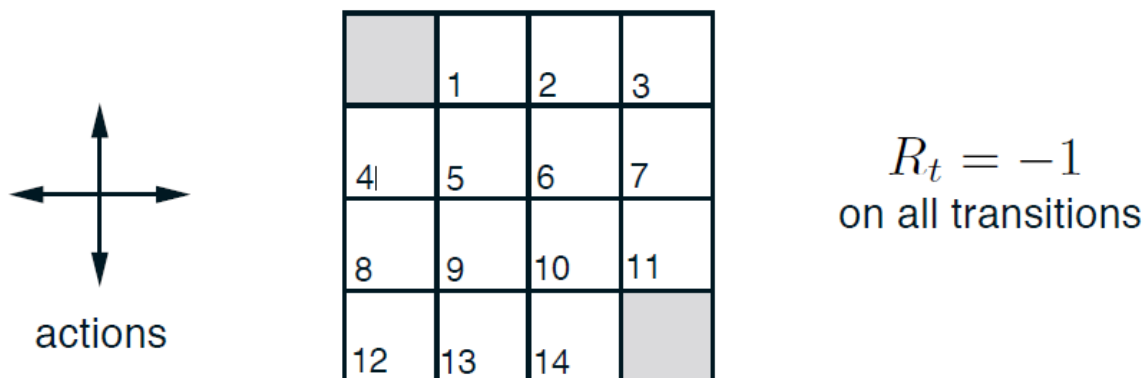


强化学习基础篇（二十二）DP小型网格问题

该问题基于《Reinforcement Learning: An Introduction》在第四章的例4.1。

1、问题描述

考虑下面的这个4*4的网格图



非终止状态集合 $S = 1, 2, \dots, 14$ 。每个状态有四种可能的动作, $A = up, down, left, right$ 。每个动作会导致状态转移, 但当动作会导致智能体移出网格时, 状态保持不变。比如, $p(6, -1 | 5, right) = 1$, $p(7, -1 | 7, right) = 1$ 和对于任意 $r \in R$, 都有 $p(10, r | 5, right) = 0$ 。这是一个无折扣的分幕式任务。在到达终止状态之前, 所有动作的收益均为-1。终止状态在图中以阴影显示 (尽管图中显示了两个格子, 但实际仅有一个终止状态)。对于所有的状态 s, s' 以及动作 a , 期望的收益函数均为 $r(s, a, s') = -1$ 。假设智能体采取等概率随机策略 (所有动作等可能执行), 我们需要计算在迭代策略评估中价值函数序列的收敛情况。

2、实现过程

2.1、环境定义

首先导入库函数以及定义环境信息:

```
1 import matplotlib
2 import matplotlib.pyplot as plt
3 import numpy as np
4 from matplotlib.table import Table
5
6 # 定义网格世界大小
7 WORLD_SIZE = 4
8
9 # 把动作定义为对x, y坐标的增减改变
10 # left, up, right, down
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 # 该问题中每个动作选择的概率为0.25
16 ACTION_PROB = 0.25
17 # 定义画图会用到的动作
18 ACTIONS_FIGS = ['←', '↑', '→', '↓']
```

然后定义左上角和右下角两个坐标 (0,0) 与 (WORLD_SIZE - 1, WORLD_SIZE - 1) 两个点为terminal

```
1 def is_terminal(state):
2     '''
3     返回是否为terminal
4     '''
5     x, y = state
6     return (x == 0 and y == 0) or (x == WORLD_SIZE - 1 and y == WORLD_SIZE - 1)
```

2.2、定义动作执行过程

```
1 def step(state, action):
2     # 当到达terminal时, 下一步状态不变, 奖励为0
3     if is_terminal(state):
4         return state, 0
5     # 计算下一个状态
6     next_state = (np.array(state) + action).tolist()
7     x, y = next_state
8     # 当运动未知超出格子世界则在原位置不变
9     if x < 0 or x >= WORLD_SIZE or y < 0 or y >= WORLD_SIZE:
10         next_state = state
11
12     reward = -1
13     return next_state, reward
```

2.3、辅助函数

以下辅助函数主要用于画图

```
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        tb.add_cell(i, j, width, height, text=val,
12                    loc='center', facecolor='white')
13
14    # Row and column labels...
15    for i in range(len(image)):
16        tb.add_cell(i, -1, width, height, text=i+1, loc='right',
17                    edgecolor='none', facecolor='none')
18        tb.add_cell(-1, i, width, height/2, text=i+1, loc='center',
19                    edgecolor='none', facecolor='none')
20    ax.add_table(tb)
```

