

强化学习作业2代码题

18300290007 加兴华

Q3

3. (10 pts) Reproduce the figure on pg. 21 of Lecture 4 for the test of RM on root finding.

由于lecture 4第21页没有图，我复现了lecture 3第21页。

RM算法代码实现：

```
1  import numpy as np
2  import random
3  def RM(x0, h, noise_SD=0, alpha_type=1, steps=100):
4      '''
5      [Inputs]:
6      x0                初始条件
7      h                x=h(x) 一式中的h函数
8      noise_SD         高斯噪声标准差
9      steps            迭代次数
10     alpha_k          迭代步长类型 (1:ak=1; 2:ak=1/k; 3:ak=1/k^2)
11     [Outputs]:
12     X                记录每次迭代x的列表
13     '''
14     X=[x0]
15     alpha=[0]
16     if alpha_type==1:
17         alpha+=[1 for i in range(1, steps+1)]
18     elif alpha_type==2:
19         alpha+=[1/i for i in range(1, steps+1)]
20     elif alpha_type==3:
21         alpha+=[1/(i**2) for i in range(1, steps+1)]
22     for i in range(1, steps+1):
23         epsilon=np.random.randn()*noise_SD
24         X.append(X[-1] + alpha[i] * (h(X[-1]) - X[-1]+epsilon))
25     return X
```

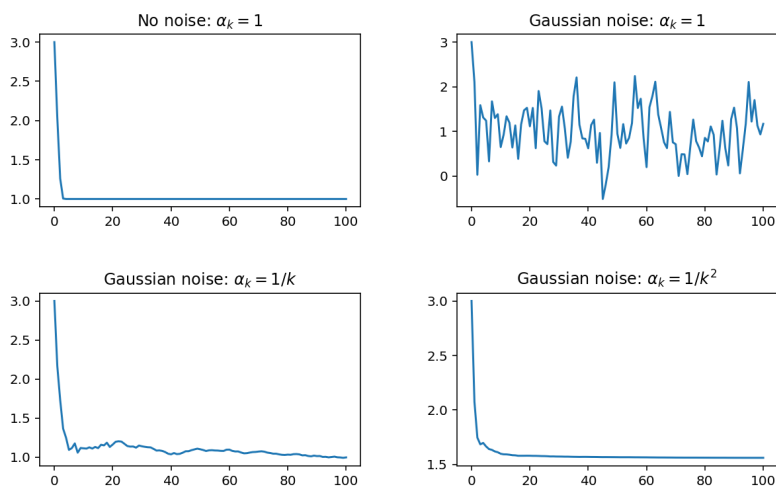
题目：寻找 $f(x) = \tanh(x) - 1$ 的根，初值条件 $x_1 = 3$

```
1  ## Cell内打印图片
2  %matplotlib inline
3  ## 设置显示标准维retina
4  %config InlineBackend.figure_format = 'retina'
5  import matplotlib.pyplot as plt
6  import numpy as np
```

```

7 plt.figure(figsize=(10,6))
8
9
10 ## 设置模型
11 def h(x):
12     return x - np.tanh(x-1)
13 x0=3
14
15 ## 调用RM算法
16 n = np.arange(101)
17 x_case1 = RM(x0, h)
18 x_case2 = RM(x0, h, 0.5)
19 x_case3 = RM(x0, h, 0.5, 2)
20 x_case4 = RM(x0, h, 0.5, 3)
21
22
23 plt.subplots_adjust(left=None, bottom=None, right=None, top=None, wspace=0.3, hspace=0.5)
24 plt.subplot(2, 2, 1)
25 plt.plot(n, x_case1)
26 plt.title('No noise:  $\alpha_k = 1$ ')
27
28 plt.subplot(2, 2, 2)
29 plt.plot(n, x_case2)
30 plt.title('Gaussian noise:  $\alpha_k = 1$ ')
31
32 plt.subplot(2, 2, 3)
33 plt.plot(n, x_case3)
34 plt.title('Gaussian noise:  $\alpha_k = 1/k$ ')
35
36 plt.subplot(2, 2, 4)
37 plt.plot(n, x_case4)
38 plt.title('Gaussian noise:  $\alpha_k = 1/k^2$  ');

