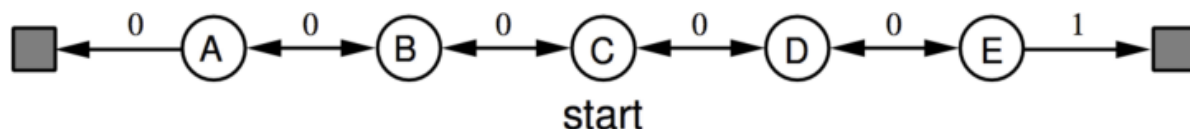


# 强化学习基础篇（十九）TD与MC在随机游走问题应用

为了比较讨论一下TD与MC方法，本位简单探索在一个随机游走的示例中，TD与MC的差异。

## 1、随机行走（Random Walk）问题设定



**状态空间：**如上图A、B、C、D、E为中间状态，C同时作为起始状态。灰色方格表示终止状态；

**行为空间：**除终止状态外，任一状态可以选择向左、向右两个行为之一；

**即时奖励：**右侧的终止状态得到即时奖励为1，左侧终止状态得到的即时奖励为0，在其他状态间转化得到的即时奖励是0；

**状态转移：**100%按行为进行状态转移，进入终止状态即终止；

**衰减系数：**1；

**给定的策略：**随机选择向左、向右两个行为。

**问题：**对这个MDP问题进行预测，也就是评估随机行走这个策略的状态。

## 2、初始化

首先导入需要用到的库函数

```
1 # 导入库函数
2 import numpy as np
3 import matplotlib
4 import matplotlib.pyplot as plt
5 from tqdm import tqdm
6 ## 解决matplotlib中文画图乱码问题
7 from pylab import mpl
8 mpl.rcParams['font.sans-serif'] = ['SimHei']
```

按照问题的设定，状态空间一共为七个： $S \in \{left\_terminate, A, B, C, D, E, right\_terminate\}$ 。

由于这个任务没有折扣，所以每个状态的真实价值是从这个状态开始并终止与最右侧的概率。因此中心状态的真实价值为 $v_{\pi}(C) = 0.5$ 。状态A-E的真实价值分别为：1/6, 2/6, 3/6, 4/6以及5/6。

```

1  # 定义七个状态，初始化A-E的值为0.5，右边终点值为1.
2  VALUES = np.zeros(7)
3  VALUES[1:6] = 0.5
4  VALUES[6] = 1
5  # 由于这个任务没有折扣，所以每个状态的真实价值是从这个状态开始并终止与最右侧的概率。
6  # 因此中心状态的真实价值为0.5。
7  # 状态A-E的真实价值分别为：1/6, 2/6, 3/6, 4/6, 5/6
8  TRUE_VALUE = np.zeros(7)
9  TRUE_VALUE[1:6] = np.arange(1, 6) / 6.0
10 TRUE_VALUE[6] = 1
11 # 定义向左与向右两个动作
12 ACTION_LEFT = 0
13 ACTION_RIGHT = 1

```

### 3、定义TD方法的实现

该代码实现遵循表格型TD(0)算法伪代码：

#### Tabular TD(0) for estimating $v_\pi$

Input: the policy  $\pi$  to be evaluated  
 Algorithm parameter: step size  $\alpha \in (0, 1]$   
 Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$   
 Loop for each episode:  
   Initialize  $S$   
   Loop for each step of episode:  
    $A \leftarrow$  action given by  $\pi$  for  $S$   
   Take action  $A$ , observe  $R, S'$   
    $V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$   
    $S \leftarrow S'$   
 until  $S$  is terminal

```

1  def temporal_difference(values, alpha=0.1, batch=False):
2      # 定义开始点为C，即state=3
3      state = 3
4      # 定义轨迹列表
5      trajectory = [state]
6      # 定义奖励列表
7      rewards = [0]
8      while True:
9          old_state = state
10         # 通过一个二项分布，随机选择一个动作，并按照动作更新状态
11         if np.random.binomial(1, 0.5) == ACTION_LEFT:
12             state -= 1
13         else:
14             state += 1
15         # 按照问题定义，处理右边终点，其余的奖励都是0。
16         reward = 0
17         # 将state状态加入trajectory列表中
18         trajectory.append(state)
19         # 进行TD更新
20         if not batch:
21             values[old_state] += alpha * (reward + values[state] -
22             values[old_state])
23         # 遇到终结点则结束该次的episode。

