

第十六章 强化学习RL(Reinforcement Learning)

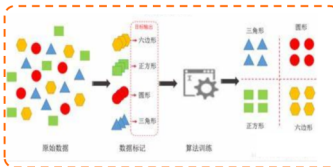
王星

中国人民大学统计学院

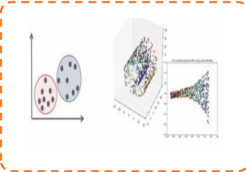
June 1, 2022

学习的分类

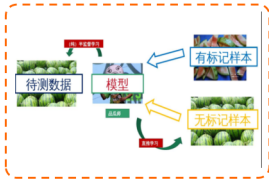
有指导学习



无指导学习



半指导学习



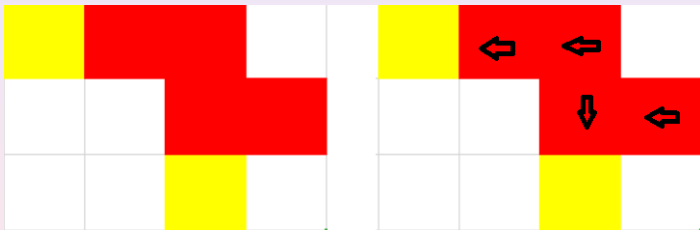
增强学习



- 强化学习简介
- K-摇臂赌博机 (ϵ 贪心+Softmax)
- 模型学习-马尔科夫决策过程
- 策略迭代
- 策略评估
- 时序差分
- Q值学习

16.1 问题来源1/4

场景：考虑一个格子游戏，每个格子等概率可以向4个方向移动，每次移动一步，收益为-1，移动到黄色出口结束游戏。若当前移动会导致出界，则移动后位置不变，要求最佳策略？



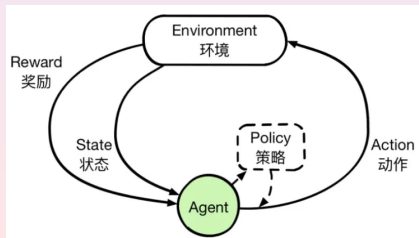
最终可以解出的策略如右图。

强化学习的定义2/4

- 在上述序列决策问题中，有一个智能体，智能体的行动和结果不确定，其行动具有随机性，有生存代价，也有停止条件。
- 智能体执行了某个动作后，环境将转换到一个新的状态，对于该新的状态环境会给出奖励信号（正奖励或者负奖励）。随后，智能体根据新的状态和环境反馈的奖励，按照一定的策略执行新的动作。智能体通过策略学习，可以知道自己在什么状态下，应该采取什么样的动作使得自身获得最大奖励。
- **定义：**强化学习(Reinforcement Learning) 是机器学习中的一个领域，也称为增强学习，评价学习。
- **定义：**强化学习是介于监督学习和非监督学习的另外一种学习方式。它的原理是从错误中学习，找到规律，不断逼近目标(Goal)
- 一个智能体(Agent)不断地在环境(Environment)的交互中学习，通过环境给予的反馈（奖励Reward）不断优化状态(State)- 动作(Action)完成特定目标（比如取得最大奖励目标）。

强化学习的基本元素3/4

- **特点1**: 强化学习关注智能体与环境之间的交互
- **特点2**: 反复实验(trial and error) 和延迟奖励(delayed reward) 是强化学习最重要的两个特征。
- **智能体**: 强化学习的本体, 感知外界作为学习者或者决策者;
- **环境**: 强化学习智能体以外的一切, 主要由状态集合组成;
- **状态**: 记作 s_t , 表示环境的数据, 状态集则是环境中所有可能的状态;
- **动作**: 记作 a_t , 智能体可以做出的动作, 动作集则是智能体可以做出的所有动作;



- **状态转移**:记作 $p(s'|s, a)$, 智能体根据当前状态 s 做出动作 a 后, 下一时刻环境处于不同状态 s' 的概率;
- **即时奖励**:智能体在 t 时刻执行一个动作后, 在 $t + 1$ 时刻获得的正/负反馈信号 r_{t+1} , 奖励集则是智能体可以获得的所有反馈信息, 该奖励与前一时刻的动作、动作和动作之后下一时刻的状态, $r : S \times A \times S' \rightarrow R$;
- **策略**:强化学习是从环境状态到动作的映射学习, 称该映射关系为策略。通俗的理解, 即智能体如何选择动作的思考过程称为策略。
- **目标**:智能体自动寻找在连续时间序列里的最优策略。**最优策略**通常指最大化长期累积奖励 V 如下:

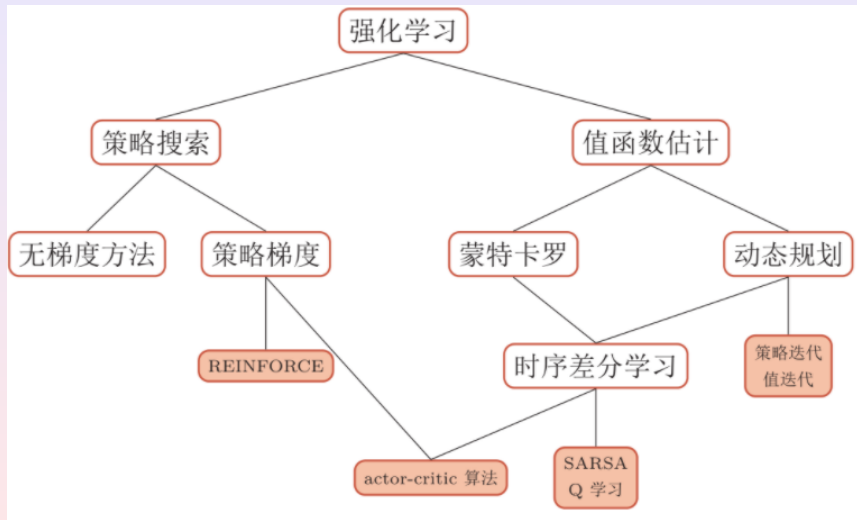
$$V = r_1 + \cdots + r_t$$

强化学习的分类1/2

强化学习可以分为两类：基于价值的学习；基于概率的学习

- 基于价值的学习：输出值的下一步要采取的各种动作的价值，根据价值最高的原则选择动作，经常应用于不连续的动作行为
 - 基于价值的方法（Value Based）：没有策略但是有价值函数:Q-learning,Sarsa算法
 - 参与评价方法（Actor Critic）：既有策略也有价值函数
- 基于概率的学习：通过感官分析所处的环境，输出下一步要采取的各种动作的概率，根据概率选择行动，经常应用于连续的动作行为，常见的算法包括Policy Gradients算法等
 - 无模型的方法: 有策略和价值函数，没有模型，在无模型训练方法中，代理对环境一无所知，代理需要对环境的所有区域进行探索(explore)，意味着会花费大量时间去探索低奖励回报的区域，这样会浪费大量的时间去做无用功，但在大多数情况又不得不面对这样的情形。
 - 基于模型的方法：有策略和价值函数，也有模型；研发者知道环境中某些区域或者全部环境的情况下，研发者可将这些情况告知代理，就像一张地图，会告知代理哪些区域回报极低，避免探索低回报区域，加快训练速度，节约资源。

强化学习的分类



16.2单步学习-K-摇臂赌博机1/4

- 与一般监督学习不同，强化学习任务的最终奖赏是在多步动作之后才能观察到，这里如果仅仅考虑最大化单步奖赏最佳策略依赖于两方面的信息：
 - 需知道每个动作带来的奖赏：
 - 获得奖赏最大的动作；如每个动作对应的奖赏是一个确定值，那么尝试一遍所有的动作就能确定奖赏最大的动作。然而现实中一个动作的奖赏值是来自于一个概率分布，仅通过一次尝试并不能确切地获得平均奖赏值，这就需要探索。
- 一方面，为了从环境中获取尽可能多的知识，需要让Agent进行探索，另一方面，为获得较大的奖励，还需要让agent对已知的信息加以利用。
- 强化学习所面临的“探索-利用窘境” (Exploration-Exploitation dilemma)。显然，欲累积奖赏最大，则必须在探索与利用之间达成较好的折中。强化学习问题中存在的一个重要挑战即是如何权衡探索-利用之间的关系。这也是得到最优策略的算法纷纷努力的主要内容。