以下辅助函数用户策略描述:

```
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] == [0,0]:
23            val = "terminal"
24        if [i, j] == [WORLD_SIZE - 1,WORLD_SIZE - 1]:
25            val = "terminal"
26
27        tb.add_cell(i, j, width, height, text=val,
28                    loc='center', facecolor='white')
29
30    # Row and column labels...
31    for i in range(len(optimal_values)):
32        tb.add_cell(i, -1, width, height, text=i+1, loc='right',
33                    edgecolor='none', facecolor='none')
34        tb.add_cell(-1, i, width, height/2, text=i+1, loc='center',
35                    edgecolor='none', facecolor='none')
36
37    ax.add_table(tb)

```

2.4、使用迭代策略评估算法估算状态值函数

使用迭代策略评估算法估算状态值函数，遵循的算法如下：

Iterative Policy Evaluation, for estimating $V \approx v_\pi$

Input π , the policy to be evaluated

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation

Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

$\Delta \leftarrow 0$

Loop for each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

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

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$

这里同时考虑了是否使用In-Place动态规划（In-place dynamic programming）。

在基于同步动态规划的值迭代算法中，存储了两个值函数的备份，分别是 $v_{new}(s)$ 和 $v_{old}(s)$ 。

$$v_{new}(s) = \max_a (r + \gamma \sum_{s' \in S} p(s'|s, a) v_{old}(s'))$$

即在计算过程中，通过赋值的方式使旧的状态值作为下一次计算新的状态值。

而In-place动态规划 (In-Place Dynamic Programming, IPDP) 则是去掉旧的状态值 $v_{old}(s)$ ，只保留最新的状态值 $v_{new}(s)$ ，在更新的过程中可以减少存储空间的浪费。

$$v(s) = \max_a (r + \gamma \sum_{s' \in S} p(s'|s, a) v(s'))$$

直接原地更新下一个状态值 $v(s)$ ，而不像同步迭代那样需要额外存储新的状态值 $v_{new}(s)$ 。在这种情况下，按何种次序更新状态值有时候会更具有意义。

```
1 def compute_state_value(in_place=True, discount=1.0):
2     # 初始化状态值函数为0
3     new_state_values = np.zeros((WORLD_SIZE, WORLD_SIZE))
4     # 在中间几个迭代进行可视化绘图
5     draw_iteration=[0,1,2,3,10]
6     iteration = 0
7     while True:
8         if iteration in draw_iteration:
9             draw_image(np.round(new_state_values, decimals=2))
10        # 判断是否使用In-Place动态规划 (In-place dynamic programming)
11        if in_place:
12            state_values = new_state_values
13        else:
14            state_values = new_state_values.copy()
15            old_state_values = state_values.copy()
16
17        # 遍历所有状态
18        for i in range(WORLD_SIZE):
19            for j in range(WORLD_SIZE):
20                value = 0
21                # 遍历所有动作,按DP算法更新
22                for action in ACTIONS:
23                    (next_i, next_j), reward = step([i, j], action)
24                    value += ACTION_PROB * (reward + discount *
state_values[next_i, next_j])
25                    new_state_values[i, j] = value
26
27        # 误差小于门限则停止更新
28        max_delta_value = abs(old_state_values - new_state_values).max()
29        if max_delta_value < 1e-4:
30            draw_image(np.round(new_state_values, decimals=2))
31            break
32
33        iteration += 1
34
35    return new_state_values, iteration
```

3、实验结果

运行如下代码测试在使用In-place动态规划 (In-Place Dynamic Programming, IPDP) 的结果：

```
1 _, asyncn_iteration = compute_state_value(in_place=True)
```

结果整理后如下：

1 | In-place: 113 iterations

K=0

	1	2	3	4
1	0.0	0.0	0.0	0.0
2	0.0	0.0	0.0	0.0
3	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0	0.0

	1	2	3	4
1	terminal	↖↑↗↓	↖↑↗↓	↖↑↗↓
2	↖↑↗↓	↖↑↗↓	↖↑↗↓	↖↑↗↓
3	↖↑↗↓	↖↑↗↓	↖↑↗↓	↖↑↗↓
4	↖↑↗↓	↖↑↗↓	↖↑↗↓	terminal

K=1

	1	2	3	4
1	0.0	-1.0	-1.25	-1.31
2	-1.0	-1.5	-1.69	-1.75
3	-1.25	-1.69	-1.84	-1.9
4	-1.31	-1.75	-1.9	0.0

	1	2	3	4
1	terminal	←	←	←
2	↑	↖↑	↑	↑
3	↑	←	↖↑	↓
4	↑	←	→	terminal

K=2

	1	2	3	4
1	0.0	-1.94	-2.55	-2.73
2	-1.94	-2.81	-3.24	-3.4
3	-2.55	-3.24	-3.57	-3.22
4	-2.73	-3.4	-3.22	0.0