```



多次运行,观察发现前3张图与课件较为一致,但最后一张差异较大,我自己每次运行差异也较大,可能是步长衰减过快和高斯噪声波动性过大这两个因素组合导致的。

Q5

5. (10 pts) Test TD(0) for policy evaluation (pg. 5 in Lecture 5) by applying the algorithm on the following single episode **repeatedly**:

$$A, a_0, r = 0, B, a_1, r = 0, A, a_0, r = 0, B, a_2, r = 1, T,$$

where T is the terminal state (that is, $v(T) = 0$). Try different values for the discount factor γ and 1) report the state values $v(A)$ and $v(B)$ that the algorithm converges to; 2) try to find some pattern from your results and guess what kind of solution the algorithm converges to. [Assume $v(A) = v(B) = 0$ for the initial guess and you can use a small constant stepsize in the algorithm or use the reciprocal of the visit times as the stepsize]

Algorithm 1: TD(0) Policy Evaluation

Initialization: $v_0(s) = 0 \forall s \in \mathcal{S}$, target policy π and initial state s_0 .

for $t = 0, 1, 2, \dots$ **do**

 Sample a tuple $(s_t, a_t, r_t, s_{t+1}) \sim \pi$ from s_t
 $v_{t+1}(s_t) = v_t(s_t) + \alpha_t(s_t)([r_t + \gamma v_t(s_{t+1})] - v_t(s_t))$

end

TD(0)评估阶段代码实现 (为适配本题, 我将采样过程从TD0算法中分离出来)

```
1  import random
2  def TD0_eval(samples, times=5, gamma=0.9):
3      '''
4      [Inputs]:
5      samples          采样的轨迹列表
6      times            重复TD0次数
7      gamma            衰减系数
8      samples
9      [Outputs]:
10     states           状态名列表
11     values           状态值列表
12     '''
13     ## 统计出现过的状态名
14     ### 轨迹形如sarsarsars... 因此每隔3个提取一次
15     states=[]
16     for sample in samples:
17         states+= sample[::3]
18     states=sorted(list(set(states)))    ### 去除重复
19     ## 生成状态的相关列表
```

```

20     values=[0 for i in states] ### 状态值
21     freqs=[0 for i in states]    ### 每个状态已经更新过的次数，用于更新步长
22     ## 从轨迹列表中随机采样一条进行TD0更新
23     for t in range(times):
24         sample=random.choice(samples)
25         for i in range(len(sample)):
26             if sample[i] not in states: continue    ### 跳过非状态元素
27             index=states.index(sample[i])
28             freqs[index]+=1
29             if i==len(sample)-1: continue    ### 跳过轨迹最后一个状态
30             index_next_state=states.index(sample[i+3])
31             alpha=1/freqs[index]
32             values[index]+=alpha*(sample[i+2]+gamma*values[index_next_state]-
values[index])
33     return states, values

```

1)

```

1  import pandas as pd
2  import numpy as np
3  times=50000
4  samples=[['A', 'a0', 0, 'B', 'a1', 0, 'A', 'a0', 0, 'B', 'a2', 1, 'T']]
5  states, values=TD0_eval(samples, times, 0.1)
6  df = pd.DataFrame([values+[0.1, times]], columns = states+['衰减率', '重复次数'])
7  for gamma in np.linspace(0.2, 0.9, 8):
8      states, values=TD0_eval(samples, times, gamma)
9      df.loc[len(df)]=values+[gamma, times]
10 df.style.hide(axis='index').format({'衰减率': '{:.1f}', '重复次数': '{:.0f}'})

```

A	B	T	衰减率	重复次数
0.050247	0.502511	0.000000	0.1	50000
0.102028	0.510197	0.000000	0.2	50000
0.157033	0.523538	0.000000	0.3	50000
0.217298	0.543414	0.000000	0.4	50000
0.285467	0.571256	0.000000	0.5	50000
0.365199	0.609295	0.000000	0.6	50000
0.461823	0.661014	0.000000	0.7	50000
0.583471	0.731927	0.000000	0.8	50000
0.743060	0.830973	0.000000	0.9	50000

2)

在不同的 γ 下, $V(A) \approx \gamma * V(B)$;

可以猜测收敛状态值有如下形式: $V(B) = f(\gamma), V(A) = f(\gamma) * \gamma$

Q6

6. (10 pts) Implement MC Learning with ϵ -Greedy Exploration (pg. 23 of Lecture 4) and Off-policy MC Learning (pg. 30 of Lecture 4) and test them on the 10 gridworld problem shown in Figure 1.