- ϵ -贪心算法基于一个概率对探索和利用进行折中：每次尝试时，以 ϵ 概率进行探索，即以均匀概率随机选取一个摇臂；以 $1 - \epsilon$ 的概率进行利用，即选择当前平均奖赏最高的摇臂(若有多个，则随机选取一个)。

$$\pi^\epsilon(s) = \begin{cases} \pi(s), & \text{按概率 } 1 - \epsilon, \\ \text{随机选择 } \mathcal{A} \text{ 中的动作}, & \text{按概率 } \epsilon. \end{cases}$$

- 若摇臂奖赏不确定较大，则需要探索的时间比较长，需要较大的€值；若摇臂奖赏确定性较集中，则需要探索时间较短，需要较小的€值。 $\epsilon = \frac{1}{\sqrt{T}}$

```
输入: 摇臂数  $K$ ;  
      奖赏函数  $R$ ;  
      尝试次数  $T$ ;  
      探索概率  $\epsilon$ .  
  
过程:  
1:  $r = 0$ ;  
2:  $\forall i = 1, 2, \dots, K : Q(i) = 0, \text{count}(i) = 0$ ;  
3: for  $t = 1, 2, \dots, T$  do  
4:   if  $\text{rand}() < \epsilon$  then  
5:      $k =$  从  $1, 2, \dots, K$  中以均匀分布随机选取  
6:   else  
7:      $k = \arg \max_i Q(i)$   
8:   end if  
9:    $v = R(k)$ ;  
10:   $r = r + v$ ;  
11:   $Q(k) = \frac{Q(k) \times \text{count}(k) + v}{\text{count}(k) + 1}$ ;  
12:   $\text{count}(k) = \text{count}(k) + 1$ ;  
13: end for  
输出: 累积奖赏  $r$ 
```

- 若各摇臂的平均奖赏相当，选取各摇臂的概率也一样；若某些摇臂的平均奖赏明显高于其他摇臂，则被选中的概率也明显更高。
- Softmax 算法中摇臂概率的分配是基于 Boltzmann 分布：
- 其中， $Q(i)$ 记录当前摇臂的平均奖赏； $\tau > 0$ 称为“温度”， τ 越小平均奖赏高的摇臂被选取的概率越高； τ 越趋于 0 时，Softmax 将趋于“仅利用”； τ 趋于无穷大时，Softmax 则趋于“仅探索”。

输入：摇臂数 K ；
奖赏函数 R ；
尝试次数 T ；
温度参数 τ 。

过程：

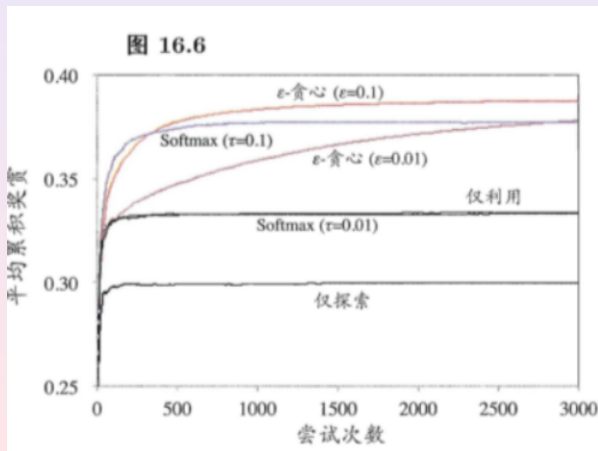
```
1:  $r = 0$ ;  
2:  $\forall i = 1, 2, \dots, K : Q(i) = 0, \text{count}(i) = 0$ ;  
3: for  $t = 1, 2, \dots, T$  do  
4:    $k =$  从  $1, 2, \dots, K$  中根据式(16.4)随机选取  
5:    $v = R(k)$ ;  
6:    $r = r + v$ ;  
7:    $Q(k) = \frac{Q(k) \times \text{count}(k) + v}{\text{count}(k) + 1}$ ;  
8:    $\text{count}(k) = \text{count}(k) + 1$ ;  
9: end for
```

