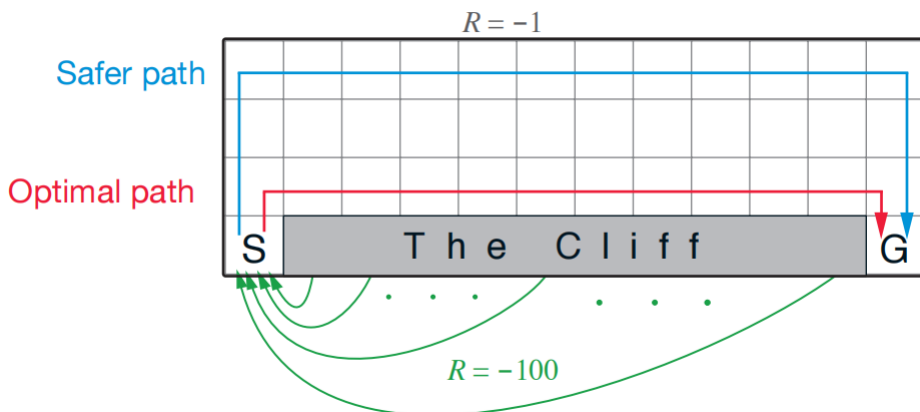


第六章作业

题目 1. 例 6.6、练习 6.12：在悬崖边上行走

如右图所示是一个 4×12 的网格图，起点在左下角，终点在右下角，最下面一行除去左右端点外均为悬崖。每一幕开始时，智能体从起点开始移动，每一步有上下左右四个方向可以选择，如果到达悬崖或达到终点，结束当前幕。当



智能体掉下悬崖时收益为 -100 ，其他每步收益均为 -1 ，使用 sarsa 和 Q Learning 求解该问题，比较在不同 ε 下，sarsa 和 Q Learning 的策略有何不同。

解答. 设置步长 $\alpha = 0.9$ ，折扣率为 $\gamma = 1$ ，分别测试 $\varepsilon = 0.1, 0.01, 0$ 下 sarsa 和 Q Learning 在每一幕下收益之和的变换情况，其中 $\varepsilon = 0$ 即使用贪心方式选择策略。

通过平均 1000 次试验结果，得到以下收益之和变换情况（如图 2 所示），并输出不同 ε 下的最优策略（如图 1 所示），可以发现，在 $\varepsilon = 0$ 时，两者使用了相同的最优贪心策略，并且收益之和曲线几乎完全重合，即在悬崖边上行走；但当 ε 增大时，Q Learning 仍采用最优策略，但是 sarsa 逐渐转换为安全策略，靠远离悬崖的路线上行走。

