

Reinforcement Learning

Task3

姚冠宇

2020211260070

先展示最终结果（程序可以直接运行）：

Sarsa和Q-learning各自迭代了500轮

Sara

最优路线：

```
['^^', 'VV', '<<', '>>']
```

分别对应上、下、左、右

```
Sarsa算法最终收敛得到的路径为：
00 00 00 00 00 00 00 00 00 00 00 00
>> >> >> >> >> >> >> >> >> >> VV
^^ 00 00 00 00 00 00 00 00 00 00 00 VV
^^ ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** EE
```

最优路线的步数：

```
Sarsa算法最终收敛得到的步数为： 20
```

最终得到每个状态的动作表格（一个位置可能存在多个策略故采用这种方式打印）：

```
Sarsa算法最终收敛得到的策略为：
000> ^000 000> 000> 000> 000> 000> 0V00 000> 000> 000> 0V00
000> 000> 000> 000> 000> 000> 000> 000> 000> 000> 0V00
^000 ^000 ^000 000> ^000 ^000 000> 00<0 ^000 000> 000> 0V00
^000 ***** ***** ***** ***** ***** ***** ***** ***** EEEE
```

Q-learning;

最优路线:

```
['^^', 'VV', '<<', '>>']
```

分别对应上、下、左、右

```
Q-learning算法最终收敛得到的路径为:  
00 00 00 00 00 00 00 00 00 00 00 00  
00 00 00 00 00 00 00 00 00 00 00 00  
>> >> >> >> >> >> >> >> >> >> VV  
^^ ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** ** EE
```

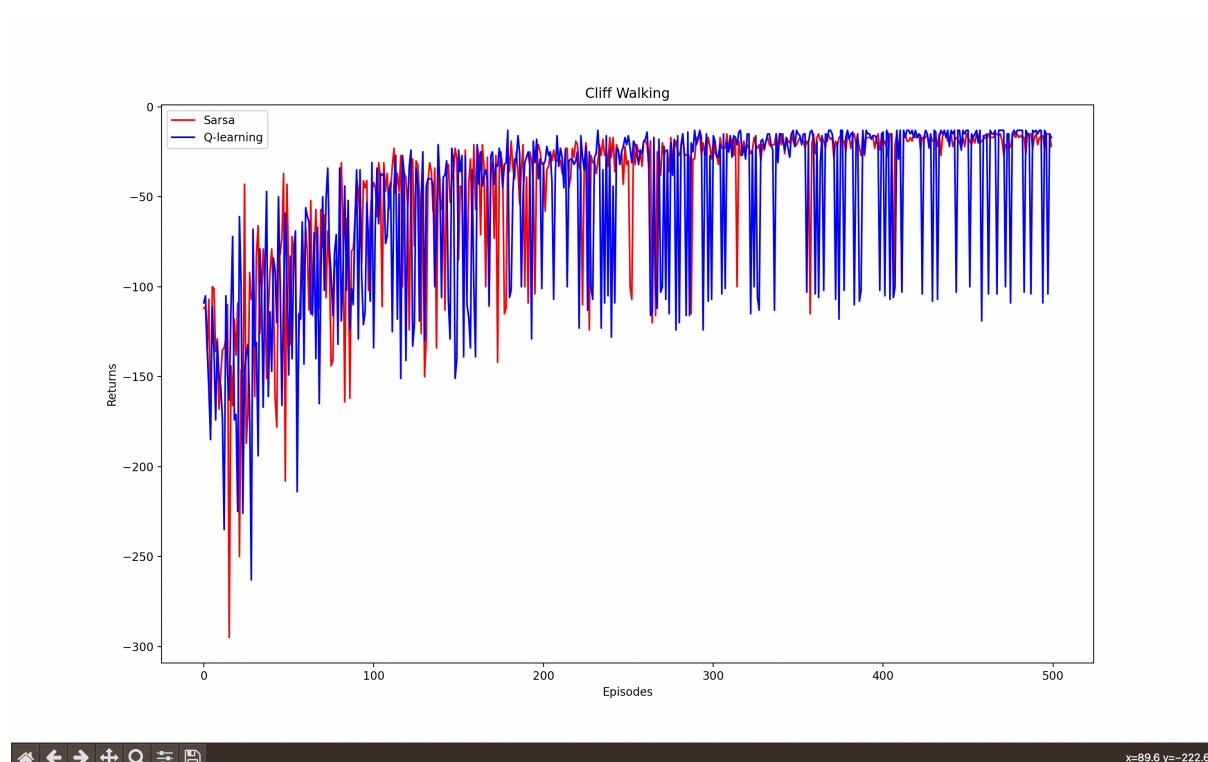
最优路线步数:

```
Q-learning算法最终收敛得到的步数为: 13
```

