

# Assignment 3

Peng Tang 517020910038

In experiment, we are going to find the best solution for Cliff Walking problem by using on-policy learning (Sarsa) and off-policy learning(Q-learning) and compare their performance with different parameters.

## 1 Analysis

---

### 1.1 Sarsa

---

Sarsa is an on-policy learning, which means we need to use same policy to take actions and update Q-values. We can implement Sarsa for this problem as following steps.

- We set the start state and take actions by using  $\epsilon$ -greedy policy(shown as following equation) until reaching the terminal state.

$$\pi(a|s) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{4} & a = \arg \max_a Q(s, a) \\ \frac{\epsilon}{4} & a \neq \arg \max_a Q(s, a) \end{cases}$$

- In this process, we use the  $\epsilon$ -greedy policy to update the Q-values by using the following equation.

$$Q(s, a) \leftarrow Q(s, a) + \alpha(R + \lambda Q(s', a') - Q(s, a))$$

where  $a'$  is taken by using  $\epsilon$ -greedy;  $R$  is  $-100$  if stepping into the region marked "The Cliff", otherwise  $-1$ .

- We do above two steps for enough times until all Q-values converge.
- Then we can find the best path from the start state to the goal state by using greedy policy.

### 1.2 Q-learning

---

Q-learning is an off-policy learning, which means we need to use different policies to take actions and update Q-values. We can implement Q-learning for this problem as following steps.

- We set the start state and take actions by using  $\epsilon$ -greedy policy(shown as following equation) until reaching the terminal state.

$$\pi(a|s) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{4} & a = \arg \max_a Q(s, a) \\ \frac{\epsilon}{4} & a \neq \arg \max_a Q(s, a) \end{cases}$$

- In this process, we use the greedy policy to update the Q-values by using the following equation.

$$Q(s, a) \leftarrow Q(s, a) + \alpha(R + \lambda \max_{a'} Q(s', a') - Q(s, a))$$

where  $R$  is  $-100$  if stepping into the region marked "The Cliff", otherwise  $-1$ .

- We do above two steps for enough times until all Q-values converge.
- Then we can find the best path from the start state to the goal state by using greedy policy.

## 2 Code

### 2.1 Sarsa Class

```

1  class Sarsa(object):
2      def __init__(self, discount, maxIteration, epsilon, alpha) -> None:
3          self.discount = discount
4          self.maxIteration = maxIteration
5          self.epsilon = epsilon
6          self.alpha = alpha
7
8          self.Q_value = np.random.random((4,12,4))
9          self.Q_value[3][11][0] = self.Q_value[3][11][1]=self.Q_value[3][11]
[2]=self.Q_value[3][11][3]=0
10         self.terminal = [3,11]
11
12         self.actions = [[-1,0],[1,0],[0,-1],[0,1]]
13         self.cliffstate = [[3,1],[3,2],[3,3],[3,4],[3,5],[3,6],[3,7],[3,8],
[3,9],[3,10]]
14
15     def findOptimalPath(self):
16
17         iteration = 0
18         cost = []
19         while iteration < self.maxIteration:
20             self.alpha = 0.99*self.alpha
21             state = [3,0]
22             nextstate = [0,0]
23
24             tempcost = 0
25             action = self.getaction(state)
26             while state !=self.terminal:
27                 nextstate[0] = max(0,min(3,state[0]+self.actions[action]
[0]))
28                 nextstate[1] = max(0,min(11,state[1]+self.actions[action]
[1]))
29
30                 reward = -1
31                 if nextstate in self.cliffstate:
32                     reward = -100
33                     nextstate[0] = 3
34                     nextstate[1] = 0
35
36                 tempcost += reward
37                 nextaction = self.getaction(nextstate)
38
39                 self.Q_value[state[0]][state[1]][action] = (1-
self.alpha)*self.Q_value[state[0]][state[1]][action] + self.alpha*
(reward+self.discount*self.Q_value[nextstate[0]][nextstate[1]][nextaction])
40
41                 state[0] = nextstate[0]
42                 state[1] = nextstate[1]
43                 action = nextaction

```

```

44         iteration += 1
45         cost.append(tempcost)
46
47         state = [3,0]
48         actionlist = []
49         realaction = ['up', 'down', 'left', 'right']
50         while state!=self.terminal:
51             action = np.argmax(self.Q_value[state[0]][state[1]])
52             actionlist.append(realaction[action])
53             state[0] = max(0,min(3,state[0]+self.actions[action][0]))
54             state[1] = max(0,min(11,state[1]+self.actions[action][1]))
55
56             if state in self.cliffstate:
57                 state[0] = 3
58                 state[1] = 0
59
60         return actionlist, cost
61
62     def getaction(self,state):
63         greedyAction = np.argmax(self.Q_value[state[0]][state[1]])
64         p = np.random.random()
65         if p < 1-self.epsilon:
66             return greedyAction
67         else:
68             return np.random.randint(4)

