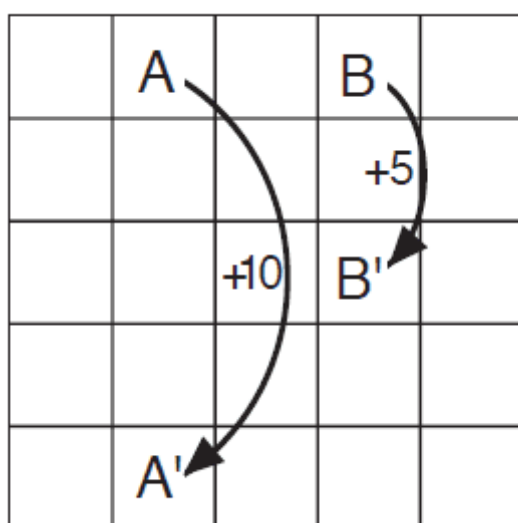


强化学习基础篇（二十一） 网格问题

该问题基于《Reinforcement Learning: An Introduction》在第三章中给出了一个简单的例子: Gridworld, 以帮助我们理解有限MDPs。

1、网格问题描述

如下图所示的长方形网格代表的是一个简单的有限MDP。网格中的格子代表的是环境中的状态。在每个格子中，有四个可选的动作：东、南、西和北，每个动作都会使智能体在对应的方向上前进一步。如果采取的动作使得其不在网格内，则智能体会在原位置不移动，并且得到一个值为-1的收益。除了将智能体从特定的状态A和B中移走的动作外，其他动作的收益都是0。在状态A中，四种动作都会得到一个+10的收益，并且把智能体移动到A'。在状态B中，四种动作都会得到一个+5的收益，并且把智能体移动到B'。



2、解析

这是一个典型的MDP问题，有限的状态空间，有限的动作空间。但是，我们面临两个棘手的问题：

- 根据Bellman等式，当前状态的价值依赖于下一个状态的价值，这是一个递归的过程。从何处开始，在何处结束？对于每一个状态都有不同的episode。
- 状态空间和动作空间如何表达？状态空间的表达容易想到使用格子的坐标来表示，比如格子A的状态可表示为坐标0,10,1，那么动作空间如何表达比较合适呢？

2.1、Bellman等式的迭代求解

根据bellman求解MDP问题一般有两种思路：第一，如果MDP的状态转移矩阵是已知的，则可以精确的数值求解：求解线性方程组而已。。第二，Bellman等式本身是一个迭代公式，因此可以通过迭代的方式逐步逼近状态价值函数。这种方式对状态转移矩阵不敏感，即无论是否知道状态转移矩阵都可以使用，也非常适合于编程实现。

迭代求解状态价值函数的基本步骤为：

1. 初始化状态价值矩阵，一般初始化为0。这是计算的起点，此时每个状态的价值函数都是0。
2. 根据Bellman等式，计算每个状态的新的价值函数。注意这一步**不是一个递归的过程**，因为上一轮的状态价值矩阵是已知的，在这一轮直接可以根据上一轮的状态价值矩阵计算新的价值函数即可，即：

$$v_{k+1}(s) = \sum_a \pi(a | s) \sum_{r, s'} p(s', r | s, a) [r + \gamma v_k(s')]$$

其中的 k 代表了迭代的轮次（episode），注意不要和每个轮次（episode）中的时刻（步骤）的 t 混淆了。

计算每个轮次迭代的误差，达到设定的误差范围即可停止迭代，此时即获得了最终的状态价值函数。

这种迭代求解状态价值函数的方法称为“Policy Evaluation Prediction”，即对于任意给定的策略 π ，都可以通过迭代的方式逐步逼近求解状态价值函数。但是，必须确保这种迭代求解的方法对于任意策略 π 都是收敛的，

证明如下：

通过观察上式，并对照Bellman等式

$$v_{\pi}(s) = \sum_a \pi(a | s) \sum_{r, s'} p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