```

```

23         if state == 6 or state == 0:
24             break
25         rewards.append(reward)
26     return trajectory, rewards

```

- 在决定随机行走的动作过程中，这里使用了二项分布  $X \sim b(n, p)$ , 这里即选择的为  $X \sim b(1, 0.5)$

```

1 | np.random.binomial(1, 0.5)

```

- TD的更新过程为:  $V(S) \leftarrow V(S) + \alpha(R + \gamma V(S') - V(S))$

```

1 | values[old_state] += alpha * (reward + values[state] - values[old_state])

```

## 4、定义MC方法的实现

```

1 | def monte_carlo(values, alpha=0.1, batch=False):
2 |     # 定义开始点为C, 即state=3
3 |     state = 3
4 |     # 定义轨迹列表
5 |     trajectory = [3]
6 |
7 |     # 如果最终是在左边介绍, 那么回报是0。
8 |     # 如果最终是在右边介绍, 那么回报是1。
9 |     while True:
10 |         # 通过一个二项分布, 随机选择一个动作, 并按照动作更新状态
11 |         if np.random.binomial(1, 0.5) == ACTION_LEFT:
12 |             state -= 1
13 |         else:
14 |             state += 1
15 |         trajectory.append(state)
16 |
17 |         if state == 6:
18 |             returns = 1.0
19 |             break
20 |         elif state == 0:
21 |             returns = 0.0
22 |             break
23 |     # 在episode完成后进行MC更新。
24 |     if not batch:
25 |         for state_ in trajectory[:-1]:
26 |             # MC update
27 |             values[state_] += alpha * (returns - values[state_])
28 |     return trajectory, [returns] * (len(trajectory) - 1)

```

- MC更新的方式遵循:  $V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$

```

1 | values[state_] += alpha * (returns - values[state_])

```

## 5、计算价值函数

这里考虑在episode为[0,1,10,100]四种情况下的估计价值，我们会将运行一次TD(0)所得到的价值估计值和真实值进行比较

```

1 def compute_state_value():
2     episodes = [0, 1, 10, 100]
3     current_values = np.copy(VALUE)
4     plt.figure(1)
5     for i in range(episodes[-1] + 1):
6         if i in episodes:
7             plt.plot(current_values, label=str(i) + ' episodes (幕的次数)')
8             temporal_difference(current_values)
9     plt.plot(TRUE_VALUE, label='true valuesd(真实价值)')
10    plt.xlabel('state (状态)')
11    plt.ylabel('estimated value (估计价值)')
12    plt.legend()

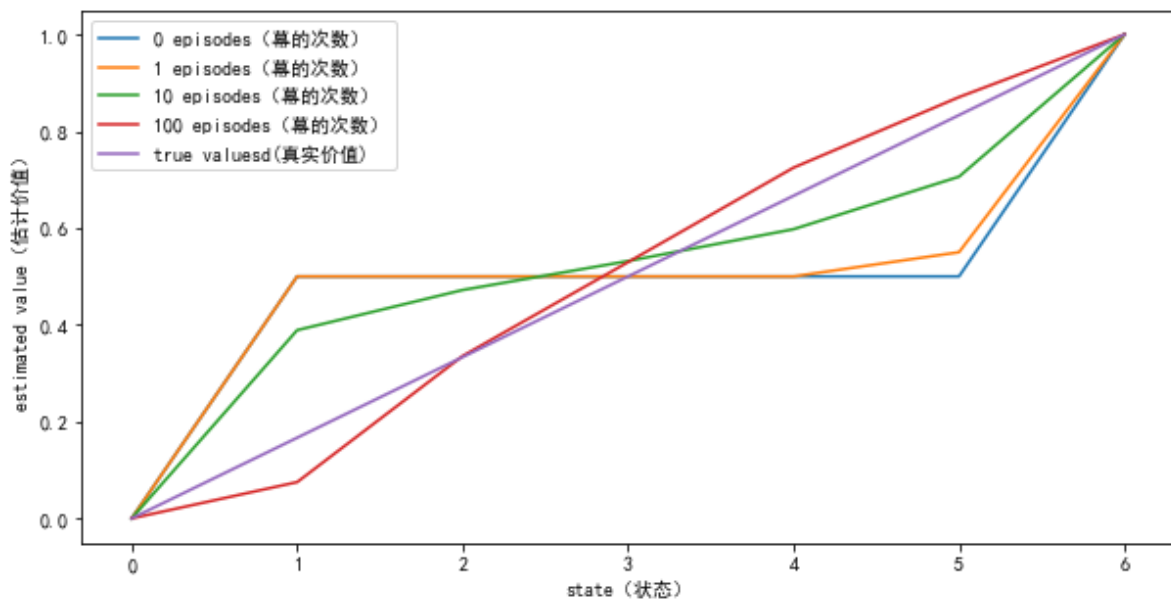
```