```

- We use this class to implement Sarsa.
- In this class, **discount** is the discount factor( $\lambda$  in the equation); **maxIteration** is the number of episodes we sample; **alpha** is the step size( $\alpha$  in the equation) and we decrease it every episode by multiplying 0.99; **epsilon** is the  $\epsilon$  in the  $\epsilon$ -greedy policy; **Q\_value** is used to record the Q-values and we initialize it randomly; **terminal** is the goal; **actions** includes four actions, up, down, left and right; **cliffstate** includes all states marked "The Cliff".
- The **findOptimalPath** function is used to sample and update Q-values and find the optimal path from the start state to the goal state at the end. We implement this function as our analysis and we use the **cost** to record the total rewards of every episode. And the function will return the actions we need to take from the start state and the costs of all episodes.
- The **getaction** function is used to choose a action by using the  $\epsilon$ -greedy policy.

## 2.2 Q\_learning Class

```

1 class Q_learning(object):
2     def __init__(self, discount, maxIteration, epsilon, alpha) -> None:
3         self.discount = discount
4         self.maxIteration = maxIteration
5         self.epsilon = epsilon
6         self.alpha = alpha
7
8         self.Q_value = np.random.random((4,12,4))
9         self.Q_value[3][11][0] = self.Q_value[3][11][1]=self.Q_value[3][11]
10        [2]=self.Q_value[3][11][3]=0
11        self.terminal = [3,11]
12
13        self.actions = [[-1,0],[1,0],[0,-1],[0,1]]

```

```

13         self.cliffstate = [[3,1],[3,2],[3,3],[3,4],[3,5],[3,6],[3,7],[3,8],
14                             [3,9],[3,10]]
15
16     def findOptimalPath(self):
17
18         iteration = 0
19         cost = []
20         while iteration < self.maxIteration:
21             self.alpha = 0.99*self.alpha
22             state = [3,0]
23             nextstate = [0,0]
24             tempcost = 0
25             while state !=self.terminal:
26                 action = self.getaction(state)
27                 nextstate[0] = max(0,min(3,state[0]+self.actions[action]
28                                     [0]))
29                 nextstate[1] = max(0,min(11,state[1]+self.actions[action]
30                                     [1]))
31
32                 reward = -1
33                 if nextstate in self.cliffstate:
34                     reward = -100
35                     nextstate[0] = 3
36                     nextstate[1] = 0
37
38                 tempcost += reward
39                 self.Q_value[state[0]][state[1]][action] = (1-
40 self.alpha)*self.Q_value[state[0]][state[1]][action] + self.alpha*
41 (reward+self.discount*np.max(self.Q_value[nextstate[0]][nextstate[1]]))
42
43                 state[0] = nextstate[0]
44                 state[1] = nextstate[1]
45
46             iteration += 1
47             cost.append(tempcost)
48
49         state = [3,0]
50         realaction = ['up','down','left','right']
51         actionlist = []
52         while state!=self.terminal:
53             action = np.argmax(self.Q_value[state[0]][state[1]])
54             actionlist.append(realaction[action])
55             state[0] = max(0,min(3,state[0]+self.actions[action][0]))
56             state[1] = max(0,min(11,state[1]+self.actions[action][1]))
57
58             if state in self.cliffstate:
59                 state[0] = 3
60                 state[1] = 0
61
62         return actionlist, cost
63
64     def getaction(self,state):
65         greedyAction = np.argmax(self.Q_value[state[0]][state[1]])
66         p = np.random.random()
67         if p < 1-self.epsilon:
68             return greedyAction
69         else:
70             return np.random.randint(4)

```

This class is similar to the Sarsa and the difference is the way to update the Q-values in **findOptimalPath** function. We use greedy policy rather than  $\epsilon$ -greedy in this class.