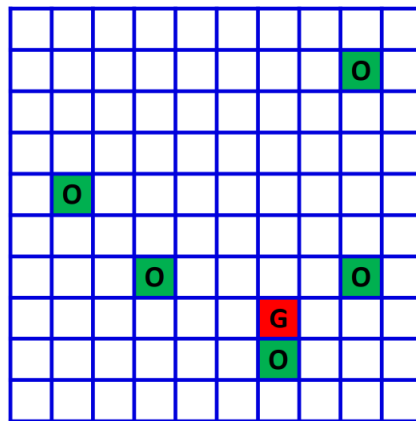


Figure 1: Hit obstacle grid:-10; reach goal state (from other states): 10. Goal state is the terminal state, that is, if the agent leaves the goal state no matter what action it takes, it will return to the goal state with reward 0. The other settings are the same as the one in Homework I.

```
1  ## 模型参数
2  ### action坐标分解及名称列表
3  dx = [-1, 0, 1, 0, 0]
4  dy = [0, -1, 0, 1, 0]
5  acts = ['up', 'left', 'down', 'right', 'stay']
6  arrow= ['↑', '←', '↓', '→', '-']
7
8  ### reward函数
9  H, W = (10, 10)
10 reward = np.zeros([H, W])
11 obstacles = [(4, 1), (6, 3), (1, 8), (6, 8), (8, 6)]
12 target = [(7, 6)]
13 for x, y in obstacles:
14     reward[x, y] = -10
15 for x, y in target:
16     reward[x, y] = 10
17
18 np.random.seed(10086)
```

ϵ - greedy 版本代码实现**Algorithm 6:** MC Learning with ϵ -Greedy Exploration**Initialization:** $N(s, a) = 0, Q(s, a) = 0, \forall s, a, \pi_0$ **for** $k = 0, 1, 2, \dots$ **do** Initialize s_0 and sample an episode following π_k :

$$(s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T) \sim \pi_k$$

 $G \leftarrow 0$ **for** $t = T-1, T-2, \dots, 0$ **do** $G \leftarrow \gamma G + r_t$ **if** (s_t, a_t) does not appear in $(s_0, a_0, s_1, a_1, \dots, s_{t-1}, a_{t-1})$ **then** $N(s_t, a_t) \leftarrow N(s_t, a_t) + 1$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \frac{1}{N(s_t, a_t)} (G - Q(s_t, a_t))$$

 Update policy of visited state via ϵ_k -greedy:

$$\pi_{k+1}(a|s_t) = \begin{cases} 1 - \epsilon_k + \frac{\epsilon_k}{|\mathcal{A}|} & \text{if } a = \underset{a'}{\operatorname{argmax}} Q(s_t, a'), \\ \frac{\epsilon_k}{|\mathcal{A}|} & \text{otherwise.} \end{cases}$$

end **end****end**

个人优化:

假设轨迹长度=1000, $\gamma = 0.9$, 终点回报对于轨迹前面的 $Q(s,a)$ 影响非常小, 反而更可能被障碍的惩罚影响, 因此不如采样更多次并舍弃特别长的轨迹.

```

1  import numpy as np
2  import random
3  import tqdm
4  def MC_egreedy(times=1000, gamma=0.9, threshold=500):
5      '''
6      [Inputs]:
7      times          重复次数
8      gamma          衰减率
9      threshold      轨迹长度上限, 超过就放弃本次采样
10     [Outputs]:
11     pi              策略矩阵
12     Q               动作值3维矩阵
13     '''
14     ## 初始化相关矩阵
15     Q=np.zeros((len(acts), H, W))    ### |A|个块, 每块H行W列
16     N=np.zeros((len(acts), H, W))
17     pi=np.random.randint(0, 5, (H, W))    ### 泛化前的策略矩阵
18     pi[np.array(target)[: , 0], np.array(target)[: , 1]]=4    ### 终点不动
19
20     for k in tqdm.tqdm(range(times)):