以下为运行结果，中间紫色线条为真实价值。我们可以看到在100幕后，估计值就非常接近于真实值了。这里我们使用了默认的步长参数 $\alpha = 0.1$ 。

```

1 plt.figure()
2 compute_state_value()
3 plt.show()

```



## 6、不同状态下平均经验均方根误差

上面只是简单比较了TD在运行过程中估计价值的变化，接下来我们考虑不同步长参数设置的情况下，MC与TD在不同步长参数下平均经验均方根误差变化情况。

这里我们将比较TD的三种步长参数 [0.15, 0.1, 0.05] 以及MC的四种步长参数 [0.01, 0.02, 0.03, 0.04]，他们在100幕运行过程

```

1 def rms_error():
2     # 设置TD与MC的步长参数
3     td_alphas = [0.15, 0.1, 0.05]
4     mc_alphas = [0.01, 0.02, 0.03, 0.04]
5     # 设定总episode数量
6     episodes = 100 + 1
7     runs = 100
8     # 遍历每个alpha设置

```

```

9     for i, alpha in enumerate(td_alphas + mc_alphas):
10         total_errors = np.zeros(epochs)
11         if i < len(td_alphas):
12             method = 'TD'
13             linestyle = 'solid'
14         else:
15             method = 'MC'
16             linestyle = 'dashdot'
17         # 这里整个过程一共运行100次，每次都是100幕，最后会对结果进行平均。
18         for r in tqdm(range(runs)):
19             errors = []
20             current_values = np.copy(VVALUES)
21             for i in range(0, epochs):
22                 # 计算当次幕下当前估计值和真实值之间的均方根误差。
23                 errors.append(np.sqrt(np.sum(np.power(TRUE_VALUE -
24                 current_values, 2)) / 5.0))
25                 if method == 'TD':
26                     temporal_difference(current_values, alpha=alpha)
27                 else:
28                     monte_carlo(current_values, alpha=alpha)
29                 total_errors += np.asarray(errors)
30             total_errors /= runs
31             plt.plot(total_errors, linestyle=linestyle, label=method + ', alpha
32             = %.02f' % (alpha))
33             plt.xlabel('epochs')
34             plt.ylabel('RMS')
35             plt.legend()

