

### **强化学习:**

强化学习是机器学习的一个分支, 研究智能体如何在与环境的交互中学习最优行为策略, 以最大化累积奖励。智能体通过试错的方式, 根据环境的反馈 (奖励或惩罚) 来调整自己的行为, 从而逐步提高决策能力。

### **奖励:**

奖励是环境在每个时间步给予智能体的标量反馈信号, 用于指示当前行为的好坏。智能体的目标是最大化其在长期内获得的累计奖励, 奖励可以是正的, 也可以是负的。

### **状态:**

状态是对环境在某一时刻的完整或部分描述。智能体根据当前的状态来决定采取何种行动。状态包含了智能体作出决策所需的所有相关信息。

### **行动:**

行动是智能体在给定状态下可以执行的操作。智能体通过选择不同的行动来影响环境, 并从环境中获得反馈。

### **策略:**

策略是智能体从状态到行动的映射, 它定义了智能体在任何给定状态下应该如何行动。

### **值函数:**

值函数是衡量一个状态或一个状态-动作对的长期期望回报的函数。它表示从某个状态开始 (状态值函数  $V(s)$ ) 或在某个状态下执行某个动作后 (动作值函数  $Q(s, a)$ ), 智能体能够获得的未来累积奖励的期望。

## **Q-Learning**

Learning 是一种无模型、异策略 (pff-policy) 的时间差分强化学习算法。它通过学一个动作价值函数 ( $Q(s, a)$ ) 来找到最优策略,  $Q(s, a)$  表示在状态  $s$  下执行动作  $a$  所能获得的未来最大期望回报。Q-Learning 的更新规则是基于贝尔曼最优方程的。

## **TD 算法的定义**

时序差分学习算法是强化学习的一类核心算法, 通过结合动态规划和蒙特卡洛方法的优点, 实现高效、在线、无模型的值函数估计与策略优化。其核心思想是利用当前估计与未来估计的差异逐步更新值函数, 从而避免等待完整回合结束即可

实现增量学习。

### 深度 Q 网络:

Q-Learning 的深度学习扩展，用神经网络拟合 Q 值函数，解决高维状态空间问题。

### 策略梯度方法:

直接优化策略函数，通过梯度上升最大化期望回报

### TD( $\lambda$ )算法:

TD( $\lambda$ )算法是强化学习中一种结合了多步预测和资格迹的时序差分学习方法，其中  $\lambda$  是控制权重分配的衰减参数。他通过调整蒙特卡洛方法和单步 TD 学习之间的权衡，实现更高效的价值函数更新。

#### Q-learning算法中Q函数的定义及Q(s,a)的含义解析

##### 1. Q函数的定义

Q函数（动作价值函数）是Q-learning算法的核心，其数学定义为：

$$Q(s, a) = \mathbb{E}[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots | S_t = s, A_t = a]$$

其中：

- $s$  表示当前状态
- $a$  表示采取的动作
- $\gamma$  为折扣因子 ( $0 \leq \gamma \leq 1$ )
- $R$  表示即时奖励
- 期望值  $\mathbb{E}$  考虑所有可能的未来状态轨迹

### Q (s, a) 的具体含义:

Q (s, a) 表示在状态  $s$  下执行动作  $a$  后，智能体在整个 episode 中能获得的预期累积奖励（考虑折扣因子）。

### 支持向量机 (SVM) 与强化学习算法的核心区别

**SVM** 是监督学习算法，依赖标注数据（输入-输出对），强化学习是交互式学习，通过环境反馈（状态-动作-奖励）优化策略。

**SVM** 目标是最大化分类间隔，最小化分类误差，强化学习目标是最大化长期累积奖励。

**SVM** 需要静态数据集，强化学习需要动态交互数据。

**SVN** 输出分类边界或回归函数，强化学习输出策略或值函数。

### **QDN 和 Sarsa 的区别:**

**DQN** 是异策略算法，学习目标策略与行为策略可不同，**Sarsa** 是同策略算法，学习策略与行为策略必须一致。

**DQN** 动作选择方式采用直接学习最优策略，**Sarsa** 更新时依赖当前策略选择下一动作，更保守。

**DQN** 通过行为策略（如随机动作）探索，但目标 **Q** 值计算忽略探索，**Sarsa** 策略本身包含探索，更新时考虑探索带来的风险。

**DQN** 适合离散动作空间（如游戏 AI），追求最优策略。**Sarsa** 适合安全性要求更高（如机器人避障）的场景，避免过度乐观估计。

### 一步时间差分 (TD(0)) 算法的状态值函数更新公式

公式表达：

$$V(s_t) \leftarrow V(s_t) + \alpha [r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

