau} R(\tau) p_\theta(\tau)$ 是策略梯度方法的目标函数, 可采用梯度上升法求解最优参数 θ 。

策略梯度方法

在分析策略梯度算法的过程中，需要用到对数函数求导公式：

$$\nabla \log f(x) = \frac{\nabla f(x)}{f(x)}. \quad (68)$$

此公式在策略函数推导过程中具有重要的作用，可以改写成：

$$\nabla f(x) = f(x) \nabla \log f(x). \quad (69)$$

将函数 $f(x)$ 替换成轨迹 τ 发生的概率 $p_{\theta}(\tau)$, 可得:

$$\nabla p_\theta(\tau) = p_\theta(\tau) \nabla \log p_\theta(\tau), \quad (70)$$

即

$$\frac{\nabla p_{\theta}(\tau)}{p_{\theta}(\tau)} = \nabla \log p_{\theta}(\tau). \quad (71)$$

周炜星 谢文杰

第4章：强化学习入门

华东理工大学金融学系

策略梯度方法

公式 $\nabla \bar{R}_\theta = \mathbb{E}_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)]$ 为策略梯度方法中更新策略函数参数的核心公式，如此转化公式的作用是将目标函数梯度的计算问题，转化成通过**蒙特卡罗采样完成梯度估计**：

$$\nabla \bar{R}_\theta = \mathbb{E}_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)] \approx \frac{1}{n} \sum_{i=1}^n R(\tau_i) \nabla \log p_\theta(\tau_i). \quad (73)$$

因此，策略梯度方法的策略函数参数更新公式可写作：

$$\theta \leftarrow \theta + \alpha \nabla \bar{R}_\theta, \quad (74)$$

其中 α 为学习率，决定了策略函数参数更新的步长大小。

策略梯度方法

策略梯度方法的核心公式需要计算 $\nabla \log p_{\theta}(\tau)$ ，对式(65)求梯度可得：

$$\begin{aligned}\nabla \log p_{\theta}(\tau) &= \nabla \left(\log p(s_0) + \sum_{t=1}^T \log \pi_{\theta}(a_t | s_t) + \sum_{t=1}^T \log p(s_{t+1} | s_t, a_t) \right) \\ &= \nabla \log p(s_0) + \nabla \sum_{t=1}^T \log \pi_{\theta}(a_t | s_t) + \nabla \sum_{t=1}^T \log p(s_{t+1} | s_t, a_t) \\ &= \nabla \sum_{t=1}^T \log \pi_{\theta}(a_t | s_t) \\ &= \sum_{t=1}^T \nabla \log \pi_{\theta}(a_t | s_t)\end{aligned}$$

策略梯度方法

在公式推导过程中有两个关键之处，分别为：

$$\nabla \log p(s_0) = 0 \quad (76)$$

以及

$$\nabla \sum_{t=1}^T \log p(s_{t+1}|s_t, a_t) = 0. \quad (77)$$

状态稳定分布函数 $p(s_0)$ 和环境状态转移函数 $p(s_{t+1}|s_t, a_t)$ 由环境模型决定，不会因为策略函数的变化而变化，因而与策略函数参数 θ 无关，其梯度均为0。在公式推导中，只有策略函数 $\pi_\theta(a_t|s_t)$ 包含了参数 θ 。

策略梯度方法

将式（75）代入目标函数梯度公式（73），可得：

$$\begin{aligned}
 \nabla \bar{R}_\theta &= \sum_{\tau} R(\tau) \nabla p_\theta(\tau) \\
 &= \mathbb{E}_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)] \\
 &\approx \frac{1}{n} \sum_{i=1}^n R(\tau^i) \nabla \log p_\theta(\tau^i) \\
 &= \frac{1}{n} \sum_{i=1}^n \sum_{t=1}^{T_i} R(\tau^i) \nabla \log \pi_\theta(a_t^i | s_t^i),
 \end{aligned} \tag{78}$$