```

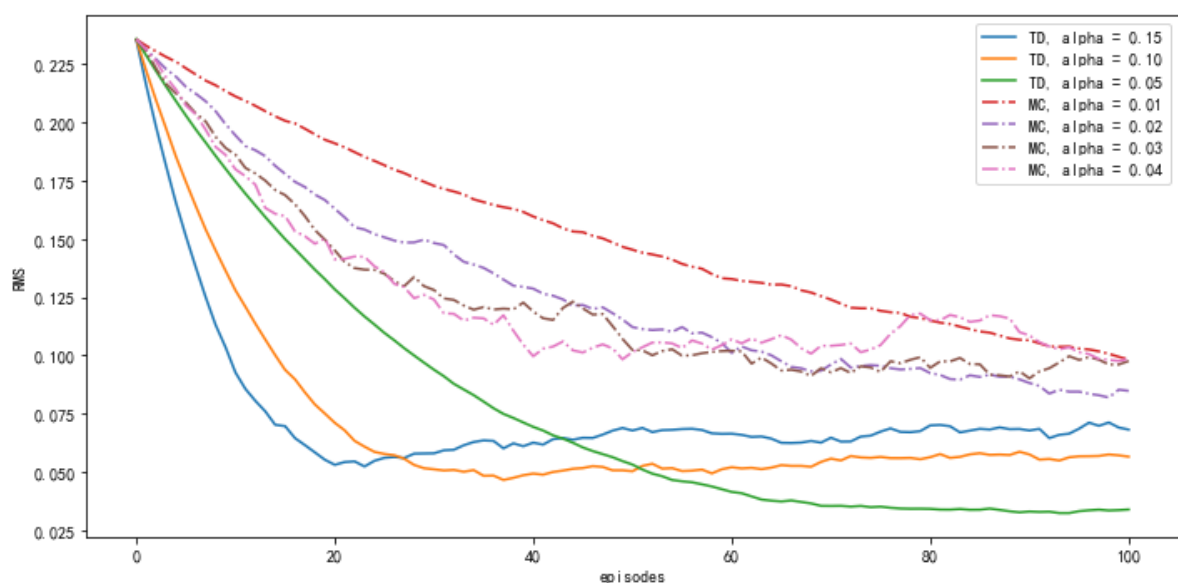
运行如下代码：

```

1 plt.figure(figsize=(10,5))
2 rms_error()
3 plt.tight_layout()

```

运行的最终结果如下：



我们可以看到对于不同的 $\alpha$ 取值，两种方法的学习曲线。图中显示的性能衡量指标是学到的价值函数和真实价值函数的均方根（RMS）误差。图中显示的误差是在5个状态上的平均误差，并在100次运行中取平均的结果。在所有情况下，对于所有 $s$ ，近似价值函数都被初始化为中间值 $V(s) = 0.5$ 。在这个任务中，TD方法一直比MC方法要好。

## 7、批量更新的随机游走

在随机游走问题中，批量更新版本的 $TD(0)$ 和常数 $\alpha$  MC方法的过程是这样的：每经过新的一幕序列之后，之前所有幕的数据就放视为一个批次。算法 $TD(0)$ 常数 $\alpha$  MC方法不断地使用这些批次进行逐次更新。这里 $\alpha$ 要设置得足够小以使价值函数能够收敛。最后将所得的价值函数与 $v_\pi$ 进行比较，绘制5个状态下的平均均方根误差（以整个实验的100次的独立重复为基础）的学习曲线。

```
1 def batch_updating(method, episodes, alpha=0.001):
2     # 整个实验进行100次独立重复运行
3     runs = 100
4     total_errors = np.zeros(episodes)
5     for r in tqdm(range(0, runs)):
6         current_values = np.copy(VALUE)
7         errors = []
8         # trajectories需要记录所有episode以及奖励
9         trajectories = []
10        rewards = []
11        for ep in range(episodes):
12            # 执行TD(0)
13            if method == 'TD':
14                trajectory_, rewards_ = temporal_difference(current_values,
15 batch=True)
16            # 执行MC
17            else:
18                trajectory_, rewards_ = monte_carlo(current_values,
19 batch=True)
20            trajectories.append(trajectory_)
21            rewards.append(rewards_)
22            while True:
23                # 持续不断得将到目前为止所有的trajectories都用于训练。
24                updates = np.zeros(7)
25                for trajectory_, rewards_ in zip(trajectories, rewards):
26                    for i in range(0, len(trajectory_) - 1):
27                        if method == 'TD':
28                            updates[trajectory_[i]] += rewards_[i] +
29 current_values[trajectory_[i + 1]] - current_values[trajectory_[i]]
30                        else:
31                            updates[trajectory_[i]] += rewards_[i] -
32 current_values[trajectory_[i]]
33                    updates *= alpha
34                    # 当接近收敛时才停止
35                    if np.sum(np.abs(updates)) < 1e-3:
36                        break
37                    # 进行批量更新
38                    current_values += updates
39                # 计算rms
40                errors.append(np.sqrt(np.sum(np.power(current_values -
41 TRUE_VALUE, 2)) / 5.0))
42            total_errors += np.asarray(errors)
43        total_errors /= runs
44    return total_errors
```

