

问题：

Jack管理着一家有两个场地的小型租车公司（分别称为first location和second location，每租出一部车，Jack可赚10美金。为了提高车子的出租率，Jack在夜间调配车辆，即把车子从一个场地调配到另外一个场地，成本是2美金每辆。假设每个场地每天租车和还车的数量是泊松随机变量，即其数值是 n 的概率为 $\frac{\lambda^n}{n!} e^{-\lambda}$ ，其中 λ 为期望。假设场地1和场地2租车的 λ 分别为3和4，还车的分别 λ 为3和2。为了简化问题起见，我们假设每个场地最多可停20部车（如果归还的车辆超出了20部，我们假设超出的车辆无偿调配到了别的地方，比如总公司），并且每个场地每天最多调配5部车子。

请问Jack在每个场地应该部署多少部车子？每天晚上如何调配？

问题分析：

- 状态空间：1号租车点和2号租车点分别拥有的可供租赁的车辆，不大于20辆；

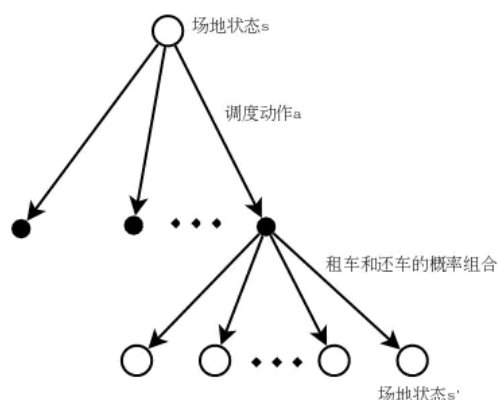
状态数量 $21 \times 21 = 441$

- 行为空间：下班后从一个租车点转移到另一个租车点的车辆，不大于5两；

动作集合

$$A = \{(-5, 5), (-4, 4), (-3, 3), (-2, 2), (-1, 1), (0, 0), (1, -1), (2, -2), (3, -3), (4, -4), (5, -5)\}$$

场地状态和动作和租车换车的关系如下：



- 即时奖励：每租出去一辆车获得10美金；
- 转移概率：租车和换车的数量是服从泊松分布的随机变量，参数见题目描述；

· 折扣因子； 0.9

· 一个分析1号租车点早晨有10辆车，晚上有0辆车收益的例子：

令收租行为 $A_{rent,return}$ ，则早晨1号租车点有10辆车，晚上有0辆车的收租行为是：

$$A_{rent,return} = \begin{bmatrix} 10 & 0 \\ 11 & 1 \\ \dots & \dots \\ 20 & 10 \end{bmatrix} \quad (1)$$

因为收租事件是相互独立的事件且服从泊松分布，所以要计算某个行为出现的概率直接将两者发生的概率相乘。结合条件概率公式还要与傍晚剩0辆车的概率相除，计算出 $P(A_{rent,return}|S' = 0)$ 。

这样子傍晚剩0辆车的收益期望为：

则一天的收益期望可以写为：

$$R(S = 10|S' = 0) = 10 \begin{bmatrix} \frac{P(A_{rent}=10)P(A_{return}=0)}{P(S'=0)} \\ \frac{P(A_{rent}=11)P(A_{return}=1)}{P(S'=0)} \\ \dots \\ \frac{P(A_{rent}=20)P(A_{return}=10)}{P(S'=0)} \end{bmatrix}^T \quad (2)$$

其中 $P(S' = 0) = \sum P(A_{rent})P(A_{return})$

加权平均计算：

$$R(S = 10) = P(S' = 0, 1, 2, \dots, 20)R^T(S = 10|S' = 0, 1, 2, \dots, 20) \quad (3)$$

将所有状态按照上述方法计算后可以得到两个租车点的奖励： $[R_1(S), R_2(S)]$

代码：

```
1 import matplotlib
2 import matplotlib.pyplot as plt
3 import numpy as np
4 from numpy.core.numeric import _ones_like_dispatcher
5 import seaborn as sns
6 from scipy.stats import poisson
```

```

7
8 MAX_CARS = 20
9
10 MAX_MOVE_CARS = 5
11
12 AVG_RENT_A = 3
13 AVG_RENT_B = 4
14 AVG_RETURN_A = 3
15 AVG_RETURN_B = 2
16
17 DISCOUNT = 0.9
18
19 RENT_BONUS = 10
20
21 MOVE_COST = 2
22
23 actions = np.arange(-MAX_MOVE_CARS, MAX_MOVE_CARS + 1)
24
25 POISSON_UPPER_BOUND = 11
26
27 poisson_cache = dict()
28
29 def poisson_probability(n, lam):
30     global poisson_cache
31     key = n * 10 + lam
32     if key not in poisson_cache:
33         poisson_cache[key] = poisson.pmf(n, lam)
34     return poisson_cache[key]
35
36 # 计算状态价值
37 def cal_v(state, action, state_value, returned_cars):
38     """
39     @state: 每个地点的车辆数
40     @action: 移动
41     @state_value: 状态价值矩阵
42     @returned_cars: 还车数目
43     """
44     returns = 0.0
45     returns -= MOVE_COST * abs(action)
46
47     NUM_OF_CARS_A = min(state[0] - action, MAX_CARS)
48     NUM_OF_CARS_B = min(state[1] + action, MAX_CARS)
49