可以看出，**式1只是对式2的简单变量替换**，由 v_{π} 替换为 v_k ，而式2在 $\gamma < 1$ 或者episode能够在有限步内结束时即收敛，因此式1在同样的条件下也会收敛，即迭代求解状态价值函数的方法的收敛条件为：或者 $\gamma < 1$ ，或者所解决的MDP问题的episode是有限的。这两个条件几乎总是能够满足其中之一的。

2.2、动作空间的表示方式

动作对状态的影响表现为状态坐标的改变，因此将动作定义为对状态坐标的偏移量是个合理的方案，这样可以直接将状态和动作相加即可获得“下一个状态”。以动作up为例，其对状态的影响为横坐标不变，纵坐标增加+1，因此可以将up定义为[0,1]。同理，动作down可以定义为[0,-1]。

3. 实现过程

3.1、导入库文件

```
1 import matplotlib
2 import matplotlib.pyplot as plt
3 import numpy as np
4 from matplotlib.table import Table
```

3.2、定义环境

```
1 WORLD_SIZE = 5 # 定义问题中格子的数量
2 # 定义问题中A,A',B,B'的位置
3 A_POS = [0, 1]
4 A_PRIME_POS = [4, 1]
5 B_POS = [0, 3]
6 B_PRIME_POS = [2, 3]
7 # 定义折扣参数
8 DISCOUNT = 0.9
9
10 # 把动作定义为对x, y坐标的增减改变
11 ACTIONS = [np.array([0, -1]), # 向上
12            np.array([-1, 0]), # 向左
13            np.array([0, 1]), # 向下
14            np.array([1, 0])] # 向右
15 # 定义画图会用到的动作
16 ACTIONS_FIGS=[ '←', '↑', '→', '↓']
17 # 该问题中每个动作选择的概率为0.25
```

3.3、定义环境转移过程

```

1  def step(state, action):
2      """每次走一步
3      state:当前状态, 坐标的list, 比如[1,1]
4      action:当前采取的动作, 是对状态坐标的修正
5      函数返回值, 下一个状态(坐标的list)和reward
6      """
7      # 当在A位置的时候, 会转移到A', 并得到奖励+10.
8      if state == A_POS:
9          return A_PRIME_POS, 10
10     # 当在B位置的时候, 会转移到B', 并得到奖励+5.
11     if state == B_POS:
12         return B_PRIME_POS, 5
13
14     # 计算下一个状态
15     next_state = (np.array(state) + action).tolist()
16     x, y = next_state
17     # 当运动未知超出格子世界则在原位置不变, 并且得到-1的收益, 其他情况得到的收益都是0.
18     if x < 0 or x >= WORLD_SIZE or y < 0 or y >= WORLD_SIZE:
19         reward = -1.0
20         next_state = state
21     else:
22         reward = 0
23     return next_state, reward

```

3.4、绘制价值函数结果

```

1  def draw_image(image):
2      fig, ax = plt.subplots()
3      ax.set_axis_off()
4      tb = Table(ax, bbox=[0, 0, 1, 1])
5
6      nrows, ncols = image.shape
7      width, height = 1.0 / ncols, 1.0 / nrows
8
9      # Add cells
10     for (i, j), val in np.ndenumerate(image):
11
12         # add state labels
13         if [i, j] == A_POS:
14             val = str(val) + " (A)"
15         if [i, j] == A_PRIME_POS:
16             val = str(val) + " (A')"
17         if [i, j] == B_POS:
18             val = str(val) + " (B)"
19         if [i, j] == B_PRIME_POS:
20             val = str(val) + " (B')"
21
22         tb.add_cell(i, j, width, height, text=val,
23                     loc='center', facecolor='white')
24
25     # Row and column labels...
26     for i in range(len(image)):

```

```

27         tb.add_cell(i, -1, width, height, text=i+1, loc='right',
28                     edgecolor='none', facecolor='none')
29         tb.add_cell(-1, i, width, height/2, text=i+1, loc='center',
30                     edgecolor='none', facecolor='none')
31
32     ax.add_table(tb)