其中， n 为随机采样的完整轨迹数量，变量上标 i 为轨迹编号， T_i 为第 i 条轨迹的长度。

蒙特卡罗策略梯度算法伪代码

Algorithm 13: 蒙特卡罗策略梯度（REINFORCE）算法伪代码

Input: 状态空间 \mathcal{S} , 动作空间 \mathcal{A} , 折扣系数 γ 以及环境 Env , 可微分策略函数 $\pi_{\theta}(a|s)$, 学习率 α

Output: 最优策略 π^*

1 初始化策略函数的参数 θ

2 **for** $n = 0, 1, 2, 3, \dots$ **do**

3 % 每次循环针对一条轨迹

4 初始化状态 s_0 , 生成一条轨迹

5

$$\tau = \{s_0, a_0, r_0, s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_T, a_T, r_T\} \quad (4.92)$$

6 **for** $t = 0, 1, 2, 3, \dots, T$ **do**

7 计算当前时间步开始到轨迹结束的累积回报 G_t

$$\theta \leftarrow \theta + \alpha G_t \nabla \log \pi_{\theta}(a_t|s_t) \quad (4.93)$$

8 **if** s 为终止状态 **then**

9 开始下一条轨迹采样

10 **if** θ 收敛 **then**

11 停止迭代

强化学习智能交易系统框架

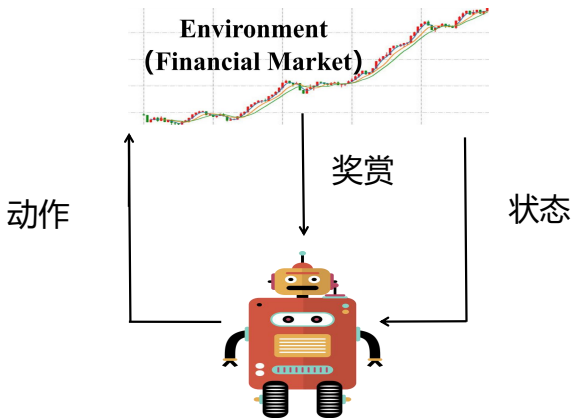


图 4: 基于强化学习的智能交易系统框架示意图

智能交易环境模型编程

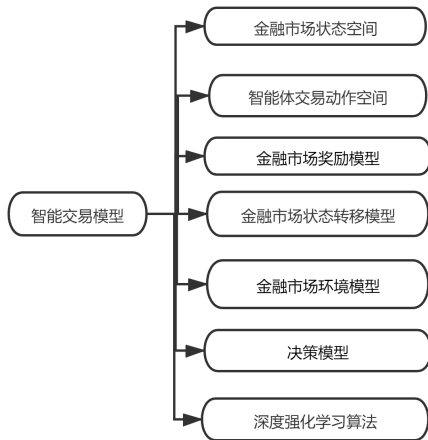
智能交易模型一直是量化金融的热门方向。相较于游戏环境的强化学习智能体建模，金融市场智能体建模更具挑战性。金融市场的复杂度远远高于一般的游戏环境系统。

基于离散马尔科夫决策过程 (Markov Decision Process, MDP) 的金融市场环境模型可以表示为六元组 (S, A, P, R, γ, H) , 其中:

- S 表示金融市场环境状态集合,
- A 表示智能体动作集合,
- $P: S \times A \times S \rightarrow [0, 1]$ 是金融市场环境状态转移函数,
- $R: S \times A \times S \rightarrow \mathcal{R}$ 是金融市场环境回报函数, \mathcal{R} 为连续的区间, $R(s_t, a_t, s_{t+1}) \in \mathcal{R}$, $R_{\max} \in \mathbb{R}^+$ (e.g., $[0, R_{\max}]$),
- $\gamma \in [0, 1)$ 是折扣系数,
- H 是投资期限。

