

Giorgio Mendoza

RBE595-S24-S04

Temporal Difference Programming Exercise

## ✓ Objective:

This code implements SARSA and Q-learning algorithms in the Cliff Walking task, a grid-based environment where an agent aims to reach a goal from a start point while avoiding a "cliff" that triggers large negative rewards. The environment setup defines the grid, start and goal positions, cliff locations, and a step function to update the agent's state and reward based on its actions. SARSA, an on-policy method, and Q-learning, an off-policy method, are compared through simulation, with their learning progress visualized in a plot of cumulative rewards per episode.

## SARSA (State-Action-Reward-State-Action):

An on-policy algorithm that updates its Q-values based on the action actually taken according to the current policy, which includes exploration. This approach makes SARSA more cautious, as it learns from the actual outcomes of its actions.

## Q-learning:

An off-policy algorithm that updates its Q-values based on the maximum potential reward from the next state, independent of the agent's current exploration strategy. This can lead to more aggressive exploration and occasionally higher negative rewards due to riskier paths taken.

## Epsilon-Greedy Policy:

A common exploration strategy where, with probability  $\epsilon$ , the agent chooses a random action, and with probability  $1-\epsilon$ , it chooses the best-known action. This policy is used by both algorithms to balance exploration and exploitation.

## Training Function:

The train function runs simulations of the specified algorithm over a series of episodes. It accumulates rewards received in each episode and updates the Q-values according to the algorithm's rules. The function outputs a list of total rewards per episode, which serves as a measure of the agent's learning progress.

## Plotting:

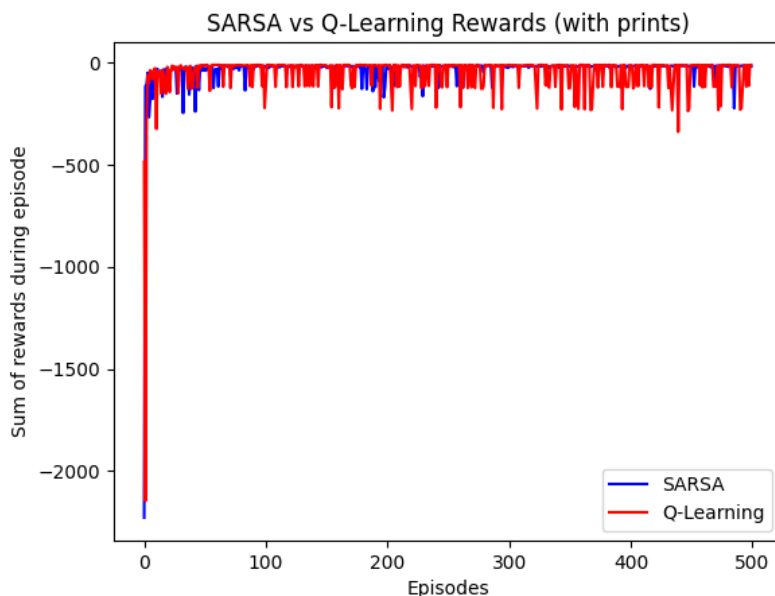
At the end, we generate a plot comparing the sum of rewards during each episode for both algorithms. This visual comparison helps in analyzing the efficiency and learning behavior of SARSA and Q-learning across episodes.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 #parameters
5 num_episodes = 500
6 epsilon = 0.1
7 alpha = 0.5
8 gamma = 1
9 actions = ['up', 'right', 'down', 'left']
10 width, height = 12, 4
11 start, goal = (3, 0), (3, 11)
12 cliff = [(3, i) for i in range(1, 11)]
13
14 #initialize Q-values
15 Q_sarsa = {state: {action: 0 for action in actions} for state in [(x, y) for x in range(height) for y in range(width)]}
16 Q_ql = {state: {action: 0 for action in actions} for state in [(x, y) for x in range(height) for y in range(width)]}
17
18 #epsilon-greedy policy
19 def epsilon_greedy(Q, state):
20     if np.random.rand() < epsilon:
21         return np.random.choice(actions)
22     else:
23         return max(Q[state], key=Q[state].get)
24
25 #step function for environment
26 def step(state, action):
27     next_state = tuple(np.array(state) + np.array({'up': (-1, 0), 'right': (0, 1), 'down': (1, 0), 'left': (0, -1)}[action]))
28     next_state = (max(0, min(next_state[0], height - 1)), max(0, min(next_state[1], width - 1)))
29     if next_state in cliff:
30         return start, -100, False
```

```

30     return next_state, -100, False
31     if next_state == goal:
32         return goal, -1, True
33     return next_state, -1, False
34
35 def train(Q, algorithm):
36     rewards = []
37     for episode in range(num_episodes):
38         total_reward = 0
39         state = start
40         if algorithm == 'sarsa':
41             action = epsilon_greedy(Q, state)
42             step_count = 0
43             while True:
44                 if algorithm == 'q_learning':
45                     action = epsilon_greedy(Q, state)
46                 next_state, reward, done = step(state, action)
47                 total_reward += reward
48                 if algorithm == 'sarsa':
49                     next_action = epsilon_greedy(Q, next_state)
50                     Q[state][action] += alpha * (reward + gamma * Q[next_state][next_action] - Q[state][action])
51                     action = next_action
52                 elif algorithm == 'q_learning':
53                     best_next_action = max(Q[next_state], key=Q[next_state].get)
54                     Q[state][action] += alpha * (reward + gamma * Q[next_state][best_next_action] - Q[state][action])
55                 state = next_state
56                 step_count += 1
57                 if done:
58                     break
59             rewards.append(total_reward)
60
61         #print statement for each episode
62         #print(f"Episode {episode + 1}/{num_episodes}, Steps: {step_count}, Total Reward: {total_reward}")
63
64     return rewards
65
66 rewards_sarsa_prints = train(Q_sarsa, 'sarsa')
67 rewards_ql_prints = train(Q_ql, 'q_learning')
68
69 plt.plot(rewards_sarsa_prints, label='SARSA', color='blue')
70 plt.plot(rewards_ql_prints, label='Q-Learning', color='red')
71 plt.xlabel('Episodes')
72 plt.ylabel('Sum of rewards during episode')
73 plt.title('SARSA vs Q-Learning Rewards (with prints)')
74 plt.legend()
75 plt.show()

```



## Results:

The plot above illustrates the learning progression of two reinforcement learning algorithms, SARSA and Q-learning, applied to the Cliff Walking task. The graph shows the sum of rewards per episode over 500 episodes.

Initially, both algorithms experience large negative rewards, likely due to the agent frequently falling off the cliff. As learning progresses, SARSA demonstrates a conservative approach, reflected by a gradual increase in rewards, indicating a strategy that avoids the cliff. Q-learning,

however, shows a more aggressive learning curve with occasional dips in rewards, suggesting it is finding a more optimal, albeit riskier, path that occasionally results in falling.

Over time, both algorithms improve, but Q-learning exhibits a more pronounced fluctuation in rewards due to its explorative strategy, while SARSA maintains a more consistent performance.