Giorgio Mendoza

RBE595-S24-S04

Monte-carlo Programming Excercise

Calculate the optimal action value function and optimal policy using the On-policy first visit Monte-carlo.

## ⌄ Setup:

First, we start by defining the environment's state space (0-5)and possible (-1 for left, 1 for right). States 0 and 5 are terminal states.

Initializations:

Q-values:

A table (dictionary) of Q-values is initialized for each state-action pair, which estimates the expected rewards for taking a certain action in a given state.

Returns:

A table to store the returns (cumulative rewards) for each state-action pair across episodes.

Policy:

An ε-soft policy is initialized, ensuring some level of randomness (exploration) in action selection but with a preference for the best-known action.

Transition and Reward Function:

The step function models the environment's dynamics, determining the next state and reward given a current state and action. It handles the stochastic transitions and assigns rewards for reaching the charging station (state 0) or finding a can (state 5).

Policy Execution and Learning:

The choose_action function selects an action based on the current policy's probability distribution for a given state. The generate_episode function simulates an entire episode from a start state until a terminal state is reached, using the choose_action and step functions.

The update_policy function calculates the returns for each state-action pair encountered in the episode, updates the Q-values, and improves the policy by making it greedier with respect to the updated Q-values.

Training Loop:

The code runs multiple episodes (100 in one loop, 10000 in another loop), updating the policy and Q-values after each episode. This iterative process is the core of the learning mechanism.

Recording and Visualization:

The Q-values are recorded after each episode, allowing for tracking their progression over time. The plot_Q_values function generates a plot to visually represent the evolution of Q-values, providing insights into the learning process and the convergence of the algorithm.

Output:

The policy and Q-values are printed after the training loop to show the learned policy's probabilities and the estimated values for each action in each state.

```
1  import random
2  import matplotlib.pyplot as plt
3
4
5  # Define the state space
6  states = [0, 1, 2, 3, 4, 5]
7
8  # Define the action space
9  actions = [-1, 1]  # -1 for left, 1 for right
10
11 # Initialize Q-values and returns
12 Q = {state: {action: 0.0 for action in actions} for state in states}
13 returns = {state: {action: [] for action in actions} for state in states}
14
15 # Set the value of epsilon for the ε-soft policy
16 epsilon = 0.1
17
18 # Initialize policy to an ε-soft policy
19 policy = {state: {action: epsilon / len(actions) for action in actions} for state in states}
20 for state in states:
```

```python
20  for state in states:
21      if state not in [0, 5]:  # No policy needed for terminal states, but we'll initialize them anyway
22          best_action = random.choice(actions)
23          policy[state][best_action] += (1.0 - epsilon)
24
25  # Define the transition probabilities and rewards in the step function
26  def step(state, action):
27      if state == 0 or state == 5:
28          return state, 0  # No reward for terminal states
29      next_state_probabilities = {
30          state + action: 0.8,
31          state: 0.15,
32          state - action: 0.05
33      }
34      next_state = random.choices(list(next_state_probabilities.keys()), weights=list(next_state_probabilities.values()), k=1)[0]
35      reward = 0
36      if state == 4 and action == 1 and next_state == 5:
37          reward = 5
38      elif state == 1 and action == -1 and next_state == 0:
39          reward = 1
40      return next_state, reward
41
42  # Function to choose an action based on the ε-soft policy
43  def choose_action(state, policy):
44      action_probabilities = policy[state]
45      return random.choices(list(action_probabilities.keys()), weights=list(action_probabilities.values()), k=1)[0]
46
47  # Function to generate an episode
48  def generate_episode(start_state, policy):
49      episode = []
50      state = start_state
51      while state not in [0, 5]:
52          action = choose_action(state, policy)
53          next_state, reward = step(state, action)
54          episode.append((state, action, reward))
55          if next_state in [0, 5]:  # Terminal state reached
56              episode.append((next_state, None, 0))  # Append terminal state
57              break
58          state = next_state
59      return episode
60
61  # Function to update Q-values and policy
62  def update_policy(episode, policy, Q, returns, gamma=1.0):
63      G = 0
64      for state, action, reward in reversed(episode):
65          if action is not None:  # Skip terminal state
66              G = gamma * G + reward
67              returns[state][action].append(G)
68              Q[state][action] = sum(returns[state][action]) / len(returns[state][action])
69      # Update policy
70      for state in states:
71          if state not in [0, 5]:
72              best_action = max(Q[state], key=Q[state].get)
73              for action in actions:
74                  if action == best_action:
75                      policy[state][action] = 1 - epsilon + (epsilon / len(actions))
76                  else:
77                      policy[state][action] = epsilon / len(actions)
78
79  # Run multiple episodes and update policy
80  for i in range(100):
81      start_state = 3  # Always start from state 3
82      episode = generate_episode(start_state, policy)
83      update_policy(episode, policy, Q, returns)
84
85  # Display the updated policy and Q-values
86  print("Policy after 100 episodes:")
87  for state in policy:
88      print(f"State {state}: {policy[state]}")
89
90  print("\nQ-values after 100 episodes:")
91  for state in Q:
92      print(f"State {state}: {Q[state]}")
93
94
95  # Initialize a dictionary to record Q-values at each episode for each state-action pair
96  Q_history = {state: {action: [] for action in actions} for state in states}
97
98  # Run multiple episodes and update policy
99  for i in range(100):
100     start_state = 3  # Always start from state 3
101     episode = generate_episode(start_state, policy)
102     update_policy(episode, policy, Q, returns)
```

```
103
104     # Record the current Q-values
105     for state in Q:
106         for action in Q[state]:
107             Q_history[state][action].append(Q[state][action])
108
109 # Function to plot Q-values
110 def plot_Q_values(Q_history):
111     for state in Q_history:
112         for action in Q_history[state]:
113             plt.plot(Q_history[state][action], label=f"State {state}, Action {action}")
114     plt.xlabel('Episodes')
115     plt.ylabel('Q-value')
116     plt.title('Q-values Over Time')
117     plt.legend()
118     plt.show()
119
120 # Call the function to plot Q-values
121 plot_Q_values(Q_history)
```
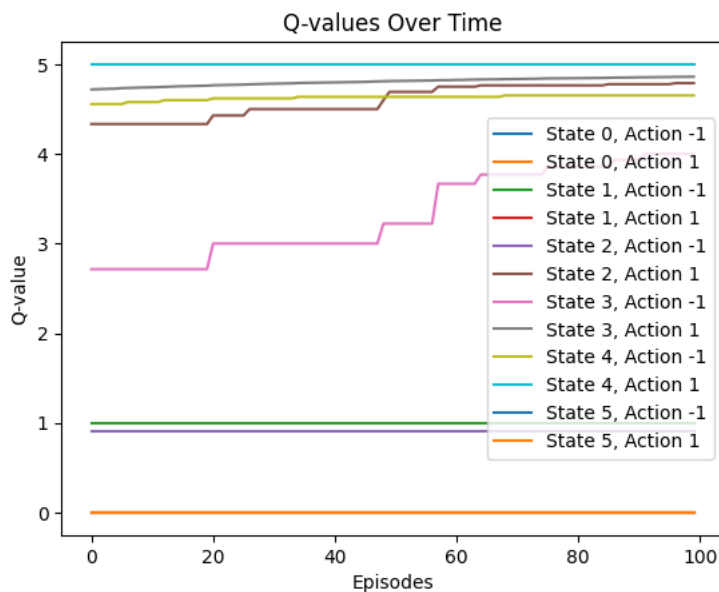
```
Policy after 100 episodes:
State 0: {-1: 0.05, 1: 0.05}
State 1: {-1: 0.9500000000000001, 1: 0.05}
State 2: {-1: 0.05, 1: 0.9500000000000001}
State 3: {-1: 0.05, 1: 0.9500000000000001}
State 4: {-1: 0.05, 1: 0.9500000000000001}
State 5: {-1: 0.05, 1: 0.05}

Q-values after 100 episodes:
State 0: {-1: 0.0, 1: 0.0}
State 1: {-1: 1.0, 1: 0.0}
State 2: {-1: 0.9, 1: 4.333333333333333}
State 3: {-1: 2.7142857142857144, 1: 4.717557251908397}
State 4: {-1: 4.555555555555555, 1: 5.0}
State 5: {-1: 0.0, 1: 0.0}
```