即根据 Softmax 函数

$$P(k) = \frac{e^{\frac{Q(k)}{\tau}}}{\sum_{i=1}^K e^{\frac{Q(i)}{\tau}}}$$

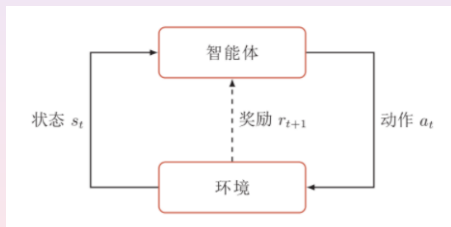
输出：累积奖赏 r

- 离散状态空间，离散动作空间上多步强化学习任务可以将每个步骤看成是一个 K -摇臂赌博机，这样做事比较简单的一种方式，它没有考虑马尔科夫决策过程的结构



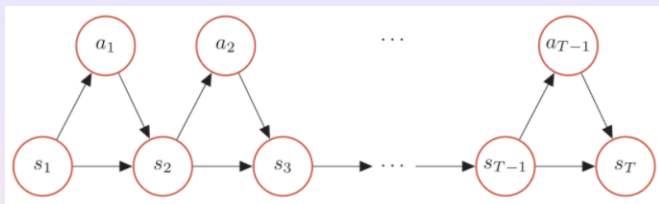
16.3 多步学习-马尔科夫决策过程1/3

马尔科夫决策过程（Markov Decision Process, MDP）适用于一种环境完全可观测的强化学习环境，即马尔科夫性质，结果部分随机、部分在决策者的控制下的决策过程建模的数学框架。在智能体与环境进行交互的过程中，根据经验调整策略最优化决策序列的过程。考虑多步学习四元组 $E = \{S, A, P, R\}$ 均为已知。



$$p(s_{t+1}|s_t, a_t, \dots, s_0, a_0) = p(s_{t+1}|s_t, a_t);$$

16.2 马尔可夫决策过程2/3



- 马尔可夫决策过程是一条“轨迹” (trajectory)

$\tau = s_0 \rightarrow a_0 \rightarrow s_1 \rightarrow r_1 \rightarrow a_1 \rightarrow \dots, a_{t-1} \rightarrow s_t \rightarrow r_t, \dots$

- 策略 $\pi(a|s)$

$$\begin{aligned} p(\tau) &= p(s_0, a_0, s_1, a_1, \dots) \\ &= p(s_0) \prod_{t=0}^{T-1} \pi(a_t | s_t) p(s_{t+1} | s_t, a_t) \end{aligned}$$

有模型的学习:策略评估

在模型已知的前提下,可以对任意策略的进行评估。一般常使用以下两种值函数来评估某个策略的优劣:

- **状态值函数 (V)**: 从状态 s 出发, 使用策略 π 所带来的累积奖赏;
- **状态动作值函数 (Q)**: 从状态 s 出发, 执行动作 a 后再使用 π 策略所带来的累积奖赏。

根据累积奖赏的定义, 可以引入 T 步累积奖赏与 γ 折扣累积奖赏:

- 期望累积奖赏函数
 - **T 步累积奖赏**: $V_T(s) = E[\frac{1}{T} \sum_{t=1}^T r_t | s_0 = x_0]$;
 - **γ 折累积奖赏**: $V_\gamma(s) = E[\sum_{t=0}^{+\infty} \gamma^t r_{t+1} | s_0 = x_0]$
- 状态值函数:
 - **T 步累积奖赏**: $Q_T(s) = E[\frac{1}{T} \sum_{t=1}^T r_t | s_0 = x_0, a_0 = a]$;
 - **γ 折累积奖赏**: $Q_\gamma(s) = E[\sum_{t=0}^{+\infty} \gamma^t r_{t+1} | s_0 = x_0, a_0 = a]$

有模型学习的策略改进

对某个策略的积累奖赏进行评估后发现非最优，改进最优状态动作值函数：

由于最优值函数的累积奖赏值已达最大，因此可对前面的 Bellman 等式(16.7)和(16.8)做一个改动，即将对动作的求和改为取最优：

$$\begin{cases} V_T^*(x) = \max_{a \in A} \sum_{x' \in X} P_{x \rightarrow x'}^a \left(\frac{1}{T} R_{x \rightarrow x'}^a + \frac{T-1}{T} V_{T-1}^*(x') \right); \\ V_\gamma^*(x) = \max_{a \in A} \sum_{x' \in X} P_{x \rightarrow x'}^a (R_{x \rightarrow x'}^a + \gamma V_\gamma^*(x')). \end{cases} \quad (16.13)$$

