

## 第二章作业

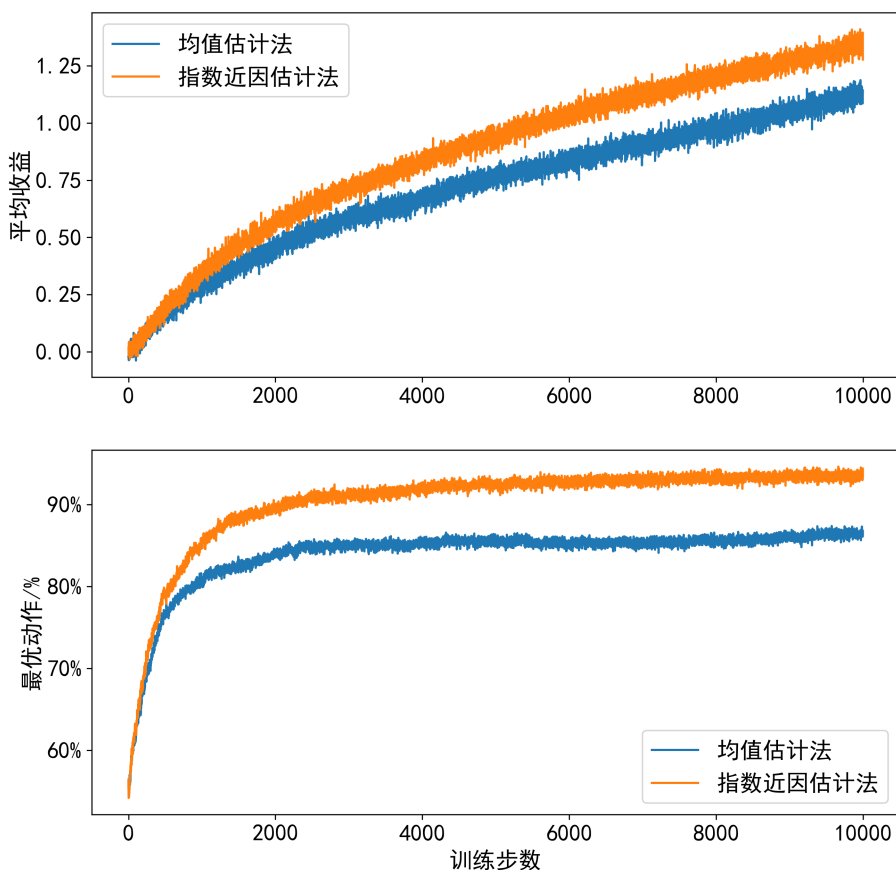
**题目 1. 练习 2.2** 考虑一个  $k = 4$  的多臂赌博机问题，记做 1, 2, 3, 4. 使用  $\varepsilon$ -贪心动作选择，基于均值的动作价值估计，初始估计  $Q_1(a) = 0, \forall a = \{1, 2, 3, 4\}$ ，假设动作及收益的序列为

$$A = \{1, 2, 2, 2, 3\}, R = \{-1, 1, -2, 2, 0\}$$

其中某些时刻中，可能发生了进行了探索 ( $\varepsilon$ -时刻)，即随机选择一个动作. 请问上述时间序列中，有哪些时刻肯定进行了探索？那些时刻可能进行了探索？

**解答.** 上述序列对应的每个时刻的动作价值  $Q_t(a)$  变换如下表所示，从该表容易得出，第 4, 5 步的选择一定是进行了探索，而所有时刻都可能进行探索，因为在探索中，每个动作被选择到的概率是相同的.

$t$	1	2	3	4	5	6
$Q_t(1)$	0	-1	-1	-1	-1	-1
$Q_t(2)$	0	0	1	-0.5	1/3	1/3
$Q_t(3)$	0	0	0	0	0	0
$Q_t(4)$	0	0	0	0	0	0
$A_t$	1	2	2	2	3	