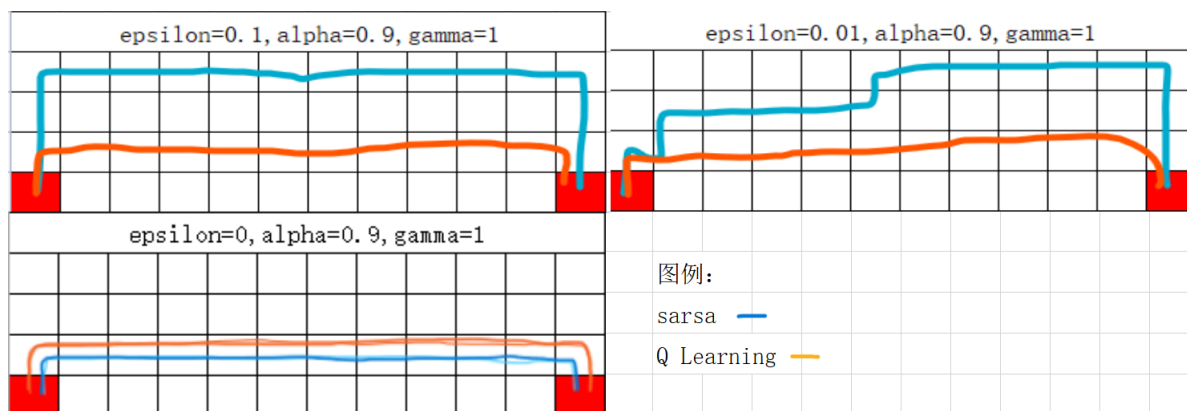


图 1: 不同 ε 下 sarsa 和 Q Learning 策略比较图

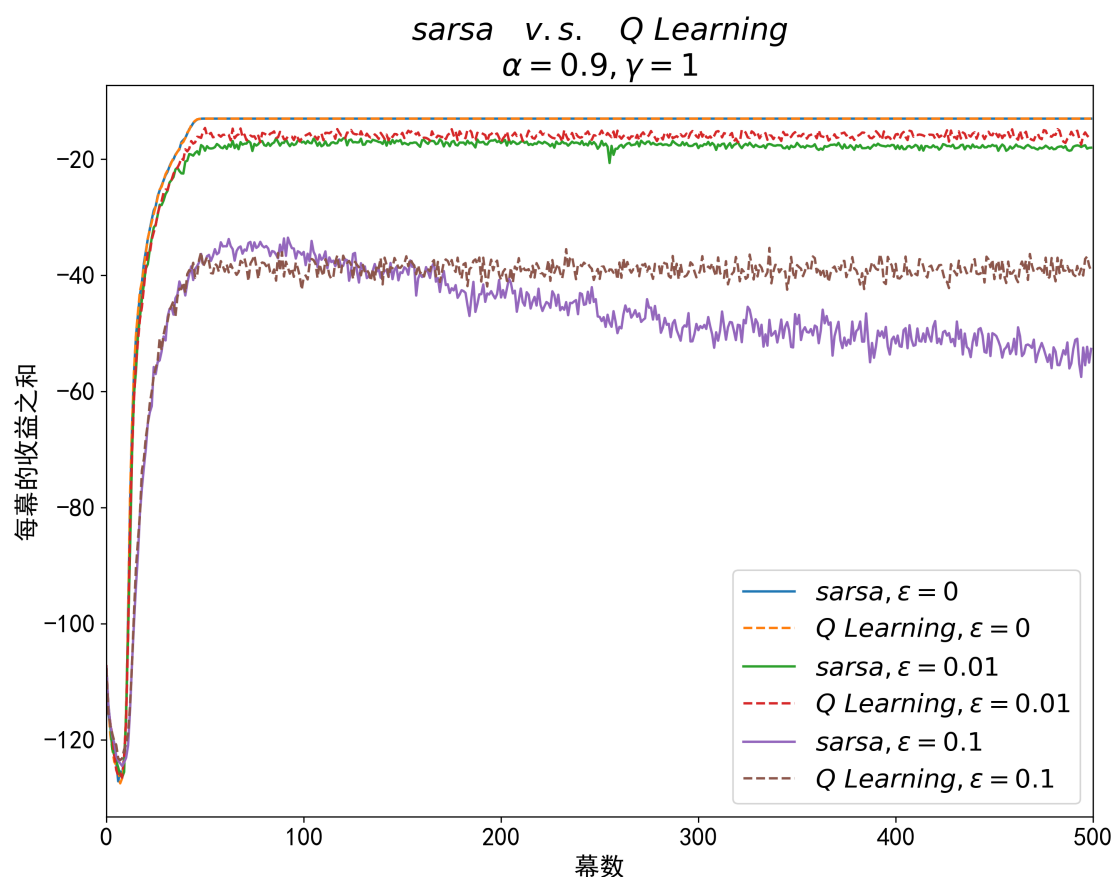


图 2: 不同 ε 下策略收益变换情况

源代码:

```

1  # -*- coding: utf-8 -*-
2  '''
3  @File      : main.py
4  @Time      : 2023/03/21 19:43:45
5  @Author    : wty-yy
6  @Version   : 1.0
7  @Blog      : https://wty-yy.space/
8  @Desc      : 《强化学习》中文书第 130 页例 6.6 复现 ( $\epsilon=0.1$ ), 并完成练习
   ↪ 6.12( $\epsilon=0$ )
9  '''
10
11 from tqdm import tqdm
12 import numpy as np
13 import matplotlib as mpl
14 import matplotlib.pyplot as plt
15
16 config = { # matplotlib 绘图配置
17     "figure.figsize": (6, 6), # 图像大小
18     "font.size": 16, # 字号大小
19     "font.sans-serif": ['SimHei'], # 用黑体显示中文
20     'axes.unicode_minus': False # 显示负号
21 }

```

```

22 plt.rcParams.update(config)
23
24 ACTION_TO_DIRECTION = [np.array(x) for x in ((0, 1), (0, -1), (1, 0), (-1, 0))]
25
26 def check_available_position(position):
27     return 0 <= position[0] < 4 and 0 <= position[1] < 12
28
29 def position_action_iterator():
30     for x in range(4):
31         for y in range(12):
32             position = np.array((x, y))
33             for action in range(4):
34                 yield position, action
35
36 def plot_savefig(title: str, fname: str):
37     plt.xlabel(" 幕数")
38     plt.ylabel(" 每幕的收益之和")
39     plt.title(title)
40     plt.legend()
41     plt.xlim(0, 500)
42     plt.tight_layout()
43     plt.savefig(fname, dpi=300)
44     plt.show()
45
46 class WalkingEnvironment:
47     def __init__(self) -> None:
48         self.position = np.array(())
49
50     def step(self, action: int):
51         direction = ACTION_TO_DIRECTION[action]
52         position_ = self.position + direction
53         assert(check_available_position(position_))
54         reward = -100 if position_[0] == 3 and 1 <= position_[1] < 11 else -1
55         terminal = True if (position_[0] == 3 and position_[1] == 11) or (reward
56             ↪ == -100) else False
57         # self.position = self.reset() if reward == -100 else position_
58         self.position = position_
59         return reward, self.position, terminal
60
61     def reset(self):
62         self.position = np.array((3, 0))
63         return self.position
64
65 class Agent:
66     alpha = None
67     def __init__(self, epsilon=0.1, gamma=1, seed=42, episode=500,
68         ↪ average_times=1000, is_plot=True) -> None:
69         self.epsilon = epsilon
70         self.gamma = gamma
71         np.random.seed(seed)
72         self.episode = episode
73         self.average_times = average_times
74         self.is_plot = is_plot