换言之，

$$V^*(x) = \max_{a \in A} Q^{\pi^*}(x, a). \quad (16.14)$$

代入式(16.10)可得最优状态-动作值函数

$$\begin{cases} Q_T^*(x, a) = \sum_{x' \in X} P_{x \rightarrow x'}^a \left(\frac{1}{T} R_{x \rightarrow x'}^a + \frac{T-1}{T} \max_{a' \in A} Q_{T-1}^*(x', a') \right); \\ Q_\gamma^*(x, a) = \sum_{x' \in X} P_{x \rightarrow x'}^a (R_{x \rightarrow x'}^a + \gamma \max_{a' \in A} Q_\gamma^*(x', a')). \end{cases} \quad (16.15)$$

上述关于最优值函数的等式，称为最优 Bellman 等式

有模型学习的策略迭代: T步奖赏累积策略迭代

- 策略迭代: 不断迭代进行策略评估和改进知道策略收敛不在改变位置

图 16.8

输入: MDP 四元组 $E = \langle X, A, P, R \rangle$;

累积奖赏参数 T .

过程:

```
1:  $\forall x \in X : V(x) = 0, \pi(x, a) = \frac{1}{|A(x)|}$ ;
2: loop
3:   for  $t = 1, 2, \dots$  do
4:      $\forall x \in X : V'(x) = \sum_{a \in A} \pi(x, a) \sum_{x' \in X} P_{x \rightarrow x'}^a (\frac{1}{t} R_{x \rightarrow x'}^a + \frac{t-1}{t} V(x'))$ ;
5:     if  $t = T + 1$  then
6:       break
7:     else
8:        $V = V'$ 
9:     end if
10:  end for
11:  $\forall x \in X : \pi'(x) = \arg \max_{a \in A} Q(x, a)$ ;
12: if  $\forall x : \pi'(x) = \pi(x)$  then
13:   break
14: else
15:    $\pi = \pi'$ 
16: end if
17: end loop
```

输出: 最优策略 π

有模型学习的值迭代：T步奖赏累积奖赏值迭代

- 值迭代：减少耗时

图 16.9

输入: MDP 四元组 $E = \langle X, A, P, R \rangle$;

累积奖赏参数 T ;

收敛阈值 θ .

过程:

1: $\forall x \in X : V(x) = 0$;

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

3: $\forall x \in X : V'(x) = \max_{a \in A} \sum_{x' \in X} P_{x \rightarrow x'}^a (\frac{1}{t} R_{x \rightarrow x'}^a + \frac{t-1}{t} V(x'))$;

4: **if** $\max_{x \in X} |V(x) - V'(x)| < \theta$ **then**

5: **break**

6: **else**

7: $V = V'$

8: **end if**

9: **end for**

输出: 策略 $\pi(x) = \arg \max_{a \in A} Q(x, a)$

免模型的学习-时序差分学习

时序差分(Temporal Difference, 简称TD) 学习结合了动态规划与蒙特卡罗方法的思想, 能做到更高效的免模型学习。对于状态-动作 (x, a) 不妨假设基于 t 个采样已估计出值函数 $Q_t^\pi(x, a) = \frac{1}{t} \sum_{i=1}^t r_i$, 在得到第 $t + 1$ 个采样 r_{t+1} 时, 有

$$Q_{t+1}^\pi(x, a) = Q_t^\pi(x, a) + \frac{1}{t+1}(r_{t+1} - Q_t^\pi(x, a)).$$

SARSA算法(State Action Reward State Action, SARSA)

```
输入: 环境  $E$ ;  
      动作空间  $A$ ;  
      起始状态  $x_0$ ;  
      奖赏折扣  $\gamma$ ;  
      更新步长  $\alpha$ .  
过程:  
1:  $Q(x, a) = 0, \pi(x, a) = \frac{1}{|A(x)|}$ ;  
2:  $x = x_0, a = \pi(x)$ ;  
3: for  $t = 1, 2, \dots$  do  
4:    $r, x' =$  在  $E$  中执行动作  $a$  产生的奖赏与转移的状态;  
5:    $a' = \pi^\epsilon(x')$ ;  
6:    $Q(x, a) = Q(x, a) + \alpha(r + \gamma Q(x', a') - Q(x, a))$ ;  
7:    $\pi(x) = \arg \max_{a'} Q(x, a')$ ;  
8:    $x = x', a = a'$   
9: end for  
输出: 策略  $\pi$ 
```

图 16.12 Sarsa 算法

无模型的学习-Q-学习

输入: 环境 E ;

动作空间 A ;

起始状态 x_0 ;

奖赏折扣 γ ;

更新步长 α .

