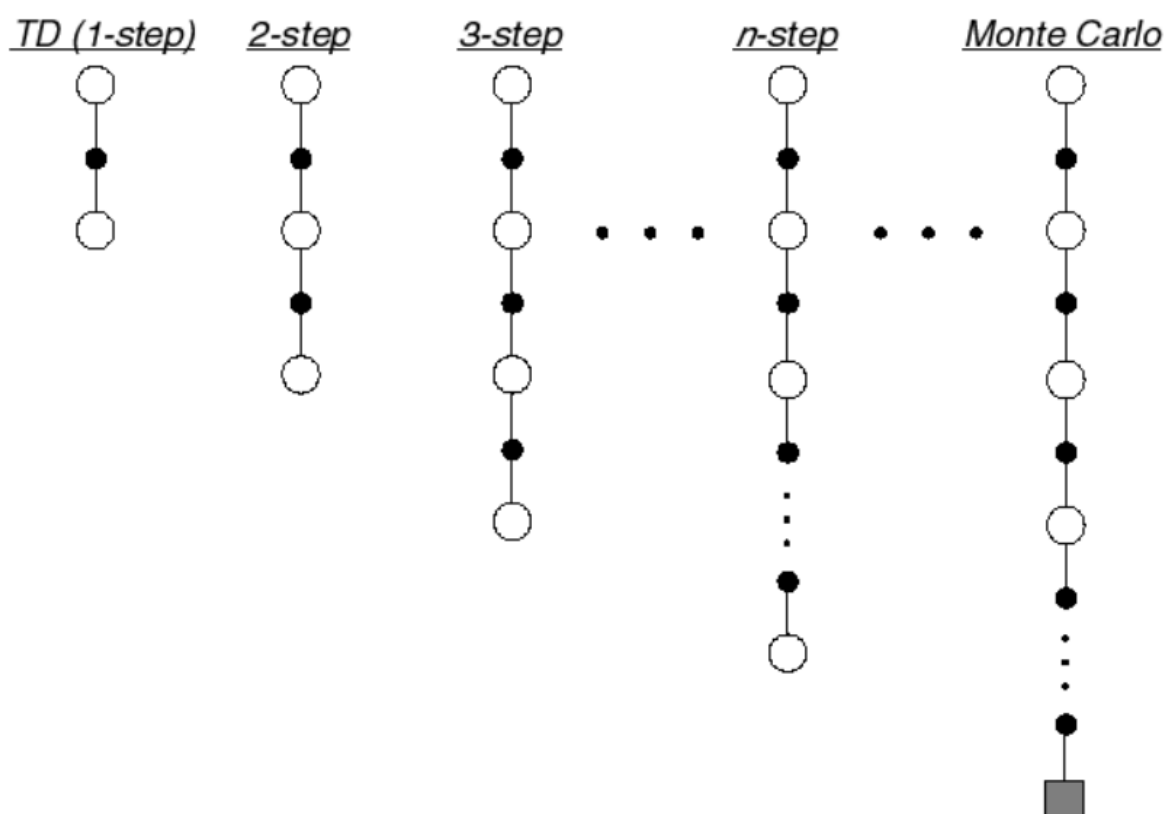


强化学习基础篇（二十五）n步时序差分预测

1、n步时序差分方法

之前在《强化学习基础篇（十七）时间差分预测》所介绍的是 $TD(0)$ 算法，其更新过程仅仅依赖于当前状态向下走一步的情况，将走一步走后的状态价值用于bootstrap更新。而蒙特卡洛方法是根据当前状态开始到终止状态的整个收益序列进行状态价值的更新。这节介绍的n步时序差分(n-step TD)是基于 $TD(0)$ 的一步更新与MC对整个序列进行更新的两个极端之间的算法。从n步时序差分方法的回溯图中，我们可以看到每个n步方法都考虑了从当前状态向下走n步的情况。



2、n步回报

如果我们考虑如下的n取值下的回报 ($n = 1, 2, \dots, \infty$)

$$\begin{aligned} n = 1 \quad (TD) \quad G_t^{(1)} &= R_{t+1} + \gamma V(S_{t+1}) \\ n = 2 \quad (TD) \quad G_t^{(2)} &= R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2}) \\ &\vdots \\ n = \infty \quad (MC) \quad G_t^{(\infty)} &= R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T \end{aligned}$$

