

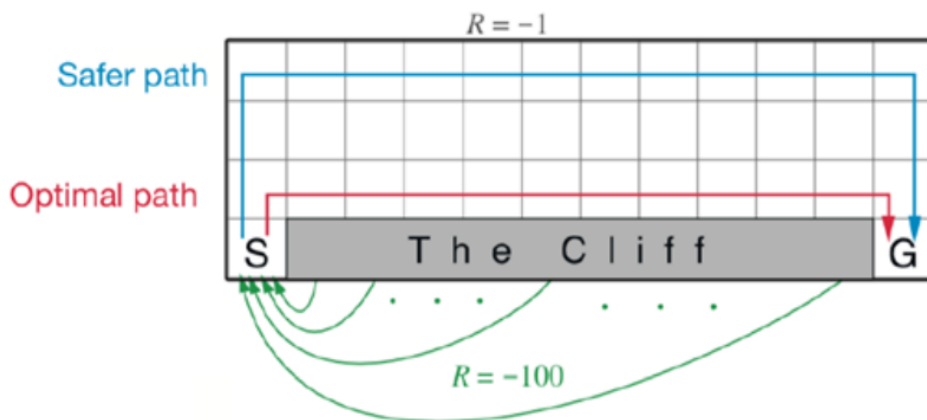
# 实验报告2

学号	姓名
20337025	崔璨明

## 1、实验内容

基于Cliff Walk例子实现SARSA、Q-learning算法，要求如下：

- This is a standard undiscounted, episodic task, with start and goal states, and the usual actions causing movement up, down, right, and left.
- Reward is -1 on all transitions except those into the region marked as Cliff. Stepping into this region incurs a reward -100 and sends the agent instantly back to the start.



## 2、核心思路

在该问题中，对每个状态进行编号，并建立相应的Q-table，在每次的迭代中，初始化状态为角色位于起点，每个状态有四个动作（上下左右），每一步转移的reward都设为-1，若采取某一动作后若到达了边界则待在原地，若到达了悬崖则返回起点，reward设为-100，并设置一个符号变量flag来判断某次迭代是否终止，若到达了终点则该次迭代终止，保留Q表，然后重新开始。

### 2-1 Q-learning

根据以上条件，采用Q-learning算法，该算法的伪代码如下：

```

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat (for each episode):
    Initialize  $S$ 
    Repeat (for each step of episode):
        Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
        Take action  $A$ , observe  $R, S'$ 
         $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ 
         $S \leftarrow S'$ ;
    until  $S$  is terminal

```

参数解释：

$\epsilon - greedy$ 是用在决策上的一种策略，比如  $\epsilon = 0.9$  时，就说明有90% 的情况会按照Q-table 的max-value选择行为，10% 的情况使用随机选择的行为。

$\alpha$ ：学习率，来决定这次的误差有多少是要被学习的

$\gamma$ ：对未来 reward 的衰减值。

## 2-2 SARAS

saras算法假设前一时刻的状态价值的值 $q_{t-1}(s_{t-1})$ 是最优的，利用当前的行动状态值 $q_{t-1}(s_t, a_t)$ 和奖励值 $r_t$ 来更新 $q_t(s_t, a_t)$ ，公式如下：

$$q_t(s_t, a_t) = q_{t-1}(s_{t-1}, a_{t-1}) + \frac{1}{N} (r_t + \gamma * q_{t-1}(s_t, a_t) - q_{t-1}(s_{t-1}, a_{t-1}))$$

算法伪代码如下：

```

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat (for each episode):
    Initialize  $S$ 
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
    Repeat (for each step of episode):
        Take action  $A$ , observe  $R, S'$ 
        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
         $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$ 
         $S \leftarrow S'; A \leftarrow A'$ ;
    until  $S$  is terminal

```

参数解释：

$\epsilon - greedy$ 是用在决策上的一种策略，比如  $\epsilon = 0.9$  时，就说明有90% 的情况会按照Q-table 的max-value选择行为，10% 的情况使用随机选择的行为。

$\alpha$ ：学习率，来决定这次的误差有多少是要被学习的

$\gamma$ ：对未来 reward 的衰减值。

### 3、关键代码展示

根据伪代码，python实现如下：

#### 3-1 Q-learning

```
# Q-learning
def q_learning():
    Q_table=np.zeros((rows*cols,4)) #Q表
    cliff_map=np.zeros((rows,cols)) #地图，用于每次存储路径
    rw_list=[] #记录Reward的变化
    for step in range(episodes):
        #初始化S
        sum=0
        pos=(3,0)# 起点
        end_flag=False# 结束标志
        cliff_map=np.zeros((rows,cols))# 记录路径
        cliff_map[3][0]=1
        while(not end_flag):#判断是否到达目的地
            loc=int(12 * pos[0] + pos[1]) #计算此时的状态
            mov=epsilon_greedy(loc,Q_table) #根据Q表用epsilon_greedy选择动作
            pos=move(pos,mov) # 移动
            cliff_map[pos[0]][pos[1]]=1
            new_loc=int(12 * pos[0] + pos[1]) #计算新状态
            new_loc,rwd, end_flag= reward(new_loc) #计算reward并判断是否结束或掉下悬崖
            sum+=rwd
            max_q_s_a=find_max(Q_table,pos) #选取当前状态最大Q表项
            # 根据公式更新Q表
            Q_table[loc][mov]=Q_table[loc][mov]+alpha*(rwd+(gamma*max_q_s_a)-
            Q_table[loc][mov])
            pos=(int(new_loc/12),int(new_loc%12)) #更新位置
            rw_list.append(sum)
        return cliff_map,Q_table,rw_list
```