## 2.3 Main Function

```
1  def main():
2
3      np.random.seed(0)
4      discount = 0.9
5      maxIteration = 700
6      epsilon = [0,0.05,0.1,0.15,0.2,0.25]
7      alpha = 0.5
8      SarsaRewardList = []
9      Q_learningRewardList = []
10     CostTime = []
11     Sarsa_ActionList = []
12     Q_learningActionList = []
13
14
15     for i in range(len(epsilon)):
16         t1 = time.time()
17         SarsaModel = Sarsa(discount=discount,
18 maxIteration=maxIteration,epsilon=epsilon[i],alpha=alpha)
19         Sarsa_Action, Sarsa_Cost = SarsaModel.findOptimalPath()
20         Q_learningModel = Q_learning(discount=discount,
21 maxIteration=maxIteration,epsilon=epsilon[i],alpha=alpha)
22         Q_learning_Action, Q_learning_Cost = Q_learningModel.findOptimalPath()
23         t2 = time.time()
24
25         CostTime.append(t2-t1)
26         Sarsa_ActionList.append(Sarsa_Action)
27         Q_learningActionList.append(Q_learning_Action)
28         SarsaRewardList.append(Sarsa_Cost)
29         Q_learningRewardList.append(Q_learning_Cost)
30
31
32     subfig = [231,232,233,234,235,236]
33     plt.figure(figsize=(20,10),dpi=90)
34     for i in range(len(epsilon)):
35         plt.subplot(subfig[i])
36         plt.plot(list(range(maxIteration)),SarsaRewardList[i],label="Sarsa")
37
38         plt.plot(list(range(maxIteration)),Q_learningRewardList[i],label="Q_learning")
39
40         plt.xlabel("Episodes")
41         plt.ylabel("Path Cost")
42         plt.legend()
43         plt.title("epsilon="+str(epsilon[i]))
44     plt.savefig('imgs/costChanges.png')
45     plt.close()
46
47     plt.plot(epsilon, CostTime)
48     plt.xlabel("Different Epsilons")
```

```

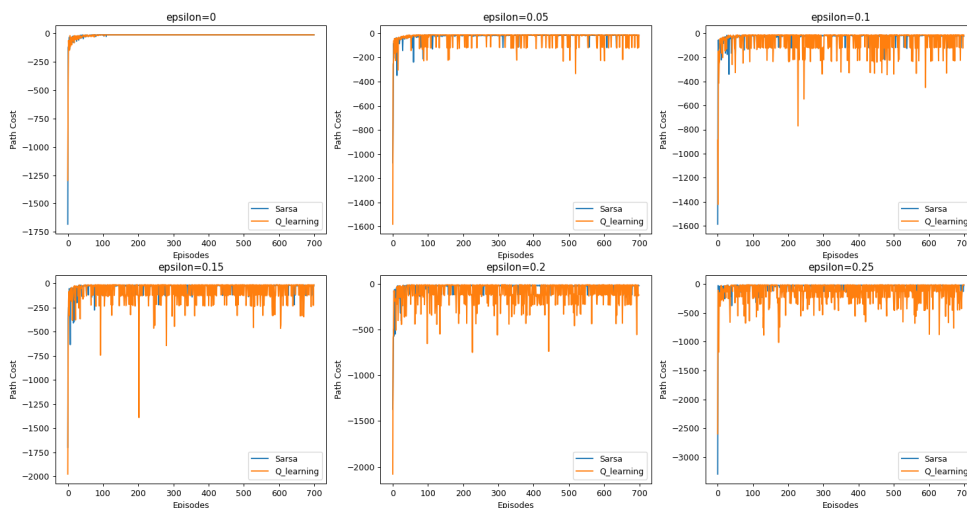
48 plt.ylabel("Cost Time(s)")
49 plt.savefig('imgs/CostTime.png')
50 plt.close()
51
52
53 print("Sarsa actions:")
54 for i in range(len(epsilon)):
55     print("epsilon="+str(epsilon[i])+" : ",Sarsa_ActionList[i])
56 print()
57 print("Q_learning actions:")
58 for i in range(len(epsilon)):
59     print("epsilon="+str(epsilon[i])+" : ",Q_learningActionList[i])

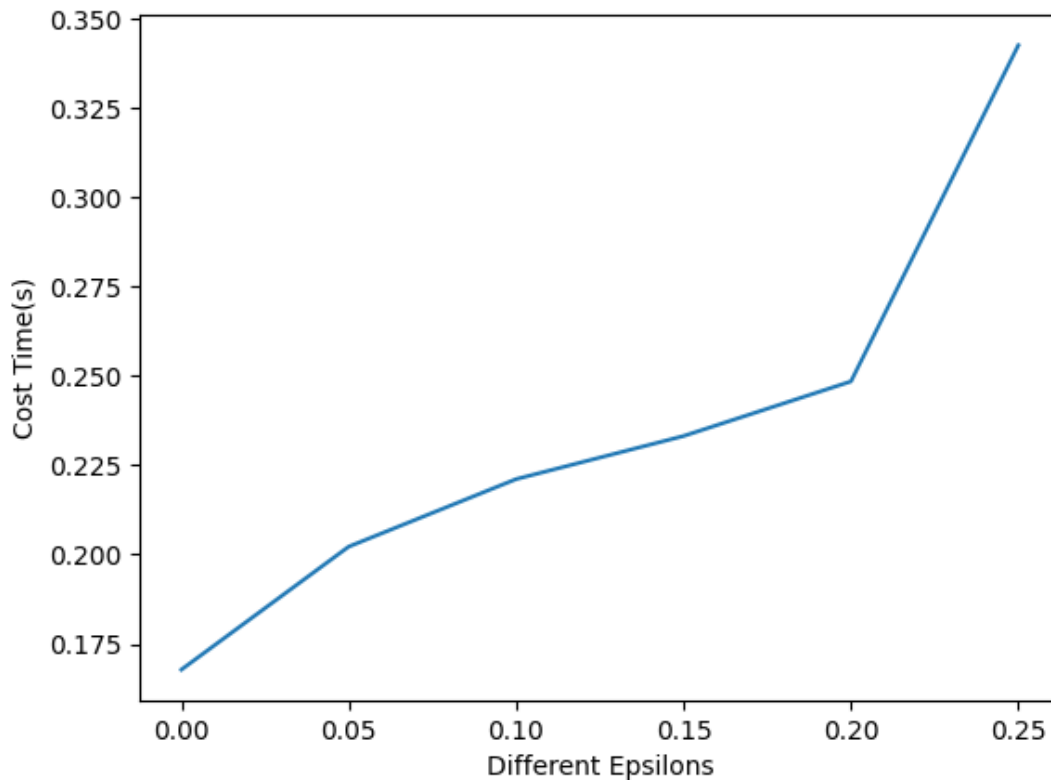
```