那么我们可以进行泛化定义n步回报为：

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

根据n步回报修改 $TD(0)$ 的更新方法为：

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t^{(n)} - V(S_t))$$

这样我们就可以得到如下的n步时序差分算法。

❗ n 步TD(0)估计 $V \approx v_\pi$

输入：策略 π

算法参数：步长 $\alpha \in (0, 1]$ ，正整数 n

对 $s \in \mathcal{S}$ ，任意初始化 $V(s)$

所有存储和访问操作（对于 S_t 和 R_t ）都可以使其索引 $\text{mod } n + 1$

对每个回合循环：

 初始化并存储 $S_0 \neq$ 终点

$T \leftarrow \infty$

 对 $t = 0, 1, 2, \dots$ 循环：

 如果 $t < T$ ，则：

 根据 $\pi(\cdot | S_t)$ 采取行动

 观察并将下一个奖励存储为 R_{t+1} ，将下一个状态存储为 S_{t+1}

 如果 S_{t+1} 是终点，则 $T \leftarrow t + 1$

$\tau \leftarrow t - n + 1$ （ τ 是状态估计正在更新的时间）

 如果 $\tau \geq 0$ ：

$$G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$$

 如果 $\tau + n < T$ ，则 $G \leftarrow G + \gamma^n V(S_{\tau+n})$

$$V(S_\tau) \leftarrow V(S_\tau) + \alpha [G - V(S_\tau)] \quad (G_{\tau:\tau+n})$$

 直到 $\tau = T - 1$

3、n步时序差分方法在随机游走上的应用

在《强化学习基础篇（十九）TD与MC在随机游走问题应用》我们实现了随机游走的问题。这里我们将原问题的6个状态调整为19个状态，下面看看通过n步回报的方法效果如何。

导入库函数定义超参数

```
1 import numpy as np
2 import matplotlib
3 matplotlib.use('Agg')
4 import matplotlib.pyplot as plt
5 from tqdm import tqdm
```

```

6
7 # 共19个状态
8 N_STATES = 19
9
10 # 定义折扣因子
11 GAMMA = 1
12
13 # 定义状态空间
14 STATES = np.arange(1, N_STATES + 1)
15
16 # 起始状态为第10个状态
17 START_STATE = 10
18
19 # 一共两个terminal状态
20 # 左边界的状态的回报为-1，右边界的状态的回报为+1
21 END_STATES = [0, N_STATES + 1]
22
23 # 设定真实价值 (true value)
24 TRUE_VALUE = np.arange(-20, 22, 2) / 20.0
25 TRUE_VALUE[0] = TRUE_VALUE[-1] = 0

```

n-steps TD算法实现

```

1 # n-steps TD 算法实现
2 # value: 输入状态价值函数
3 # n: 输入n步的值
4 # alpha: 定义步长
5 def temporal_difference(value, n, alpha):
6     # 初始化状态位置
7     state = START_STATE
8
9     # 定义一个列表存储states和rewards
10    states = [state]
11    rewards = [0]
12
13    # 进行时间跟踪
14    time = 0
15
16    # 是定时间初始为无限
17    T = float('inf')
18    while True:
19        # 进一个时间步
20        time += 1
21
22        if time < T:
23            # 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
24            if np.random.binomial(1, 0.5) == 1:
25                next_state = state + 1
26            else:
27                next_state = state - 1
28            # 按照问题定义，处理计算奖励。
29            if next_state == 0:
30                reward = -1
31            elif next_state == 20:
32                reward = 1
33            else:
34                reward = 0

```

```

35
36         # 存储下一步状态与奖励
37         states.append(next_state)
38         rewards.append(reward)
39
40         if next_state in END_STATES:
41             T = time
42
43         # get the time of the state to update
44         update_time = time - n
45         if update_time >= 0:
46             returns = 0.0
47             # 计算n步奖励
48             for t in range(update_time + 1, min(T, update_time + n) + 1):
49                 returns += pow(GAMMA, t - update_time - 1) * rewards[t]
50             # 将n步奖励增加到总回报中
51             if update_time + n <= T:
52                 returns += pow(GAMMA, n) * value[states[(update_time + n)]]
53             state_to_update = states[update_time]
54             # 更新状态值函数
55             if not state_to_update in END_STATES:
56                 value[state_to_update] += alpha * (returns -
value[state_to_update])
57             if update_time == T - 1:
58                 break
59             state = next_state

```

实验运行与绘制结果

```

1  def figure7_2():
2      # 这里要比较的n步包含了1,2,4,8..512
3      steps = np.power(2, np.arange(0, 10))
4
5      # 这里比较了三个步长
6      alphas = np.arange(0, 1.1, 0.1)
7
8      # 每次运行10个episodes
9      episodes = 10
10
11     # 实验总次数（因为结果要对这些100次取平均）
12     runs = 100
13
14     # track the errors for each (step, alpha) combination
15     errors = np.zeros((len(steps), len(alphas)))
16     for run in tqdm(range(0, runs)):
17         for step_ind, step in enumerate(steps):
18             for alpha_ind, alpha in enumerate(alphas):
19                 # print('run:', run, 'step:', step, 'alpha:', alpha)
20                 value = np.zeros(N_STATES + 2)
21                 for ep in range(0, episodes):
22                     temporal_difference(value, step, alpha)
23                     # 计算均方根误差 (RMS error)
24                     errors[step_ind, alpha_ind] +=
np.sqrt(np.sum(np.power(value - TRUE_VALUE, 2)) / N_STATES)
25                 # 对结果取平均
26                 errors /= episodes * runs
27

```

测试结果