```

```

21     ## 采样轨道
22     epsilon=0.3+0.2*np.cos(k/times*np.pi)
23     s0=(random.randint(0, H-1), random.randint(0, W-1))
24     episode=[s0]
25     while episode[-1] not in target:
26         if len(episode)>threshold: break
27         s_now=episode[-1]
28         ### 使用egreedy作为pi的beta变形
29         if random.uniform(0, 1) < epsilon:
30             a=random.randint(0, 4)
31         else:
32             a=np.int8(pi[s_now])
33         episode.append(a)
34         s_next=(s_now[0]+dx[a], s_now[1]+dy[a])
35         if not (0<=s_next[0]<=H-1 and 0<=s_next[1]<=W-1):
36             s_next=s_now
37             r=0
38         else:
39             r=reward[s_next]
40             episode.append(r)
41             episode.append(s_next)
42
43     # print(episode, '\n')
44
45     ## 更新
46     if len(episode)>threshold: continue
47     if len(episode)==1: continue
48     state=episode[-4:0:-3]
49     r=episode[-2:2:-3]
50     a=episode[-3:1:-3]
51     G = 0
52     for t in range(0, len(state)):
53         G=gamma*G+r[t]
54         x,y=state[t][0], state[t][1]
55         if (state[t], a[t]) in zip(state[t+1:], a[t+1:]): continue
56         N[a[t], x, y]+=1
57         Q[a[t], x, y]+=(G-Q[a[t], x, y])/N[a[t], x, y]
58         if Q[a[t], x, y] > Q[int(pi[x, y]), x, y]:
59             pi[x, y] = a[t]
60     return pi, Q
61 pi, _=MC_egreedy(100000, 0.9, 10000)
62 for i in range(H):
63     for j in range(W):
64         print(arrow[np.int8(pi[i, j])], end=' ')
65     print('')

```

```

1  100%|████████████████████| 100000/100000 [00:06<00:00, 14557.15it/s]
2
3  → ↓ → ↓ ↓ ↓ ↓ ← ← ↓
4  → → → ↓ ↓ ↓ ↓ ← ↓ ↓
5  → → → → ↓ ↓ ↓ ↓ ↓ ↓
6  ↑ → → → → ↓ ↓ ↓ ← ↓
7  ↓ ↓ → → → ↓ ↓ ↓ ← ←
8  ↓ ↓ → → → ↓ ↓ ← ← ↑
9  ↓ ↓ ↓ ↓ → → ↓ ← ← ↓
10 → → → → → → - ← ← ←
11 → → → → ↑ ↑ ↑ ↑ ← ←
12 → → → → ↑ ↑ ← ↑ ← ←

```

Off-policy 版本代码实现

Algorithm 7: Off-policy MC Learning

Initialization: $\forall s, a$, arbitrarily initial $Q(s, a)$, $\pi_0(s) = \operatorname{argmax}_a Q(s, a)$,
 $N(s, a) = 0$.

```

for  $k = 0, 1, 2, \dots$  do
     $b_k \leftarrow$  any soft policy, i.e.,  $b_k(a|s) > 0, \forall s, a$ 
    Initialize  $s_0$  and sample an episode following  $b_k$ :

         $(s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T) \sim b_k$ 

     $G \leftarrow 0, W \leftarrow 1$ 
    for  $t = T - 1, T - 2, \dots, 0$  do
         $G \leftarrow r_t + \gamma G$ 
        if  $(s_t, a_t)$  does not appear in  $(s_0, a_0, s_1, a_1, \dots, s_{t-1}, a_{t-1})$  then
             $N(s_t, a_t) \leftarrow N(s_t, a_t) + 1$ 
             $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \frac{1}{N(s_t, a_t)} (W \cdot G - Q(s_t, a_t))$ 
             $\pi_{k+1}(s_t) \leftarrow \operatorname{argmax}_a Q(s_t, a)$ 
        end
         $W \leftarrow W \frac{\pi_k(a_t|s_t)}{b_k(a_t|s_t)}$ 
    end
end