#### 3-2 Saras

```
# sarsa
def sarsa():
    Q_table=np.zeros((rows*cols,4))#Q表
    cliff_map=np.zeros((rows,cols))#地图，用于每次存储路径
    rw_list=[]
    for step in range(episodes):
        sum=0
        pos=(3,0)# 起点
```

```

end_flag=False# 结束标志
clif_map=np.zeros((rows,cols))
clif_map[3][0]=1
loc=int(12 * pos[0] + pos[1])
mov=epsilon_greedy(loc,Q_table)
while(not end_flag):
    loc=int(12 * pos[0] + pos[1])#当前状态
    pos=move(pos,mov)#根据动作进行移动
    clif_map[pos[0]][pos[1]]=1
    new_loc=int(12 * pos[0] + pos[1])#新的状态
    new_loc,rwd, end_flag= reward(new_loc)#计算reward
    sum+=rwd
    #选择下一个状态的新动作
    next_mov=epsilon_greedy(new_loc,Q_table)
    q_s_a=Q_table[new_loc][next_mov]
    # 根据公式更新Q表
    Q_table[loc][mov]=Q_table[loc][mov]+alpha*(rwd+(gamma*q_s_a)-
Q_table[loc][mov])
    #更新位置和动作
    mov=next_mov
    pos=(int(new_loc/12),int(new_loc%12))
    rw_list.append(sum)
return clif_map,Q_table,rw_list

```

## 4、实验结果

初始条件设置为：

```

episodes = 500 #迭代次数
epsilon = 0.05 #选择随机方向的概率
alpha = 0.5 #学习率
gamma = 0.99 #衰减值

```

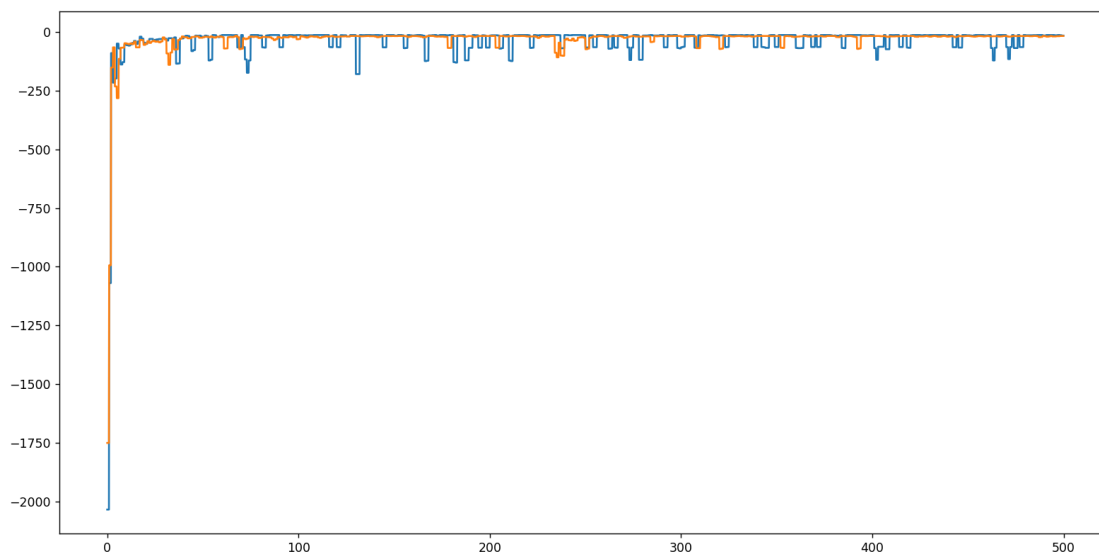
运行程序，得到分别使用Q-learning和sarsa算法的实验结果如下：

路径图：（绿色为可走区域，黄色为路径，紫色为悬崖）：

Age Group	Male (%)	Female (%)	Other (%)
18-24	100	0	0
25-34	100	0	0
35-44	100	0	0

[illegible]

在每次迭代的过程中，reward的总和的变化图如下，其中橙色的曲线为sarsa算法的reward总和变化，蓝色的曲线为q-learning算法的reward总和变化：



我们可以看到在on-line训练中，Q-learning的表现要比Sarsa差，但q-learning最后得到的结果比sarsa要好。并且通过分析可以看出两个算法的区别，Sarsa是一种on-policy算法，Q-learning是一种off-policy算法。Sarsa选取的是一种保守的策略，在更新Q值的时候已经为未来规划好了动作，对错误和死亡比较敏感。而Q-learning每次在更新的时候选取的是最大化Q的方向，而当下一个状态时，再重新选择动作，Q-learning是一种鲁莽、大胆、贪婪的算法，对于死亡和错误并不在乎。

---

## 5、实验心得

通过这次实验，我基于Cliff Walk例子实现了Q-learning算法和Sarsa算法，及时复习了课上所学的内容，并对两种算法有了更加深刻的理解，也对这两种算法的区别和各自的优点有了更深的体会。除此之外，通过将所学知识应用于实际问题中，我的代码编写能力和实践能力都有所提示，这让我受益匪浅。