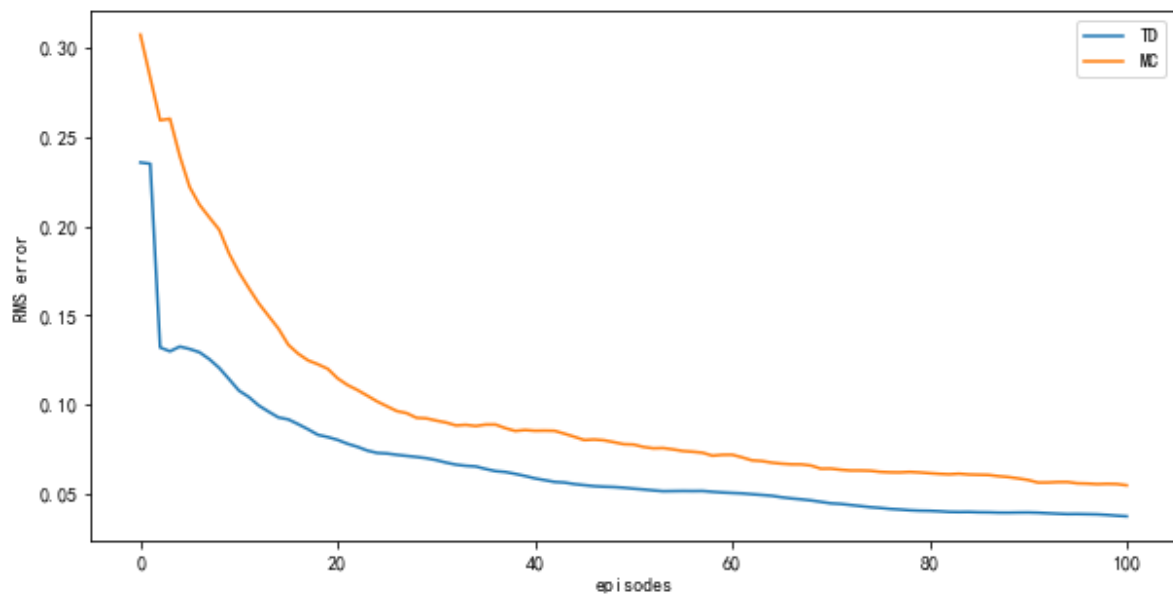
执行以下代码检查结果：

```

1 episodes = 100 + 1
2 # 运行TD (0) 的批量更新
3 td_erros = batch_updating('TD', episodes)
4 # 运行MC的批量更新
5 mc_erros = batch_updating('MC', episodes)
6
7 # 画图
8 plt.plot(td_erros, label='TD')
9 plt.plot(mc_erros, label='MC')
10 plt.xlabel('episodes')
11 plt.ylabel('RMS error')
12 plt.legend()
13 plt.show()

```

测试结果如下：



测试的实验结果可以看出，批量TD的rms始终是低于MC方法的，批量TD方法始终优于批量蒙特卡洛方法。其原因在蒙特卡洛方法只是从某些有限的方面来说是最优的，而TD方法的最优性则与预测回报这个任务更为相关。