最终得到每个状态的动作表格（一个位置可能存在多个策略故采用这种方式打印）:

```
00<0 00<0 00<0 0V00 000> 000> 000> ^000 000> 0V00 0V00 0V00  
^000 000> 000> ^000 000> 000> 000> 0V00 000> 000> 000> 0V00  
000> 000> 000> 000> 000> 000> 000> 000> 000> 000> 0V00  
^000 ***** ***** ***** ***** ***** ***** ***** ***** ***** EEEE
```

两种算法每个epoch的回报



问答题

$\pi_*(Q)$ 和 $\pi_*(S)$ 是否相同?

答: 不同, 我们可以用价值迭代的思想来理解 Q-learning, 即 Q-learning 是直接在估计 Q^* , 因为动作价值函数的贝尔曼最优方程是

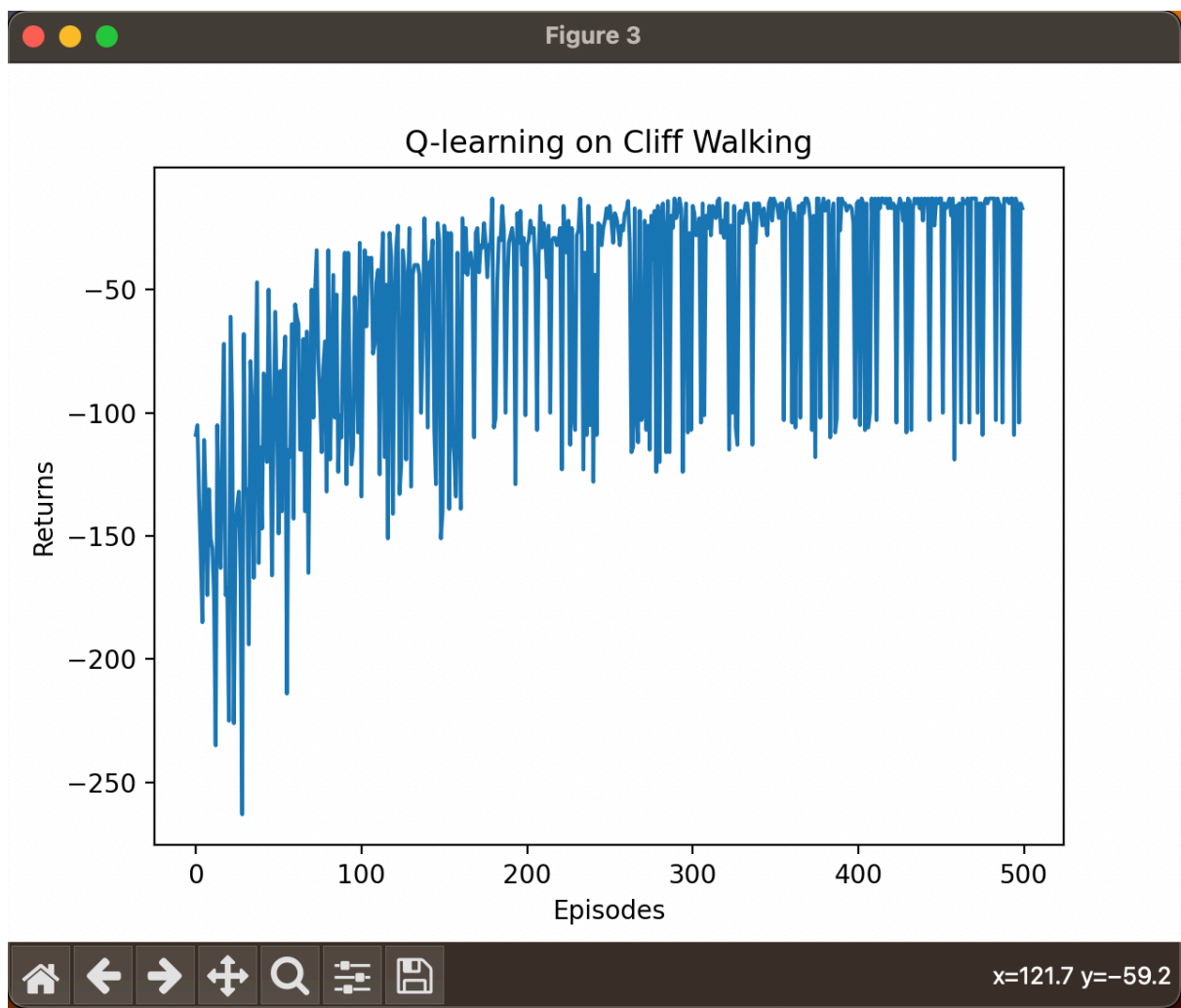
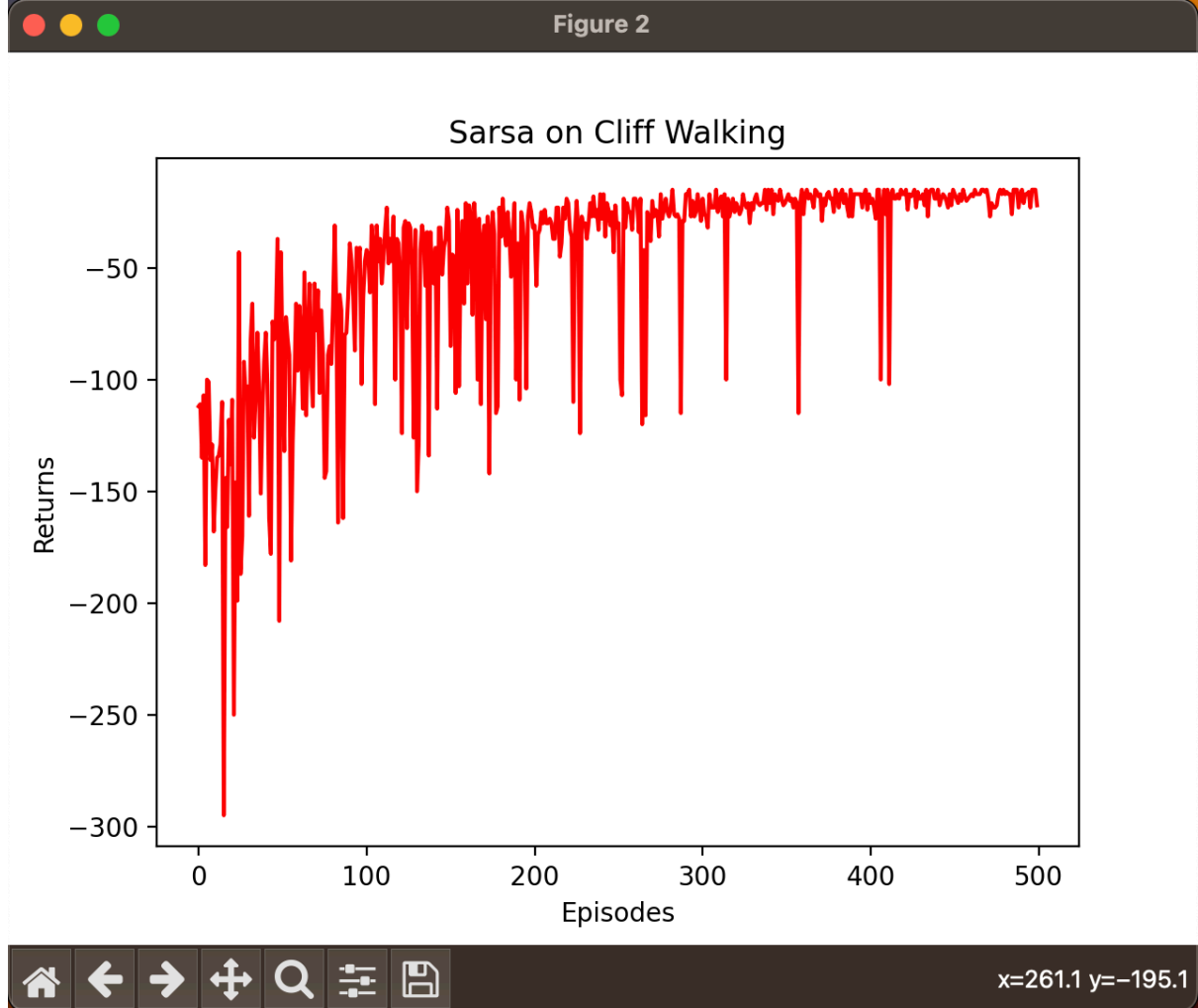
$$Q^*(s, a) = r(s, a) + \gamma \sum_{s' \in S} P(s' | s, a) \max_{a'} Q^*(s', a')$$