	1	2	3	4
1	terminal	←	←	←
2	↑	↖↑	↑	↑
3	↑	←	↖↓	↓
4	↑	←	→	terminal

K=3

	1	2	3	4
1	0.0	-2.82	-3.83	-4.18
2	-2.82	-4.03	-4.71	-4.88
3	-3.83	-4.71	-4.96	-4.26
4	-4.18	-4.88	-4.26	0.0

	1	2	3	4
1	terminal	←	←	←
2	↑	←↑	↑	↑
3	↑	←	→↓	↓
4	↑	←	→	terminal

K=10

	1	2	3	4
1	0.0	-7.83	-11.12	-12.23
2	-7.83	-10.42	-11.77	-11.86
3	-11.12	-11.77	-11.05	-8.81
4	-12.23	-11.86	-8.81	0.0

	1	2	3	4
1	terminal	←	←	←
2	↑	←↑	←	↓
3	↑	↑	→↓	↓
4	↑	→	→	terminal

K=113

	1	2	3	4
1	0.0	-14.0	-20.0	-22.0
2	-14.0	-18.0	-20.0	-20.0
3	-20.0	-20.0	-18.0	-14.0
4	-22.0	-20.0	-14.0	0.0

	1	2	3	4
1	terminal	←	←	←↓
2	↑	←↑	←↓	↓
3	↑	↑→	→↓	↓
4	↑→	→	→	terminal

运行如下代码测试在不使用In-place动态规划（In-Place Dynamic Programming, IPDP）的结果：

```
1 | values, sync_iteration = compute_state_value(in_place=False)
```

结果整理后如下：

```
1 | synchronous: 172 iterations
```

K=0

	1	2	3	4
1	0.0	0.0	0.0	0.0
2	0.0	0.0	0.0	0.0
3	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0	0.0

	1	2	3	4
1	terminal	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$
2	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$
3	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$
4	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$	terminal

↩

K=1

	1	2	3	4
1	0.0	-1.0	-1.0	-1.0
2	-1.0	-1.0	-1.0	-1.0
3	-1.0	-1.0	-1.0	-1.0
4	-1.0	-1.0	-1.0	0.0

	1	2	3	4
1	terminal	\leftarrow	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$
2	\uparrow	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$
3	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$	\downarrow
4	$\leftarrow \uparrow \rightarrow \downarrow$	$\leftarrow \uparrow \rightarrow \downarrow$	\rightarrow	terminal

↩

K=2

	1	2	3	4
1	0.0	-1.75	-2.0	-2.0
2	-1.75	-2.0	-2.0	-2.0
3	-2.0	-2.0	-2.0	-1.75
4	-2.0	-2.0	-1.75	0.0

	1	2	3	4
1	terminal	\leftarrow	\leftarrow	$\leftarrow \uparrow \rightarrow \downarrow$
2	\uparrow	$\leftarrow \uparrow$	$\leftarrow \uparrow \rightarrow \downarrow$	\downarrow
3	\uparrow	$\leftarrow \uparrow \rightarrow \downarrow$	$\rightarrow \downarrow$	\downarrow
4	$\leftarrow \uparrow \rightarrow \downarrow$	\rightarrow	\rightarrow	terminal

↩

K=3

	1	2	3	4
1	0.0	-2.44	-2.94	-3.0
2	-2.44	-2.88	-3.0	-2.94
3	-2.94	-3.0	-2.88	-2.44
4	-3.0	-2.94	-2.44	0.0

	1	2	3	4
1	terminal	←	←	←↓
2	↑	←↑	←↓	↓
3	↑	↑→	→↓	↓
4	↑→	→	→	terminal

↙

K=10

	1	2	3	4
1	0.0	-6.14	-8.35	-8.97
2	-6.14	-7.74	-8.43	-8.35
3	-8.35	-8.43	-7.74	-6.14
4	-8.97	-8.35	-6.14	0.0

	1	2	3	4
1	terminal	←	←	←↓
2	↑	←↑	←↓	↓
3	↑	↑→	→↓	↓
4	↑→	→	→	terminal

↙

K=172

	1	2	3	4
1	0.0	-14.0	-20.0	-22.0
2	-14.0	-18.0	-20.0	-20.0
3	-20.0	-20.0	-18.0	-14.0
4	-22.0	-20.0	-14.0	0.0

	1	2	3	4
1	terminal	←	←	←↓
2	↑	←↑	←↓	↓
3	↑	↑→	→↓	↓
4	↑→	→	→	terminal

↙