This function is used to call Sarsa and Q\_learning and output the results. In order to analyze the influence of different  $\epsilon$  values, we set six different values(0, 0.05, 0.1, 0.15, 0.2, 0.25) and compare their performance.

## 3 Results

By calling the **main()** function, we get the following results.





```

PS E:\大三下\强化学习\作业\homework3\code> D:\python\python.exe "E:\大三下\强化学习\作业\homework3\code\solution.py"
Sarsa actions:
epsilon=0: ['up', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'down']
epsilon=0.05: ['up', 'up', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'down', 'down', 'down']
epsilon=0.1: ['up', 'up', 'up', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'down', 'down', 'down']
epsilon=0.15: ['up', 'up', 'up', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'down', 'down', 'down']
epsilon=0.2: ['up', 'up', 'up', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'down', 'down', 'down']
epsilon=0.25: ['up', 'up', 'up', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'down', 'down', 'down']

Q_learning actions:
epsilon=0: ['up', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'down']
epsilon=0.05: ['up', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'down']
epsilon=0.1: ['up', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'down']
epsilon=0.15: ['up', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'down']
epsilon=0.2: ['up', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'down']
epsilon=0.25: ['up', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'right', 'down']
PS E:\大三下\强化学习\作业\homework3\code>

```

- From the first figure, we can see that as the  $\epsilon$  becomes bigger, the fluctuation of curves is bigger. This is because the  $\epsilon$  is bigger, we explore more randomly. What more, if  $\epsilon$  is 0, the Sarsa is same as Q-learning according our algorithm. And if  $\epsilon$  is not 0, we can see that the fluctuation of Q-learning is bigger than the Sarsa. I think it is because that for the Q-learning, the Q-values of the states near the "the Cliff" states is bigger, then the agent will step into these states more often and then it will step into the "the Cliff" states easier by taking actions with  $\epsilon$ -greedy.
- From the second figure, we can see that as the  $\epsilon$  is bigger, we need more time to sample the fixed number of episodes. The is obvious because  $\epsilon$  is bigger, we need explore more. When  $\epsilon$  is too big(maybe bigger than 0.4), we need very long time to sample.
- From the third figure, we can see that with different  $\epsilon$  values, Sarsa will give different actions while Q-learning will give same action.
  - For Sarsa, because it uses  $\epsilon$ -greedy policy to update the Q-values, then as the  $\epsilon$  is bigger, the Q-values of the states near "The Cliff" region will become smaller than usual. Thus, when we take actions by using greedy policy, we will stay away from "The Cilff" region.
  - For Q-learning, because it uses greedy policy to update the Q-values, then the Q-values of the states near the goal states will become bigger although these state is also near "The Cliff" states in this problem.

- From this problem, we can see that when  $\epsilon$  is 0, we can get the best performance. However, this is because the state space and action space are smaller, we needn't to explore. When these two spaces become bigger, we need to explore more to find the best solution.