```

```

1  import numpy as np
2  import random
3  import tqdm
4  def MC_offpolicy(times=1000, gamma=0.9, threshold=500):
5      '''
6      [Inputs]:
7      times          重复次数
8      gamma          衰减率

```



```

9     threshold          轨迹长度上限，超过就放弃本次采样
10    [Outputs]:
11    pi                  策略矩阵
12    Q                   动作值3维矩阵
13
14    '''
15    ## 初始化相关矩阵
16    Q=np.zeros((len(acts),H,W))    ### |A|个块，每块H行W列
17    N=np.zeros((len(acts),H,W))
18    pi=np.random.randint(0,5,(H,W))    ### 泛化前的策略矩阵
19    pi[np.array(target)[: ,0],np.array(target)[: ,1]]=4    ### 终点不动
20
21    for k in tqdm.tqdm(range(times)):
22        ## 采样轨道
23        epsilon=0.3+0.2*np.cos(k/times*np.pi)
24        s0=(random.randint(0, H-1), random.randint(0, W-1))
25        episode=[s0]
26        while episode[-1] not in target:
27            if len(episode)>threshold: break
28            s_now=episode[-1]
29            if random.uniform(0, 1) < epsilon:    ### explore
30                a=random.randint(0, 4)
31            else:                                ### exploit
32                a=np.int8(pi[s_now])
33            episode.append(a)
34            s_next=(s_now[0]+dx[a], s_now[1]+dy[a])
35            if not (0<=s_next[0]<=H-1 and 0<s_next[1]<=W-1):
36                s_next=s_now
37                r=0
38            else:
39                r=reward[s_next]
40            episode.append(r)
41            episode.append(s_next)
42
43        # print(episode, '\n')
44
45        ## 更新
46        if len(episode)>threshold: continue
47        if len(episode)==1: continue
48        state=episode[-4:0:-3]
49        r=episode[-2:2:-3]
50        a=episode[-3:1:-3]
51        G = 0
52        w = 1
53        for t in range(0, len(state)):

```

```

54         Echoes=pi[state[t]]
55         G=gamma*G+r[t]
56         x,y=state[t][0],state[t][1]
57         if (state[t],a[t]) not in zip(state[t+1:],a[t+1:]):
58             N[a[t],x,y]+=1
59             Q[a[t],x,y]+=(w*G-Q[a[t],x,y])/N[a[t],x,y]
60             if Q[a[t],x,y] > Q[int(pi[x,y]),x,y]:
61                 pi[x,y] = a[t]
62         if a[t] == Echoes:
63             w /= 1 - epsilon * 0.8      ### 概率比=1-epsilon(1-1/|A|)
64         else:
65             w *= 0
66     return pi,Q
67 pi,_=MC_offpolicy(100000,0.9,10000)
68 for i in range(H):
69     for j in range(W):
70         print(arrow[np.int8(pi[i,j])], end=' ')
71     print('')

```

```

1  100%|████████████████████| 100000/100000 [00:10<00:00, 9891.29it/s]
2
3  → ↓ ↓ → ↓ → → ↓ ← ↓
4  → → ↓ ↓ ↓ → ↓ ↓ ↓ ↓
5  → → ↓ → → → ↓ ← ↓ ↓
6  → → → ↓ → → ↓ ← ← ←
7  ← ↓ ↓ → ↓ ↓ ↓ ← ← ←
8  → → ↓ → → ↓ ↓ ↓ ← ↓
9  → ↑ ↓ ↓ ↓ → ↓ ← ← ↓
10 → → → → → → - ← ← ←
11 → ↑ ↑ → → ↑ ↑ ↑ ↑ ↑
12 → ↑ ↑ → ↑ ↑ ← ↑ ← ↑

```

设置重复采样100000次、衰减率0.9、舍弃长度大于10000的轨迹，两种算法的结果如输出所示，可以发现MC算法受限于采样随机性，选择的行动可能不是最优的，但是一定不会选择最差的，表现为动作总是绕过障碍。