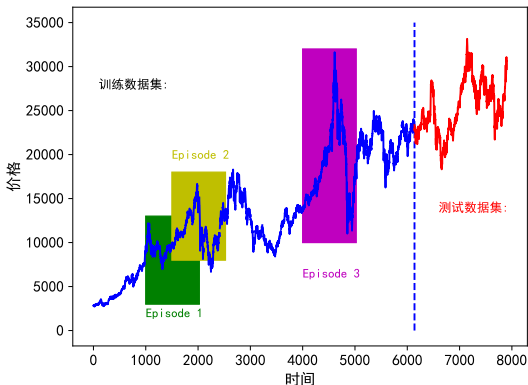
智能交易环境模型编程

在金融市场环境中，基于强化学习的智能交易系统模块：



金融市场状态空间

金融市场数据形式千变万化，如价格时间序列、交易网络、非结构化的新颖数据等，都能为个体和机构的投资决策提供市场环境信息和市场状态信息，可以作为投资智能体的决策变量。



智能体动作空间

在金融市场环境模型六元组中, A 表示智能体的动作空间, 即智能体金融交易动作。智能体基于金融市场环境变量和当前策略函数给出交易动作。一般来说, 智能体动作可以分成两类, 一类是连续型, 一类是离散型。如交易智能体的离散型动作 $a \in A$ 可表示为:

$$a = \begin{cases} 1, & \text{buy} \\ 0, & \text{hold} \\ -1, & \text{sell} \end{cases}, \quad (79)$$

离散型数值1、0和-1 分别表示买入、持有和卖出金融资产。

金融市场奖励模型

奖励函数 $R(s, a, s')$ 可以定义为投资组合总市值的变化量：

$$R(s, a, s') = v' - v, \tag{80}$$

其中，奖励函数 R 表示智能体在金融市场状态 s 下执行动作 a 并转化到下一个状态 s' 后获得的即时奖励值， v' 和 v 分别表示智能体在状态 s' 和 s 时资产市值。奖励函数 $R(s, a, s')$ 也可以定义为智能体动作前后投资组合市值的对数收益：

$$R(s, a, s') = \log(v') - \log(v). \tag{81}$$

在金融理论界和金融实务界，对数收益率更为常用。

决策模型

智能体的策略函数的功能是处理和分析金融市场环境信息，建模策略函数的主要工具是深度神经网络模型，如前馈神经网络、卷积神经网络、循环神经网络、图神经网络等。

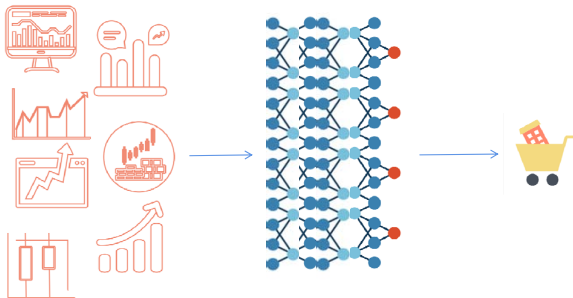


图 7: 智能体投资决策过程示意图

决策模型和强化学习算法

深度强化学习算法在强化学习理论和算法基础上融合了深度神经网络模型（DNN、CNN、RNN等）。

- Deep Q Network（DQN）
- 置信阈策略优化（Trust Region Policy Optimization, TRPO）
- 近端策略优化（Proximal Policy Optimization, PPO）
- 深度确定性策略梯度方法（Deep Deterministic Policy Gradient, DDPG）
- Twin Delayed DDPG（TD3）
- Actor-Critic算法等。

掌握的问题

- 1 什么是马尔科夫过程？
- 2 什么是马尔科夫回报过程？
- 3 什么是马尔科夫决策过程？
- 4 强化学习与动态规划的区别有哪些？
- 5 强化学习与蒙特卡罗方法的区别有哪些？
- 6 Q-learning和SARSA的区别？
- 7 基于值函数的强化学习算法有哪些优点？
- 8 基于值函数的和基于策略的强化学习的区别？