

# 强化学习实验报告

多臂赌博机与gymnasium

张恒硕

2212266

人工智能学院

2025 年 3 月 6 日



南开大学  
Nankai University



## 目 录

|     |                |   |
|-----|----------------|---|
| 1   | 实验目的           | 3 |
| 2   | 基础知识           | 3 |
| 2.1 | 探索-利用平衡策略      | 3 |
| 2.2 | 马尔可夫决策过程 (MDP) | 4 |
| 3   | 多臂赌博机          | 4 |
| 3.1 | 代码改动           | 4 |
| 3.2 | 实验结果与分析        | 4 |
| 4   | gymnasium      | 6 |
| 4.1 | 源码分析           | 6 |
| 4.2 | 实验结果与分析        | 8 |
| 5   | 实验总结           | 8 |

## 图 片

|     |                                 |   |
|-----|---------------------------------|---|
| 图 1 | 三种策略对比                          | 5 |
| 图 2 | 小车爬坡 (MountainCarContinuous-v0) | 8 |

## 表 格

|     |               |   |
|-----|---------------|---|
| 表 1 | 三种策略的 Q 值估计对比 | 5 |
|-----|---------------|---|

## 1 实验目的

1. 以多臂赌博机程序为例，了解三种探索-利用平衡策略。
2. 将多臂赌博机程序由三臂改成四臂，进行调试，对比三种动作选择策略。
3. 挑选一个gymnasium包中的环境，分析其MDP元素。

## 2 基础知识

### 2.1 探索-利用平衡策略

#### EPSILON-GREEDY策略

以 $\epsilon$ 的概率进行随机探索，以 $1 - \epsilon$ 的概率选择当前估计都最优动作。  
特点是简单高效，但依赖固定 $\epsilon$ 。

#### UPPER CONFIDENCE BOUND(UCB)策略

通过置信上限平衡探索和利用，有限不确定性高的动作。

$$UCB(a) = Q(a) + C \sqrt{\frac{\ln N}{n(a)}}$$

其中， $C$ 是控制探索强度的超参数， $N$ 是当前轮数， $n(a)$ 是动作 $a$ 被选次数。  
特点是自适应平衡探索与利用。

#### BOLTZMANN策略

基于温度参数生成概率分布。

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

其中， $\tau$ 是温度参数，其偏大时偏向探索，偏小时趋向贪心。  
特点是温度敏感，可以动态调整探索强度。

## 2.2 马尔可夫决策过程（MDP）

1. 状态 $s$ ：智能体当前的处境或信息。
2. 动作 $a$ ：智能体在每个状态下可执行的操作。
3. 转移概率 $p$ ：动作 $a$ 导致从 $s$ 转移到 $s'$ 的概率，
$$P(s' | s, a) = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a]$$
。
4. 奖励 $r$ ：在状态 $s$ 执行动作 $a$ 后，智能体获得的即时奖励。
5. 折扣因子 $\gamma$ ：衡量未来奖励的重要性。

## 3 多臂赌博机

### 3.1 代码改动

只添加了一个新的臂，并修改了相应代码。在设置 $q$ 值时，四种动作依次为1、1.5、2、2.5。改动较为简单，这里不做分析。

### 3.2 实验结果与分析

四臂赌博机的实验结果如下图图1 on the following page和下表表1 on the next page:

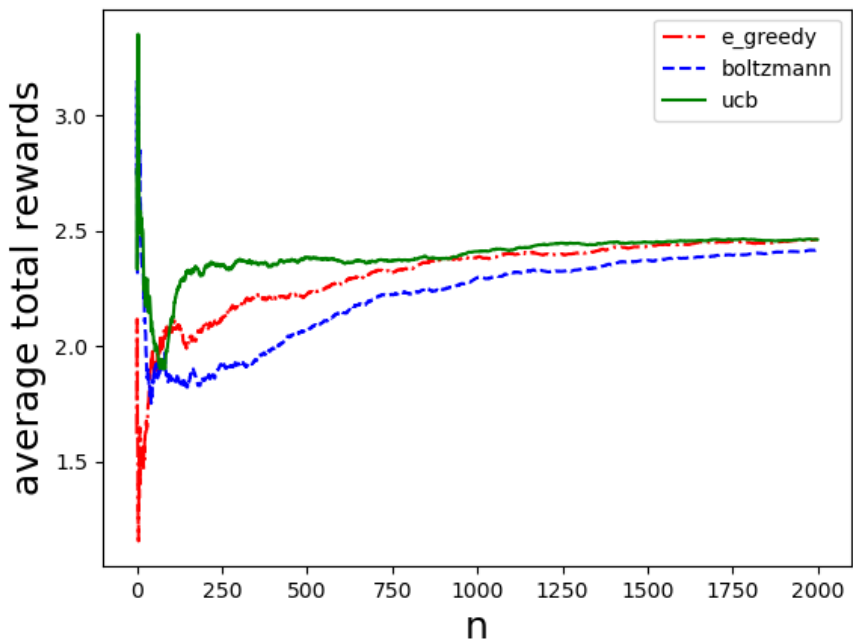


图 1: 三种策略对比

表 1: 三种策略的 Q 值估计对比

| 动作    | 实际值 | Epsilon-Greedy     | UCB                | Boltzmann          |
|-------|-----|--------------------|--------------------|--------------------|
| $a_1$ | 1   | 0.79600189<br>31   | 1.01580713<br>11   | 0.7953477<br>34    |
| $a_2$ | 1.5 | 1.44023784<br>37   | 1.58757369<br>16   | 1.4326053<br>69    |
| $a_3$ | 2   | 1.65791545<br>24   | 2.09857803<br>28   | 1.95786096<br>102  |
| $a_4$ | 2.5 | 2.51713636<br>1908 | 2.48150385<br>1945 | 2.50636821<br>1795 |