符号说明：

- $V(s_t)$ ：状态  $s_t$  的当前估计值
- $\alpha$ ：学习率 ( $0 < \alpha \leq 1$ )，控制更新步长
- $r_{t+1}$ ：从状态  $s_t$  转移到  $s_{t+1}$  时获得的即时奖励
- $\gamma$ ：折扣因子 ( $0 \leq \gamma \leq 1$ )，衡量未来奖励的重要性
- $\delta_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$ ：TD误差，反映当前估计与目标的偏差

算法特点：

1. 增量更新：无需等待回合结束，单步交互后立即更新。
2. 自举 (Bootstrapping)：利用当前值函数估计  $V(s_{t+1})$  修正前一状态  $V(s_t)$ 。
3. 无模型 (Model-free)：仅需环境交互数据，无需已知状态转移概率。

### Q-learning算法更新公式详解

1. 标准更新公式：

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

2. 公式解析：

- $Q(s,a)$ ：状态s下采取动作a的动作价值函数
- $\alpha$  (alpha)：学习率 ( $0 < \alpha \leq 1$ )，控制更新幅度
- $r_{t+1}$ ：立即奖励
- $\gamma$  (gamma)：折扣因子 ( $0 \leq \gamma < 1$ )，衡量未来奖励重要性
- $\max_a Q(s_{t+1}, a)$ ：下一状态的最大预期价值

### 马尔可夫决策过程的定义：

马尔可夫决策过程是描述智能体在环境中做序贯决策的数学模型，由五元组  $\langle S, A, P, R, \gamma \rangle$  定义。

## 马尔可夫决策过程（MDP）的定义

### 1. 基本概念

马尔可夫决策过程是描述智能体在环境中做序贯决策的数学模型，由五元组  $\langle S, A, P, R, \gamma \rangle$  定义。

### 2. 核心要素

- 状态空间  $S$ ：所有可能的环境状态（如机器人位置、游戏画面）
- 动作空间  $A$ ：智能体可执行的动作集合（如移动、开火）
- 转移函数  $P$ ：  $P(s' | s, a)$  表示在状态  $s$  执行动作  $a$  后转移到  $s'$  的概率
- 奖励函数  $R$ ：  $R(s, a, s')$  是执行动作后获得的即时奖励
- 折扣因子  $\gamma$ ：  $0 \leq \gamma < 1$ ，用于平衡当前与未来奖励的重要性

### 3. 核心特性

- 马尔可夫性：下一状态仅依赖当前状态和动作，与历史无关。
- 目标：找到最优策略  $\pi^*$ ，最大化累积折扣奖励  $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$ 。

## 强化学习中“强化”的定义：

“强化”是指智能体通过环境反馈的奖励信号来调整行为策略的学习机制。其本质是试错学习与奖励最大化的结合。

## 策略的定义”

策略是智能体在给定状态下选择动作的规则，数学表示为派： $S \rightarrow A$

## 状态——动作对的定义：

在强化学习中，状态-动作对是指智能体在特定环境状态下采取的动作，是描述智能体决策行为的基本单元。

## 奖励信号的作用：

奖励信号是环境对智能体行为的即时反馈，通常为标量值  $r_t$  属于  $R$ ，表示在状态  $s_t$  下执行动作  $a_t$  后的优劣评价。

## 强化学习中的探索策略解析

### 1. $\epsilon$ -greedy策略

- **原理**：以概率  $\epsilon$  随机选择动作（探索），否则选择当前最优动作（利用）。
- **特点**：
  - 简单易实现，超参数  $\epsilon$  控制探索强度（通常随时间衰减）。
  - 可能重复探索低价值动作。

### 2. Boltzmann探索（Softmax策略）

- **原理**：按动作价值  $Q(s, a)$  的指数分布概率选择动作：

$$P(a|s) = \frac{e^{Q(s,a)/\tau}}{\sum_{a'} e^{Q(s,a')/\tau}}$$

- $\tau$  为温度参数： $\tau \rightarrow 0$  趋近贪婪策略， $\tau \rightarrow \infty$  趋近均匀随机。
- **特点**：平衡高价值动作与潜在优动作的探索。