过程:

1: $Q(x, a) = 0, \pi(x, a) = \frac{1}{|A(x)|}$;

2: $x = x_0$;

3: **for** $t = 1, 2, \dots$ **do**

4: $r, x' =$ 在 E 中执行动作 $\pi^t(x)$ 产生的奖赏与转移的状态;

5: $a' = \pi(x')$;

6: $Q(x, a) = Q(x, a) + \alpha(r + \gamma Q(x', a') - Q(x, a))$;

7: $\pi(x) = \arg \max_{a''} Q(x, a'')$;

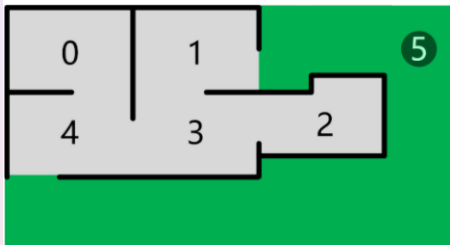
8: $x = x', a = a'$

9: **end for**

输出: 策略 π

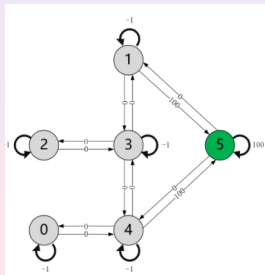
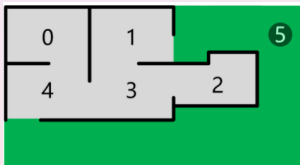
图 16.13 Q-学习算法

- **问题:**大楼有一组房间，每个房间编号为0 ~ 4，有的房间之间可互通，有的不能互通，5号区域为楼外。如下图所示：



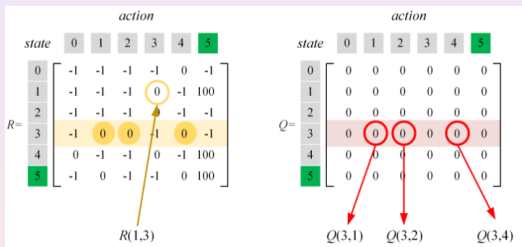
- **目标:**将Agent放在任何房间，它自己能够从目标房间5号走出大楼。

- 构造连通图，并约定：能够直接到达5号区域的边，其奖励值为100（包括5号自己到达5号）；房间之间彼此连通的边，其奖励值为0；房间之间彼此阻隔的边，其奖励值为-1；房间自己到达自己的边，其奖励值为-1。根据这样的规则，构造如下连通图（对于互相阻隔的房间，它们之间的边被省略了）：



- 定义state和action。约定将代理当前所在的区域编号作为state，将代理前往哪一个区域的编号作为action。

- 奖励 (Reward)。指“从房间1 (当前state) 前往房间3 (采取的action) 所获得的即时奖励”，因此需要查询的数据为 $R(1,3)$ 。
- 未来的价值 (Value)。指“进入房间3后 (next state) 所能获得的最大即时奖励”，需要查询的数据为： $Q(3,1)$ 、 $Q(3,2)$ 和 $Q(3,4)$ 。



- Q值的更新算式

$$Q(1,3) \leftarrow R(1,3) + 0.8 \cdot \max(Q(3,1), Q(3,2), Q(3,4))$$

$$\leftarrow 0 + 0.8 \cdot (0, 0, 0) = 0$$

- 选择另一个动作“前往区域5”，通过查询Q-table更新Q值,具体更新的是元素 $Q(1,5)$ ，此时，需要查表获得两类数据： $R(1,5)Q(5,1), Q(5,4), Q(5,5)$

$$Q(1,5) \leftarrow R(1,5) + 0.8 \cdot \max(Q(5,1), Q(5,4), Q(5,5))$$

$$\leftarrow 100 + 0.8 \cdot (0, 0, 0) = 100$$

| | | action | | | | | |
|-------|-------|--------|---|---|---|---|-----|
| | state | 0 | 1 | 2 | 3 | 4 | 5 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| $Q =$ | 1 | 0 | 0 | 0 | 0 | 0 | 100 |
| | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 3 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 4 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 5 | 0 | 0 | 0 | 0 | 0 | 0 |

| | | action | | | | | |
|-------|-------|--------|----|------|----|----|-----|
| | state | 0 | 1 | 2 | 3 | 4 | 5 |
| | 0 | 0 | 0 | 0 | 0 | 80 | 0 |
| $Q =$ | 1 | 0 | 0 | 0 | 64 | 0 | 100 |
| | 2 | 0 | 0 | 0 | 64 | 0 | 0 |
| | 3 | 0 | 80 | 51.2 | 0 | 80 | 0 |
| | 4 | 64 | 0 | 0 | 64 | 0 | 100 |
| | 5 | 0 | 0 | 0 | 0 | 0 | 0 |