可以获得以下信息：

1. 三种方法的估计值都近似实际值，但在近似程度上有差异。
2. Epsilon-Greedy中，非最大奖励动作的三种动作的选择次数接近，这是不同于另外两种策略的，这与它探索是随机选取的方式有关。

3. UCB是三种方法中收敛最快的，可以发现，它选取低奖励动作的概率较小，这与它的选取公式有关。
4. Boltzmann的收敛速度和结果都是最差的，这与 $\tau$ 值的选取有关，若能动态地实现温度曲线式的参数变化，将能取得更好的结果。

## 4 GYMNASIUM

加载了MountainCarContinuous-v0环境，其对应的源码位于continuous\_mountain\_car.py。这是一个小车爬坡的模拟过程，通过设定小车的运动，使其爬上正确方向的坡并到达指定位置。

### 4.1 源码分析

#### 状态空间

```

1 self.observation_space = spaces.Box(
2     low=self.low_state, high=self.high_state, dtype=np.float32
3 )
4 where:
5 self.low_state = np.array([-1.2, -0.07])
6 self.high_state = np.array([0.6, 0.07])

```

本部分定义小车可行的位置空间为 $[-1.2, 0.6]$ ，速度空间为 $[-0.07, 0.07]$ 。

#### 动作空间

```

1 self.action_space = spaces.Box(
2     low=self.min_action, high=self.max_action, shape=(1,), dtype=np.float32
3 )
4 where:
5 self.min_action = -1.0, self.max_action = 1.0

```

本部分限制小车动作在 $[-1.0, 1.0]$ 范围内。

## 奖励函数

```
1 reward = 100.0 if terminated else -math.pow(action[0], 2) * 0.1
```

本部分设定了小车运动过程中的奖励。当位置 $\geq 0.45$ 且速度 $\geq \text{goal\_velocity}$ 时，奖励为+100，成功抵达目标。在抵达前，每一步奖励为 $-0.1(\text{action}^2)$ ，其为负值惩罚，随动作数增加，以督促小车尽快抵达目标。

## 状态转移

本部分定义小车的状态转移为确定性转移，其速度和位置更新公式如下：

$$\text{velocity}_{t+1} = \text{clamp}(\text{velocity}_t + \text{force} * \text{power} - 0.0025 * \cos(3 * \text{position}_t), [-0.07, 0.07])$$

$$\text{position}_{t+1} = \text{clamp}(\text{position}_t + \text{velocity}_{t+1}, [-1.2, 0.6])$$

另外，还对碰撞进行了处理：若位置达到边界且速度方向与边界相反，则速度置零。

**折扣率** 源码中未显式定义。

## 初始状态

```
1 self.state = np.array([
2     self.np_random.uniform(low=low, high=high),
3     0.0
4 ])
```

本部分初始化了小车状态，其位置均匀分布在 $[-0.6, -0.4]$ ，速度固定为0。

## 4.2 实验结果与分析

运行过程如下图图2所示：

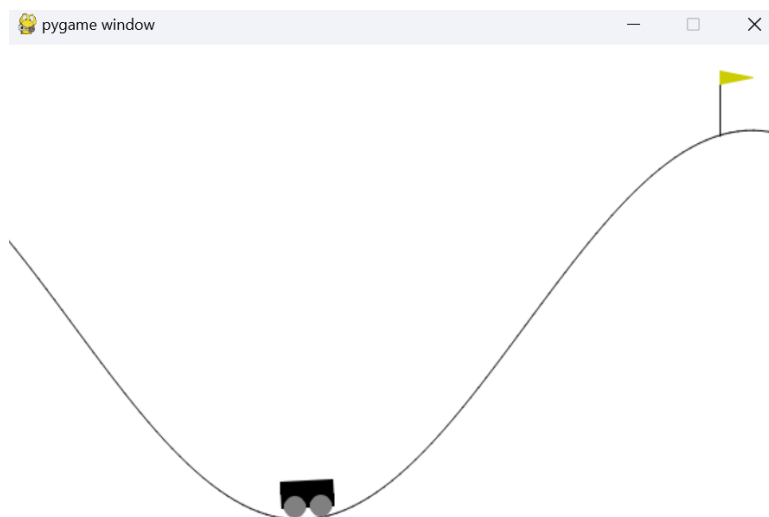


图 2: 小车爬坡 (MountainCarContinuous-v0)

结合代码和实验过程可以获得以下信息：

1. 奖励与折扣率：本实验为使小车尽快抵达目标，设定过程动作的奖励为负，且随动作数增加，这是一种变向的折扣率。
2. 违规动作处理：本实验中对于碰壁有所处理，属于理论课中介绍的禁止与惩罚两种方式的前者。
3. 合理设置状态空间和动作空间：如本实验，对小车的速度上限有所限制，这对保证实验的连续性和合理性有益。

## 5 实验总结

通过本次实验的学习，掌握了两方面知识。一方面是从多臂赌博机出发，了解了三种探索-利用平衡策略，这是强化学习中动作选择的基础策略。另一方面，通过阅读gymnasium包中特定环境的源代码，深化了对马尔可夫链的认识，并初步了解了这个常用的强化学习包。