```

3.5、绘制策略图形

```

1  def draw_policy(optimal_values):
2      fig, ax = plt.subplots()
3      ax.set_axis_off()
4      tb = Table(ax, bbox=[0, 0, 1, 1])
5
6      nrows, ncols = optimal_values.shape
7      width, height = 1.0 / ncols, 1.0 / nrows
8
9      # Add cells
10     for (i, j), val in np.ndenumerate(optimal_values):
11         next_vals=[]
12         for action in ACTIONS:
13             next_state, _ = step([i, j], action)
14             next_vals.append(optimal_values[next_state[0],next_state[1]])
15
16         best_actions=np.where(next_vals == np.max(next_vals))[0]
17         val=''
18         for ba in best_actions:
19             val+=ACTIONS_FIGS[ba]
20
21         # add state labels
22         if [i, j] == A_POS:
23             val = str(val) + " (A)"
24         if [i, j] == A_PRIME_POS:
25             val = str(val) + " (A')"
26         if [i, j] == B_POS:
27             val = str(val) + " (B)"
28         if [i, j] == B_PRIME_POS:
29             val = str(val) + " (B')"
30
31         tb.add_cell(i, j, width, height, text=val,
32                     loc='center', facecolor='white')
33
34     # Row and column labels...
35     for i in range(len(optimal_values)):
36         tb.add_cell(i, -1, width, height, text=i+1, loc='right',
37                     edgecolor='none', facecolor='none')
38         tb.add_cell(-1, i, width, height/2, text=i+1, loc='center',
39                     edgecolor='none', facecolor='none')
40
41     ax.add_table(tb)

```

3.6、计算等概率随机策略下的状态价值函数(解贝尔曼期望方程)

贝尔曼期望方程的形式为：

$$v_{\pi}(s) = \sum_a \pi(a | s) \sum_{r, s'} p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

```
1 def figure_3_2():
2     """
3     计算每个单元格的状态价值函数
4     """
5     # 状态价值函数的初值设置为0
6     value = np.zeros((WORLD_SIZE, WORLD_SIZE))
7     while True:
8         # 每一轮迭代都会产生一个new_value, 直到new_value和value很接近即收敛为止
9         new_value = np.zeros_like(value)
10
11        # 对格子世界的每个状态进行遍历
12        for i in range(WORLD_SIZE):
13            for j in range(WORLD_SIZE):
14                # 对每个状态的4个动作进行遍历
15                for action in ACTIONS:
16                    # 计算当前动作执行后的下一个状态, 以及奖励
17                    (next_i, next_j), reward = step([i, j], action)
18                    # 由于每个方向只有一个reward和s'的组合, 这里的p(s', r|s, a)=1
19                    new_value[i, j] += ACTION_PROB * (reward + DISCOUNT *
value[next_i, next_j])
20                # 当前new_value和value很接近, 即error小到一定的程度则停止。
21                error = np.sum(np.abs(new_value - value))
22                if error < 1e-4:
23                    draw_image(np.round(new_value, decimals=2))
24                    plt.show()
25                    break
26                value = new_value
27    figure_3_2()
```

输出的测试结果如下所示：

	1	2	3	4	5
1	3.31	8.79 (A)	4.43	5.32 (B)	1.49
2	1.52	2.99	2.25	1.91	0.55
3	0.05	0.74	0.67	0.36 (B')	-0.4
4	-0.97	-0.44	-0.35	-0.59	-1.18
5	-1.86	-1.35 (A')	-1.23	-1.42	-1.98

可以看出，靠近下边的格子，尤其是靠近边界的格子的状态价值偏小，这是因为靠近边界的格子很容易出界，而出界的奖励是-1；格子A的价值最高，这是很容易理解的，因为从A到A'的奖励是+10。但是为什么格子A的价值<10呢？这是因为从A到A'的奖励尽管是+10，但是A'的价值却是负值，导致A的价值会跟着损失一些。格子B的价值分析方法类似。