```

```

73
74     self.history = np.zeros(self.episode)
75     self.available_actions = np.zeros((4, 12, 4))
76     for position, action in position_action_iterator():
77         position_ = position + ACTION_TO_DIRECTION[action]
78         if check_available_position(position_):
79             self.available_actions[position[0], position[1], action] = 1
80
81     def start(self):
82         # alphas = [0.1, 0.5, 0.9]
83         alphas = [0.9]
84         for self.alpha in alphas:
85             for _ in tqdm(range(self.average_times)):
86                 self.history += (self.solve(policy='sarsa') - self.history) /
87                     ↪ (_+1)
88                 self.solve(policy='sarsa', verbose=True)
89                 self.plot_rewards('sarsa')
90                 self.history = np.zeros(self.episode)
91
92             for _ in tqdm(range(self.average_times)):
93                 self.history += (self.solve(policy='q_learning') - self.history)
94                     ↪ / (_+1)
95                 self.plot_rewards('Q\ Learning')
96                 self.solve(policy='q_learning', verbose=True)
97                 self.history = np.zeros(self.episode)
98
99         if self.is_plot:
100             plot_savefig(f"$\\epsilon={self.epsilon}$",
101                 ↪ f"$\\epsilon={self.epsilon}, alphas={alphas}, gamma={self.gamma}, av
102
103     def solve(self, policy, verbose=False):
104         history = []
105         policies = ['sarsa', 'q_learning']
106         assert(policy in policies)
107
108         q = np.random.normal(loc=0, size=(4, 12, 4))
109         for action in range(4):
110             q[3, 11, action] = 0
111         for position, action in position_action_iterator():
112             x, y = position
113             if not self.available_actions[x, y, action]:
114                 q[x, y, action] = -np.inf
115
116     def epsilon_choose(state):
117         actions = np.argwhere(self.available_actions[state[0], state[1]] >
118             ↪ 0).reshape(-1)
119         randnum = np.random.rand(1)[0]
120         if randnum <= self.epsilon:
121             random_action = np.random.randint(0, actions.shape[0])
122             return actions[random_action]
123         else:
124             return np.argmax(q[state[0], state[1]])

```

```

122
123     for _ in range(self.episode):
124         env = WalkingEnvironment()
125         state = env.reset()
126         terminal = False
127         action = epsilon_choose(state) if policy == 'sarsa' else None
128         total_reward = 0
129         total_step = 0
130         while not terminal:
131             total_step += 1
132             action = action if policy == 'sarsa' else epsilon_choose(state)
133             reward, state_, terminal = env.step(action)
134
135             if policy == 'sarsa':
136                 action_ = epsilon_choose(state_)
137             else:
138                 action_ = np.argmax(q[state_[0], state_[1]])
139             error = reward + self.gamma * q[state_[0], state_[1], action_] -
140                 ↪ q[state[0], state[1], action]
141             q[state[0], state[1], action] += self.alpha * error
142
143             state, action = state_, action_ if policy == 'sarsa' else action
144             total_reward = reward + total_reward
145             history.append(total_reward)
146
147         if verbose:
148             print(f"{policy}'s best policy:")
149             state = env.reset()
150             step_history = [tuple(state)]
151             terminal = False
152             while not terminal:
153                 action = np.argmax(q[state[0], state[1]])
154                 reward, state, terminal = env.step(action)
155                 step_history.append(tuple(state))
156             print(step_history)
157         return np.array(history)
158
159     def plot_rewards(self, name):
160         # plt.plot(self.history, label=f'${name}', \epsilon={self.epsilon},
161         ↪ \alpha={self.alpha:.1f}, \gamma={self.gamma}$ ')
162         plt.plot(self.history, label=f'${name}', \epsilon={self.epsilon}$',
163         ↪ ls='--' if name == 'sarsa' else '-')
164
165 if __name__ == '__main__':
166     plt.figure(figsize=(10, 8))
167     agent = Agent(is_plot=False)
168     epsilons = [0, 0.01, 0.1]
169     # agent.average_times = 10
170     for epsilon in epsilons:
171         agent.epsilon = epsilon
172         agent.start()
173     plot_savefig(f"$sarsa\\quad v.s.\\quad Q\\
174     ↪ Learning$\n$\\alpha={agent.alpha:.1f}, \\gamma={agent.gamma}$",

```

```
171                                     ↪ f"epsilon={epsilons}, alphas={agent.alpha}, gamma={agent.gamma}, average=  
172  
173 if __name__ == '__main__':  
174     x, y, z = np.array([-np.inf, -np.inf, -1e9])  
175     print(x, y, z)
```