### 3. Thompson采样

- **原理**：基于贝叶斯思想，维护动作价值的概率分布，每次采样一个分布实例并选择最优动作。
- **特点**：
  - 适合多臂老虎机等不确定性问题。
  - 计算成本较高，需跟踪概率分布。

### 4. 随机选择策略

- **原理**：完全随机选择动作（探索率100%）。
- **适用场景**：初期完全未知环境时快速收集数据。

## 强化学习中的值函数类型解析

### 1. 状态值函数 (State-Value Function)

- 定义：评估在策略 $\pi$ 下处于状态 $s$ 的长期价值
- 公式：  $V^\pi(s) = E_\pi[G_t \mid S_t = s]$
- 特点：反映状态的整体优劣，与具体动作无关

### 2. 状态-动作值函数 (Action-Value Function)

- 定义：评估在策略 $\pi$ 下于状态 $s$ 采取动作 $a$ 的长期价值
- 公式：  $Q^\pi(s,a) = E_\pi[G_t \mid S_t = s, A_t = a]$
- 特点：建立状态-动作对的直接价值映射

### 3. 优势函数 (Advantage Function)

- 定义：衡量特定动作相对于平均水平的优势
- 公式：  $A^\pi(s,a) = Q^\pi(s,a) - V^\pi(s)$
- 特点：正优势表示优于平均，负优势表示劣于平均

### 4. 策略值函数 (Policy-Value Function)

- 定义：评估整个策略 $\pi$ 的期望回报
- 公式：  $J(\pi) = E_{s_0}[V^\pi(s_0)]$
- 特点：用于直接优化策略参数

#### 学习率对结果的影响：

学习率的作用是控制参数更新的步长，决定新经验对旧估计的修正强度。学习率过大会导致参数震荡不稳定，难以收敛到最优解，可能错过最优策略。学习率过小会导致收敛速度极慢，易陷入局部最优解，训练资源浪费。

### 探索和利用的概念

探索：尝试新动作以发现潜在更高回报的行为（如随机选择未走过的路径）。

利用：基于当前知识选择已知最优动作（如重复执行已验证有效的策略）。

•  $\epsilon$ -greedy 算法中  $\epsilon$  值对采取随机动作和当前 Q 函数最大动作概率的影响。

以概率 $\epsilon$ 随机选择动作（探索），以  $1-\epsilon$  选择当前 Q 值最大的动作（利用）

$\epsilon=0$ ，100%选择最优动作（纯贪婪策略）

$0<\epsilon<1$ ，平衡探索（ $\epsilon$  概率随机）与利用（ $1-\epsilon$  概率最优）

$\epsilon=1$ ，100%随机动作（完全探索）

设高  $\epsilon$  值会加强探索，衰减  $\epsilon$  值侧重利用



### 1. 定义：

增量蒙特卡洛方法用于估计状态值函数 $V(s)$ ，通过逐步更新样本均值来优化估计值。其更新公式为：

$$V(s) \leftarrow V(s) + \alpha[G_t - V(s)]$$

其中：

- $V(s)$ ：状态 $s$ 的当前值估计
- $G_t$ ：从状态 $s$ 开始的完整回报（蒙特卡洛回报）
- $\alpha$ ：学习率（步长）

### 2. 推导过程：

#### a) 样本均值表达：

传统蒙特卡洛方法中，状态值估计为所有观测回报的算术平均：

$$V_n(s) = \frac{1}{n} \sum_{i=1}^n G_i$$

#### b) 递推形式转换：

将第 $n$ 次估计表示为第 $(n-1)$ 次估计的更新：

$$V_n(s) = V_{n-1}(s) + \frac{1}{n}(G_n - V_{n-1}(s))$$

这通过展开算术平均公式可得：

$$\frac{1}{n} \sum_{i=1}^n G_i = \frac{1}{n} G_n + \frac{n-1}{n} \cdot \frac{1}{n-1} \sum_{i=1}^{n-1} G_i$$

#### c) 推广到变步长：

将固定步长 $1/n$ 推广为一般步长 $\alpha$ ：

$$V(s) \leftarrow V(s) + \alpha[G_t - V(s)]$$

当 $\alpha=1/n$ 时即为传统蒙特卡洛平均。

### 3. 收敛性保证：

- 当步长满足Robbins-Monro条件：

$$\sum_{n=1}^{\infty} \alpha_n = \infty, \quad \sum_{n=1}^{\infty} \alpha_n^2 < \infty$$

- 典型选择： $\alpha_n = 1/n$ （保证收敛到真值）