而Sarsa估计的是 ϵ -贪婪策略的价值函数。这也导致Q-learning的结果一般比Sarsa算法要好。

我们把目标策略的行为打印出来后, 发现Q-learning更偏向于走在悬崖边上, 这与 Sarsa 算法得到的比较保守的策略相比是更优的。但是仔细观察 Sarsa 和 Q-learning 在训练过程中的回报曲线图, 我们可以发现, 在一个序列中 Sarsa 获得的期望回报是高于 Q-learning 的。这是因为在训练过程中智能体采取基于当前函数的-贪婪策略来平衡探索与利用, Q-learning 算法由于沿着悬崖边走, 会以一定概率探索“掉入悬崖”这一动作, 而 Sarsa 相对保守的路线使智能体几乎不可能掉入悬崖。

Sarsa 和 Q-learning 在训练过程中的回报曲线图在上一页已经附上, 文件中也附上了清晰的大图。

同时可以注意到, Q-learning 的更新并非必须使用当前贪心策略 $\arg\max_a Q(s, a)$ 得到的数据, 因为给定任意 (s, a, r, s') 都可以直接根据更新公式来更新Q。Sarsa 必须使用当前 ϵ -贪婪策略采样得到的数据, 因为它的更新中用到的 $Q(s', a')$ 的 a' 是当前策略在 s' 下的动作。因此, Sarsa是在线策略, Q-learning是离线策略。



大致讲解一下代码结构，代码里有详细的注释

CliffWalkingEnv主要是环境，输入智能体的动作给出下一个状态，**reward**和是否到达终态的**done**

Sarsa类和**Q-learning**是智能体，主要维护一个nrow*ncol行，n_action列的**Q-table**

我们初始化学习率 $\alpha=0.1$ ，折扣因子 $\gamma=1$ ，**epsilon-贪婪**中的 $\epsilon=0.1$ 。每一步选择动作价值最大的动作或随机选择动作。同时定义**best action**函数用来打印策略。定义**update**方法实现**sarsa**和**Q-learning**的**q_table**更新。

主函数进行了num_episodes=500轮次训练，每10条序列打印一下这10条序列的平均回报。训练过程中需要注意的是，**Sarsa**的输入为(s,a,r,s',a')，**Q-learning**的输入为(s,a,r,s')。因此不能直接套用**Sarsa**的训练代码用于**Q-learning**，下面讲到的**print_route()**函数亦是如此。

Sarsa算法的更新公式为

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

Q-learning算法的更新公式为：

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[R_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$$

这是Sara和Q-learning最大的不同

程序中也根据这个公式来更新q_table.

最后定义了两个打印函数print_agent()和print_route()。第一个用于打印智能

体在每个位置会选择动作（即打印策略），第二个用于打印智能体从开始位置出发会走的路线（即打印路径）。