**题目 2. 练习 2.5** 设计实验来证实使用均值估计方法取解决非平稳问题的困难，使用一个 10 臂赌博机，其中所有的  $q_*(a)$  初始时均相等，然后进行随机游走，每一步所有的  $q_*(a)$  都加上一个服从  $N(0, 0.01^2)$  的增量，分别使用均值估计方法和指数近因加权估计方法且步长  $\alpha = 0.1$  进行决策，采用  $\varepsilon$ -贪心进行动作选择，且总步数为  $T = 10000$ 。

**解答.** 如上图所示进行了 2000 次不同的多臂赌博机实验的平均结果，从中非常容易得出，指数近因估计在处理多臂赌博机问题上比均值估计要好。

---

```

1 from tqdm import tqdm
2 import numpy as np
3 import matplotlib as mpl
4 import matplotlib.pyplot as plt
5
6 config = { # matplotlib 绘图配置
7     "figure.figsize": (6, 6), # 图像大小
8     "font.size": 16, # 字号大小
9     "font.sans-serif": ['SimHei'], # 用黑体显示中文
10    'axes.unicode_minus': False # 显示负号
11 }
12 plt.rcParams.update(config)
13
14 class Bandit: # 赌博机类
15     def __init__(self, n, init_means=0) -> None:
16         # init_means: 赌博机初始均值
17         self.n = n # 赌博机臂个数
18         self.means = np.full(self.n, init_means, dtype='float') # 均值初始值
19
20     def move_state(self):
21         return np.random.normal(0, 0.01, self.n) # 随机游走变换函数
22
23     def step(self, action): # 执行一个 action, 并返回其收益
24         assert(0 <= action < self.n)
25         self.means += self.move_state()
26         mean_now = self.means[action]
27         rate = np.sum(self.means <= mean_now) / self.n
28         return np.random.normal(self.means[action], size=1)[0], rate
29
30 def train(strategy, alpha=0.1):
31     # 两种策略 average 和 recency
32     bandit = Bandit(n)
33     N = np.zeros(n)
34     Q = np.zeros(n)
35     history = [[], []] # 存储每个时刻的收益和最优动作占比
36     for t in range(T):
37         if np.random.random(1)[0] < epsilon: # 随机选取一个动作
38             action = np.random.randint(0, n)
39         else: # 贪心选择最大价值的动作
40             action = np.argmax(Q)
41         reward, rate = bandit.step(action)
42         delta = reward - Q[action]
43         N[action] += 1

```

```

44         Q[action] += delta / N[action] if strategy == 'average' else delta *
         ↪ alpha
45         history[0].append(reward)
46         history[1].append(rate)
47     return np.array(history)
48
49 def plot_reward_rate(strategy, m=2000): # m 为实验次数
50     history = np.zeros((2, T))
51     for _ in tqdm(range(m)):
52         history += train(strategy=strategy)
53     history /= m
54     ax[0].plot(range(T), history[0], label='均值估计法' if strategy=='average'
55     ↪ else '指数近因估计法')
56     ax[1].plot(range(T), history[1], label='均值估计法' if strategy=='average'
57     ↪ else '指数近因估计法')
58
59 if __name__ == '__main__':
60     n, T = 10, 10000
61     epsilon = 0.1
62     np.random.seed(42)
63     fig, ax = plt.subplots(2, 1, figsize=(10, 10))
64
65     plot_reward_rate('average', m=2000)
66     plot_reward_rate('recency', m=2000)
67
68     ax[0].set_ylabel('平均收益')
69     ax[1].set_ylabel('最优动作/%')
70     ax[1].set_xlabel('训练步数')
71     ax[1].yaxis.set_major_formatter(mpl.ticker.PercentFormatter(xmax=1,
72     ↪ decimals=0))
73     ax[0].legend()
74     ax[1].legend()
75     plt.savefig('33 页练习 2.5.png', dpi=300)
76     plt.show()

```

---