- 接下来按照上述思路不断进行迭代，最终得到Q-table,在得到了Q-table后，可以根据该矩阵令代理独自行走寻找路径了。基本想法是查找当前位置具有最大Q值的action，然后执行，如此反复，直到到达目的地为止。
- 初始化 s ，如果不是出口，循环选择Q最大的方向前进，更新 s 直到结束。
- 具体代码如下：

```
1 import numpy as np
2 import random
3
4 # Build Q-table
5 q = np.zeros((6, 6))
6 q = np.matrix(q)
7
8 # Build R-table
9 r = np.array([[[-1, -1, -1, -1, 0, -1], [-1, -1, -1, 0, -1, 100], [-1, -1, 0, -1, -1], [-1, 0, 0, -1, 0, -1],
10               [0, -1, -1, 0, -1, 100], [-1, 0, -1, -1, 0, 100]])]
11 r = np.matrix(r)
12
13 # Discount factor
14 gamma = 0.8
15
16 # Training
17 for i in range(1000):
18     # For each training episode, randomly select a state
19     state = random.randint(0, 5)
20     while state != 5:
21         # Select the action of a non-negative value in the R-table
22         r_pos_action = []
23         for action in range(6):
24             if r[state, action] >= 0:
25                 r_pos_action.append(action)
26         next_state = r_pos_action[random.randint(0, len(r_pos_action) - 1)]
27         q[state, next_state] = r[state, next_state] + gamma * q[next_state].max()
28         state = next_state
29
30 print("Finally, the content of Q-table is:\n")
31 print(q)
```

Q-学习7/8代码实现

```
33 # verification
34 for i in range(10):
35     print("\n{}-th verification:".format(i + 1))
36     state = random.randint(0, 5)
37     print('The robot is at {}'.format(state))
38     count = 0
39     while state != 5:
40         if count > 20:
41             print('fail')
42             break
43         # Choose the largest q_max
44         q_max = q[state].max()
45
46         q_max_action = []
47         for action in range(6):
48             if q[state, action] == q_max:
49                 q_max_action.append(action)
50
51         next_state = q_max_action[random.randint(0, len(q_max_action) - 1)]
52         print("the robot goes to " + str(next_state) + '.')
53         state = next_state
54         count += 1
```

Q-学习8/8验证策略

运行过程中，会输出最终Q-table中的数据，如下图所示；部分验证阶段的输出，对照房间平面图可以发现，代理成功地通过了测试。

```
Finally, the content of Q-table is:  
[[ 0.  0.  0.  0.  80.  0. ]  
 [ 0.  0.  0.  64.  0. 100. ]  
 [ 0.  0.  0.  64.  0.  0. ]  
 [ 0.  80.  51.2  0.  80.  0. ]  
 [ 64.  0.  0.  64.  0. 100. ]  
 [ 0.  0.  0.  0.  0.  0. ]]
```

```
1-th verification:  
The robot is at 2.  
the robot goes to 3.  
the robot goes to 1.  
the robot goes to 5.  
  
2-th verification:  
The robot is at 1.  
the robot goes to 5.  
  
3-th verification:  
The robot is at 1.  
the robot goes to 5.  
  
4-th verification:  
The robot is at 3.  
the robot goes to 1.  
the robot goes to 5.  
  
5-th verification:  
The robot is at 3.  
the robot goes to 1.  
the robot goes to 5.  
  
6-th verification:  
The robot is at 2.  
the robot goes to 3.  
the robot goes to 4.  
the robot goes to 5.  
  
7-th verification:  
The robot is at 3.  
the robot goes to 4.  
the robot goes to 5.  
  
8-th verification:  
The robot is at 1.  
the robot goes to 5.
```