3.7、直接求解贝尔曼方程

下面采用代码直接去求解贝尔曼方程。

```
1 def figure_3_2_linear_system():
2     '''
3     该函数通过直接求解线性方程组，得到贝尔曼方程的精确解。
4     Here we solve the linear system of equations to find the exact solution.
5     We do this by filling the coefficients for each of the states with their
6     respective right side constant.
7     '''
8     # 定义A矩阵，25*25大小，对应着25个线性方程组
9     A = -1 * np.eye(WORLD_SIZE * WORLD_SIZE)
10    b = np.zeros(WORLD_SIZE * WORLD_SIZE)
11    for i in range(WORLD_SIZE):
12        for j in range(WORLD_SIZE):
13            # 当前状态
14            s = [i, j]
15            # 该函数找到状态s的序号，即将格子展开为5*5=25后，在这个长序列中的编号。
16            index_s = np.ravel_multi_index(s, (WORLD_SIZE, WORLD_SIZE))
17            for a in ACTIONS:
18                # 获取执行动作后的下一个状态和奖励
19                s_, r = step(s, a)
20                # 找到下一个状态的序号
21                index_s_ = np.ravel_multi_index(s_, (WORLD_SIZE,
22                WORLD_SIZE))
23                # 线性方程组Ax=b的A中对应元素按照动作概率*折扣因子进行更新
24                A[index_s, index_s_] += ACTION_PROB * DISCOUNT
25                b[index_s] -= ACTION_PROB * r
26            # 求解Ax=b
27            x = np.linalg.solve(A, b)
28            draw_image(np.round(x.reshape(WORLD_SIZE, WORLD_SIZE), decimals=2))
29            plt.show()
```

输出过和通过迭代方式求解出来的结果是一致的。

	1	2	3	4	5
1	3.31	8.79 (A)	4.43	5.32 (B)	1.49
2	1.52	2.99	2.25	1.91	0.55
3	0.05	0.74	0.67	0.36 (B')	-0.4
4	-0.97	-0.44	-0.35	-0.59	-1.18
5	-1.86	-1.35 (A')	-1.23	-1.42	-1.98

3.8、求解最优价值函数（解贝尔曼最优方程）

贝尔曼最优方程的形式为：

$$\begin{aligned} q_*(s, a) &= \mathbb{E}[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r \mid s, a) [R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a')] \end{aligned}$$

```
1 def figure_3_5():
2     '''
3     求解贝尔曼最优方程
4     '''
5     # 初始值设为0
6     value = np.zeros((WORLD_SIZE, WORLD_SIZE))
7     while True:
8         # keep iteration until convergence
9         new_value = np.zeros_like(value)
10        # 遍历所有状态
11        for i in range(WORLD_SIZE):
12            for j in range(WORLD_SIZE):
13                values = []
14                # 遍历所有动作
15                for action in ACTIONS:
16                    # 执行动作, 转移到后继状态, 并获得立即奖励
17                    (next_i, next_j), reward = step([i, j], action)
18                    # value iteration
19                    # 缓存动作价值函数  $q(s, a) = r + \gamma v(s')$ 
20                    values.append(reward + DISCOUNT * value[next_i, next_j])
21                    # 根据贝尔曼最优方程, 找出最大的动作价值函数  $q(s, a)$  进行更新
22                    new_value[i, j] = np.max(values)
23        # 迭代终止条件: 误差小于  $1e-4$ 
24        if np.sum(np.abs(new_value - value)) < 1e-4:
25            draw_image(np.round(new_value, decimals=2))
26            plt.show()
27            plt.close()
28            draw_policy(new_value)
29            plt.show()
30            break
31        value = new_value
32    figure_3_5()
```

计算得到的最优价值函数如下所示：

	1	2	3	4	5
1	21.98	24.42 (A)	21.98	19.42 (B)	17.48
2	19.78	21.98	19.78	17.8	16.02
3	17.8	19.78	17.8	16.02 (B')	14.42
4	16.02	17.8	16.02	14.42	12.98
5	14.42	16.02 (A')	14.42	12.98	11.68

从最后价值函数得到的最优策略如下，图中有很多个箭头，其对应的每个动作都是最优的。

	1	2	3	4	5
1	→	←↑→↓ (A)	←	←↑→↓ (B)	←
2	↑→	↑	←↑	←	←
3	↑→	↑	←↑	←↑ (B')	←↑
4	↑→	↑	←↑	←↑	←↑
5	↑→	↑ (A')	←↑	←↑	←↑