```

```

50     # 遍历两地全部可能概率下租车数目请求
51     for rent_req_A in range(POISSON_UPPER_BOUND):
52         for rent_req_B in range(POISSON_UPPER_BOUND):
53             prob = poisson_probability(rent_req_A, AVG_RENT_A)
54             * \
55             poisson_probability(rent_req_B, AVG_RENT_B)
56
57             num_of_cars_A = NUM_OF_CARS_A
58             num_of_cars_B = NUM_OF_CARS_B
59
60             valid_rent_A = min(num_of_cars_A, rent_req_A)
61             valid_rent_B = min(num_of_cars_B, rent_req_B)
62
63             reward = (valid_rent_A + valid_rent_B) * RENT_BONUS
64             num_of_cars_A -= valid_rent_A
65             num_of_cars_B -= valid_rent_B
66
67             # 如果还车的数目为泊松分布的均值
68             if returned_cars:
69                 returned_cars_A = AVG_RETURN_A
70                 returned_cars_B = AVG_RETURN_B
71
72                 num_of_cars_A = min(num_of_cars_A +
73 returned_cars_A, MAX_CARS)
74                 num_of_cars_B = min(num_of_cars_B +
75 returned_cars_B, MAX_CARS)
76
77             # 策略评估
78             returns += prob * (reward + DISCOUNT *
79 state_value[num_of_cars_A, num_of_cars_B])
80
81             # 计算所有泊松概率分布下的还车空间
82             else:
83                 for returned_cars_A in
84 range(POISSON_UPPER_BOUND):
85                     for returned_cars_B in
86 range(POISSON_UPPER_BOUND):
87                         prob_return =
88 poisson_probability(returned_cars_A, AVG_RETURN_A) *
89 poisson_probability(returned_cars_B, AVG_RETURN_B)
90                         num_of_cars_A_ = min(num_of_cars_A +
91 returned_cars_A, MAX_CARS)

```

```

83         num_of_cars_B_ = min(num_of_cars_B +
returned_cars_B, MAX_CARS)
84         # 联合概率
85         prob_ = prob_return * prob
86         returns += prob_ * (reward + DISCOUNT)
* state_value[num_of_cars_A_, num_of_cars_B_]
87     return returns
88
89
90 def draw(constant_returned_cars = True):
91     value = np.zeros((MAX_CARS + 1, MAX_CARS + 1))
92     policy = np.zeros(value.shape, dtype=np.int)
93     iterations = 0
94
95     # 准备画布，准备多个子图
96     _, axes = plt.subplots(2, 3, figsize=(40, 20))
97     # 调整子图间距
98     plt.subplots_adjust(wspace=0.1, hspace=0.2)
99     # 将子图形成1 * 6 的列表
100    axes = axes.flatten()
101    while True:
102        # 使用seaborn的heatmap作图
103        fig = sns.heatmap(np.flipud(policy), cmap="YlGnBu",
ax=axes[iterations])
104
105        fig.set_ylabel('# cars at A', fontsize = 30)
106        fig.set_yticks(list(reversed(range(MAX_CARS + 1))))
107        fig.set_xlabel('# cars at B', fontsize = 30)
108        fig.set_title('policy{}'.format(iterations), fontsize =
30)
109
110    while True:
111        old_value = value.copy()
112        for i in range(MAX_CARS + 1):
113            for j in range(MAX_CARS + 1):
114                # 更新v(s)
115                new_state_value = cal_v([i,j], policy[i,j],
value, constant_returned_cars)
116
117                value[i,j] = new_state_value
118                max_value_change = abs(old_value - value).max()
119                print('max value change
{}'.format(max_value_change))

```

```

120         if(max_value_change < 1e-4):
121             break
122
123     policy_stable = True
124     # i,j 为 A、B两地的现有车辆
125     for i in range(MAX_CARS + 1):
126         for j in range(MAX_CARS + 1):
127             old_action= policy[i,j]
128             action_returns = []
129
130             for action in actions:
131                 if(0 <= action <= i) or (-j <=action <= 0):
132                     action_returns.append(cal_v([i,j],
action, value, constant_returned_cars))
133                 else:
134                     action_returns.append(-np.inf)
135
136             new_action = actions[np.argmax(action_returns)]
137
138             policy[i,j] = new_action
139
140             if policy_stable and old_action != new_action:
141                 policy_stable = False
142
143     print('policy stable{}'.format(policy_stable))
144     if policy_stable:
145         fig = sns.heatmap(np.flipud(value), cmap="YlGnBu",
ax=axes[-1])
146         fig.set_ylabel('# cars at first location',
fontsize=30)
147         fig.set_yticks(list(reversed(range(MAX_CARS + 1))))
148         fig.set_xlabel('# cars at second location',
fontsize=30)
149         fig.set_title('optimal value', fontsize=30)
150         break
151
152     iterations += 1
153
154     plt.savefig('C:/INFO\RL/assignment/2_Jack/fig.png')
155     plt.close()
156
157 if __name__ == '__main__':
158     draw()

```

结果图:

