Giorgio Mendoza

RBE595-S24-S04

Model-based RL Programming Exercise

The following code implements the Dyna-Q algorithm to solve the Dyna Maze problem. The key components of the code are:

## ∨  MazeEnvironment Class:

This class models the maze environment. It includes methods to:

### step:

Transition to the next state based on the agent's action. If the action leads to an obstacle or out-of-bounds, the state does not change. reset: Reset the environment to the starting state. Action Selection: The epsilon_greedy_action_selection function chooses actions based on the epsilon-greedy policy, balancing exploration and exploitation.

### Q-Value Update:

The q_learning_update function updates the Q-values using the standard Q-learning formula, considering the reward received and the estimated value of the next state.

### Planning:

The planning function simulates experiences using the learned model to further update the Q-values, which accelerates the learning process.
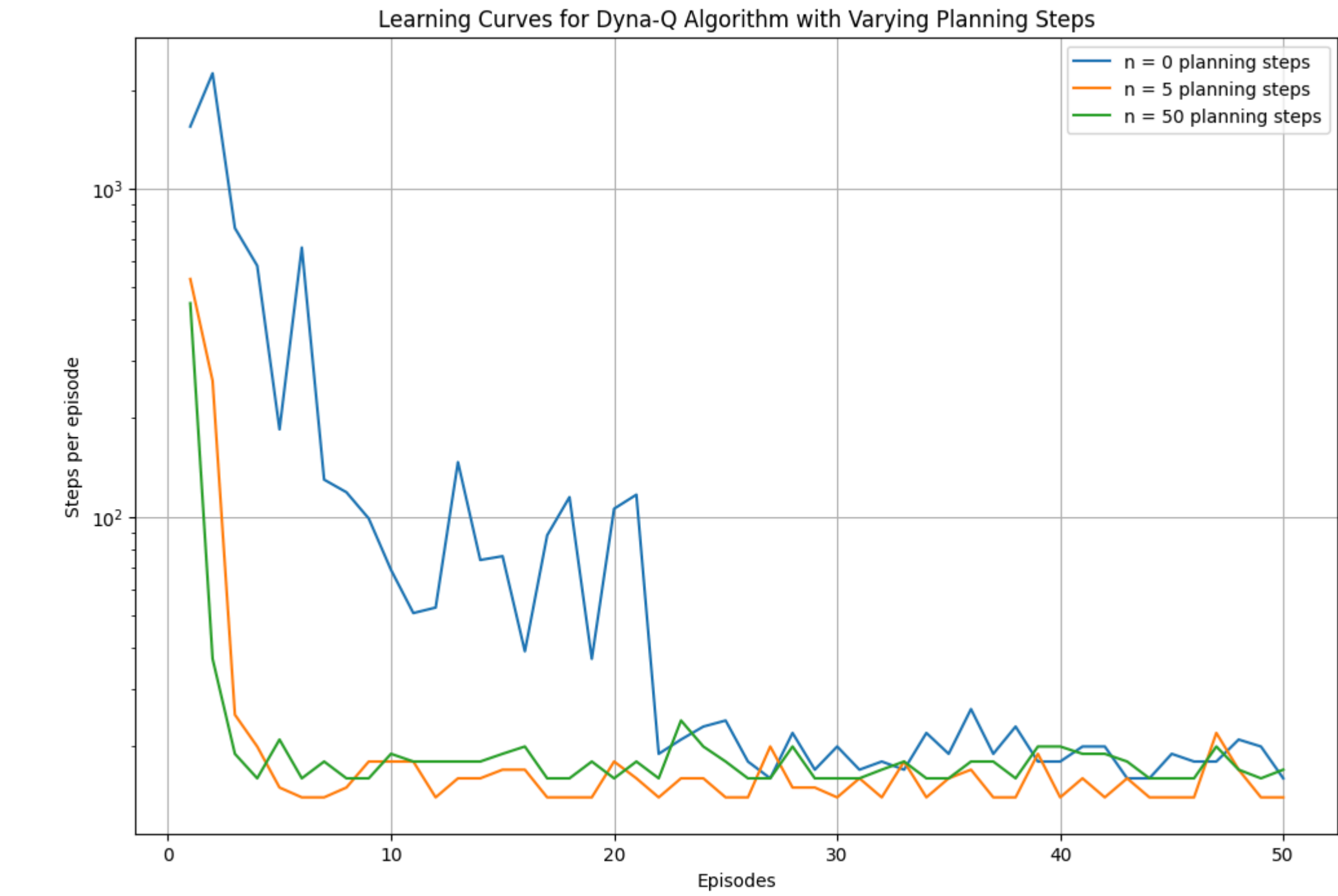
### Simulation Loop:

The simulation runs for a defined number of episodes. In each episode, the agent interacts with the environment, updates the Q-values, updates the model, and performs a number of planning steps.

### Resetting:

Between simulations with different numbers of planning steps, the Q-values and Model are reset to ensure independent learning curves.

```python
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import random
4
5 # Define the gray squares' positions (obstacles in the maze)
6 gray_squares = [(4, 2), (3, 2), (2, 2), (1, 5), (5, 7), (4, 7), (3, 7)]
7
8 # Define the start and goal positions
9 start_pos = (3, 0)   # Start at (3,0)
10 goal_pos = (5, 8)    # Goal at (5,8)
11
12 # Define the MazeEnvironment class
13 class MazeEnvironment:
14     def __init__(self, start_pos, goal_pos, gray_squares, width=9, height=6):
15         self.start_pos = start_pos
16         self.goal_pos = goal_pos
17         self.gray_squares = gray_squares
18         self.width = width
19         self.height = height
20         self.state = start_pos
21
22     def step(self, action):
23         # Define action effects
24         actions = {
25             'up': (-1, 0),
26             'down': (1, 0),
27             'left': (0, -1),
28             'right': (0, 1)
29         }
30         # Determine the next state
31         next_state = (self.state[0] + actions[action][0], self.state[1] + actions[action][1])
32
33         # Check for obstacles or out-of-bounds, and adjust the state accordingly
34         if (next_state in self.gray_squares or
35                 next_state[0] < 0 or next_state[0] >= self.height or
36                 next_state[1] < 0 or next_state[1] >= self.width):
37             next_state = self.state
38
39         # Check for the goal state
40         reward = 1 if next_state == self.goal_pos else 0
41         done = next_state == self.goal_pos
42
43         # Update the state
44         self.state = next_state if not done else self.start_pos
45         return next_state, reward, done
46
47     def reset(self):
48         self.state = self.start_pos
49         return self.state
50
51 # Initialize the Maze environment
52 env = MazeEnvironment(start_pos=start_pos, goal_pos=goal_pos, gray_squares=gray_squares)
53
54 # Define the actions available to the agent
55 actions = ['up', 'down', 'left', 'right']
56
57 # Initialize Q-values and Model
58 Q = {}
59 Model = {}
60
61 # Define epsilon-greedy action selection function
62 def epsilon_greedy_action_selection(state, Q, epsilon=0.1):
63     if np.random.rand() < epsilon:
64         return np.random.choice(actions)
65     else:
66         q_values = [Q.get((state, action), 0) for action in actions]
67         max_q = max(q_values)
68         # In case there are several actions with the same Q-value, select randomly among them
69         actions_with_max_q = [actions[i] for i, q in enumerate(q_values) if q == max_q]
70         return np.random.choice(actions_with_max_q)
71
72 # Define Q-learning update function
73 def q_learning_update(state, action, reward, next_state, Q, alpha=0.1, gamma=0.95):
74     best_next_action = epsilon_greedy_action_selection(next_state, Q, epsilon=0)
75     Q[(state, action)] = Q.get((state, action), 0) + alpha * (reward + gamma * Q.get((next_state, best_next_action), 0) - Q.get((state, action
76
77 # Define planning function for Dyna-Q
78 def planning(Q, Model, n_planning_steps, alpha=0.1, gamma=0.95):
79     for _ in range(n_planning_steps):
80         # Randomly sample a previously observed state and action
81         state, action = random.choice(list(Model.keys()))
82         next_state, reward = Model[(state, action)]
83         q_learning_update(state, action, reward, next_state, Q, alpha, gamma)
84
85 # Simulation parameters
86 num_episodes = 50   # Total number of episodes to simulate
87 num_planning_steps = [0, 5, 50]  # Number of planning steps for Dyna-Q
88 steps_per_episode = {n: [] for n in num_planning_steps}   # To record the steps taken in each episode
89
90 # Run the Dyna-Q algorithm
91 for n in num_planning_steps:  # For each number of planning steps
92     Q.clear()  # Reset Q-values for each series of simulations
93     Model.clear()  # Reset Model for each series of simulations
94     for episode in range(num_episodes):
95         state = env.reset()
96         steps = 0
97         while True:
98             action = epsilon_greedy_action_selection(state, Q, epsilon=0.1)
99             next_state, reward, done = env.step(action)
100            q_learning_update(state, action, reward, next_state, Q, alpha=0.1, gamma=0.95)
101            Model[(state, action)] = (next_state, reward)   # Update the model
102            planning(Q, Model, n, alpha=0.1, gamma=0.95)   # Planning step
103            state = next_state
104            steps += 1
105            if done:
106                break
107        steps_per_episode[n].append(steps)
108
109 # Plot the results
110 plt.figure(figsize=(12, 8))
111 for n, steps in steps_per_episode.items():
112     plt.plot(range(1, num_episodes + 1), steps, label=f'n = {n} planning steps')
113 plt.xlabel('Episodes')
114 plt.ylabel('Steps per episode')
115 plt.title('Learning Curves for Dyna-Q Algorithm with Varying Planning Steps')
116 plt.yscale('log')  # Log scale for better visibility
117 plt.legend()
118 plt.grid(True)
119 plt.show()
120
```

Learning Curves for Dyna-Q Algorithm with Varying Planning Steps

## Conclusion:

The plot illustrates the learning curves of the Dyna-Q algorithm with varying numbers of planning steps (0, 5, and 50). Each curve represents the average number of steps the agent took to reach the goal in each episode:

The blue curve (n = 0) shows the learning performance without planning. It has the slowest learning rate, with the number of steps per episode decreasing gradually as the agent learns from actual experiences.

The orange curve (n = 5) indicates faster learning due to the addition of a moderate amount of planning. The agent benefits from simulated experiences, which improve its policy more quickly than learning from real interactions alone.

The green curve (n = 50) demonstrates the most rapid learning. With extensive planning, the agent frequently updates its policy based on simulated experiences, which results in a swift decline in the number of steps needed to reach the goal.

Overall, the plot conveys a clear message: incorporating planning into the learning process can significantly speed up an agent's acquisition of an effective policy, and the more planning steps the agent performs, the faster it learns. This aligns with the fundamental principles of model-based reinforcement learning, where simulated experiences complement real interactions to enhance learning efficiency.