Q-values Over Time

Over the course of 100 episodes, the algorithm is effectively learning which actions are more valuable in the given states.

The policy is being refined to be more greedy with respect to the Q-values, as it increasingly favors actions that yield higher returns. The plot suggests that the learning process is ongoing.

While some Q-values have stabilized (notably for the most rewarding actions), others may still be in the process of being accurately estimated. This is part of the exploration-exploitation balance inherent in reinforcement learning.

```python
1  # Run multiple episodes and update policy
2  for i in range(1000):
3      start_state = 3  # Always start from state 3
4      episode = generate_episode(start_state, policy)
5      update_policy(episode, policy, Q, returns)
6
7  # Display the updated policy and Q-values
8  print("Policy after 1000 episodes:")
9  for state in policy:
10     print(f"State {state}: {policy[state]}")
11
12 print("\nQ-values after 1000 episodes:")
13 for state in Q:
14     print(f"State {state}: {Q[state]}")
15
16
17 # Initialize a dictionary to record Q-values at each episode for each state-action pair
18 Q_history = {state: {action: [] for action in actions} for state in states}
19
20 # Run multiple episodes and update policy
21 for i in range(1000):
22     start_state = 3  # Always start from state 3
23     episode = generate_episode(start_state, policy)
24     update_policy(episode, policy, Q, returns)
25
26     # Record the current Q-values
27     for state in Q:
28         for action in Q[state]:
29             Q_history[state][action].append(Q[state][action])
30
31 # Function to plot Q-values
32 def plot_Q_values(Q_history):
33     for state in Q_history:
34         for action in Q_history[state]:
35             plt.plot(Q_history[state][action], label=f"State {state}, Action {action}")
36     plt.xlabel('Episodes')
37     plt.ylabel('Q-value')
38     plt.title('Q-values Over Time')
39     plt.legend()
40     plt.show()
41
42 # Call the function to plot Q-values
43 plot_Q_values(Q_history)
```
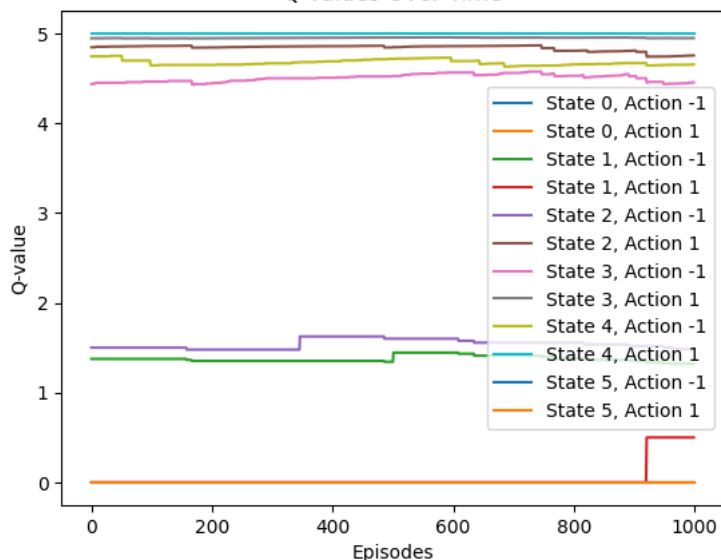
```
Policy after 1000 episodes:
State 0: {-1: 0.05, 1: 0.05}
State 1: {-1: 0.9500000000000001, 1: 0.05}
State 2: {-1: 0.05, 1: 0.9500000000000001}
State 3: {-1: 0.05, 1: 0.9500000000000001}
State 4: {-1: 0.05, 1: 0.9500000000000001}
State 5: {-1: 0.05, 1: 0.05}

Q-values after 1000 episodes:
State 0: {-1: 0.0, 1: 0.0}
State 1: {-1: 1.375, 1: 0.0}
State 2: {-1: 1.5, 1: 4.84516129032258}
State 3: {-1: 4.435897435897436, 1: 4.94521372667068}
State 4: {-1: 4.744186046511628, 1: 5.0}
State 5: {-1: 0.0, 1: 0.0}
```



Q-values Over Time

After 1,000 episodes, the Monte Carlo algorithm has significantly refined the policy for the cleaning robot, with the policy becoming more deterministic in favor of actions that lead to higher rewards.

The results indicate a successful learning process, with the algorithm likely nearing an optimal policy for this particular Markov Decision Process (MDP).

Further training could be done to confirm the stability of the policy, but at this point, it may be yielding diminishing returns in terms of new information gained.

```python
1  # Run multiple episodes and update policy
2  for i in range(10000):
3      start_state = 3  # Always start from state 3
4      episode = generate_episode(start_state, policy)
5      update_policy(episode, policy, Q, returns)
6
7  # Display the updated policy and Q-values
8  print("Policy after 10000 episodes:")
9  for state in policy:
10     print(f"State {state}: {policy[state]}")
11
12 print("\nQ-values after 10000 episodes:")
13 for state in Q:
14     print(f"State {state}: {Q[state]}")
15
16
17 # Initialize a dictionary to record Q-values at each episode for each state-action pair
18 Q_history = {state: {action: [] for action in actions} for state in states}
19
20 # Run multiple episodes and update policy
21 for i in range(10000):
22     start_state = 3  # Always start from state 3
23     episode = generate_episode(start_state, policy)
24     update_policy(episode, policy, Q, returns)
25
26     # Record the current Q-values
27     for state in Q:
28         for action in Q[state]:
29             Q_history[state][action].append(Q[state][action])
30
31 # Function to plot Q-values
32 def plot_Q_values(Q_history):
33     for state in Q_history:
34         for action in Q_history[state]:
35             plt.plot(Q_history[state][action], label=f"State {state}, Action {action}")
36     plt.xlabel('Episodes')
37     plt.ylabel('Q-value')
38     plt.title('Q-values Over Time')
39     plt.legend()
40     plt.show()
41
42 # Call the function to plot Q-values
43 plot_Q_values(Q_history)
```

```
Policy after 10000 episodes:
State 0: {-1: 0.05, 1: 0.05}
State 1: {-1: 0.05, 1: 0.9500000000000001}
State 2: {-1: 0.05, 1: 0.9500000000000001}
State 3: {-1: 0.05, 1: 0.9500000000000001}
State 4: {-1: 0.05, 1: 0.9500000000000001}
State 5: {-1: 0.05, 1: 0.05}

Q-values after 10000 episodes:
State 0: {-1: 0.0, 1: 0.0}
State 1: {-1: 1.3934426229508197, 1: 4.391666666666667}
State 2: {-1: 3.514851485148515, 1: 4.893956670467503}
State 3: {-1: 4.841370558375634, 1: 4.968505338078292}
State 4: {-1: 4.663225806451613, 1: 4.997353622229573}
State 5: {-1: 0.0, 1: 0.0}
```