

# Homework #2: Reinforcement Learning Code Analysis

Yu Tianrun

October 28, 2024

## Abstract

This code implements a Q-Learning algorithm to solve a maze navigation problem. The environment is modeled as a 4x5 grid where each cell represents an open path, an obstacle, or a goal. The primary objective of the reinforcement learning agent is to navigate from the starting position to the goal efficiently, minimizing the number of steps taken. The agent learns optimal policies by iteratively updating a Q-table, which estimates the value of taking specific actions in given states. The algorithm employs an epsilon-greedy strategy to balance exploration of new actions and exploitation of known high-reward actions. Throughout multiple training episodes, the agent refines its policy, enabling it to effectively traverse the maze and consistently reach the goal. This implementation demonstrates the fundamental principles of reinforcement learning, including state representation, action selection, reward processing, and policy optimization.

## Code with Detailed Comments

```
1 import numpy as np
2 import random
3
4 # Define the maze environment
5 class MazeEnv:
6     def __init__(self):
7         # 4x5 maze grid: 0 is free cell, -1 is obstacle, 1 is goal
8         self.maze = np.array([[0, 0, 0, 0, -1],
9                                [0, -1, 0, -1, 0],
10                               [0, -1, 0, 1, 0],
11                               [0, 0, 0, -1, 0]])
12         self.start_position = (0, 0)
13         self.reset()
14
15     def reset(self):
16         self.position = self.start_position
17         return self.position
18
19     def step(self, action):
20         # Define possible movements: right, down, left, up
21         moves = [(0, 1), (1, 0), (0, -1), (-1, 0)]
22         next_position = (self.position[0] + moves[action][0],
23                         self.position[1] + moves[action][1])
24
25         # Check if next position is within bounds and not an obstacle
26         if (0 <= next_position[0] < self.maze.shape[0] and
27             0 <= next_position[1] < self.maze.shape[1] and
28             self.maze[next_position] != -1):
29             self.position = next_position
30
31         # Check if goal is reached
32         if self.maze[self.position] == 1:
33             return self.position, 1, True # Goal reached
34         else:
35             return self.position, -0.1, False # Step penalty
36
37 # Q-Learning parameters
38 alpha = 0.1 # Learning rate
39 gamma = 0.9 # Discount factor
```

```

40 epsilon = 0.1 # Exploration rate
41 episodes = 500 # Number of training episodes
42
43 # Initialize Q-table for the 4x5 grid with 4 actions per state
44 q_table = np.zeros((4, 5, 4))
45
46 # Q-Learning algorithm
47 env = MazeEnv()
48
49 for episode in range(episodes):
50     state = env.reset()
51     done = False
52
53     while not done:
54         # Epsilon-greedy strategy for action selection
55         if random.uniform(0, 1) < epsilon:
56             action = random.choice(range(4)) # Random action (exploration)
57         else:
58             action = np.argmax(q_table[state[0], state[1]]) # Best action (exploitation)
59
60         # Take action and observe the outcome
61         next_state, reward, done = env.step(action)
62
63         # Update Q-value for the current state-action pair
64         old_value = q_table[state[0], state[1], action]
65         next_max = np.max(q_table[next_state[0], next_state[1]])
66         new_value = (1 - alpha) * old_value + alpha * (reward + gamma * next_max)
67         q_table[state[0], state[1], action] = new_value
68
69         # Transition to the next state
70         state = next_state

```

Listing 1: Reinforcement Learning Code for Maze Navigation with Detailed Comments

## Core Section with Line-by-Line Comments

```

1 # Initialize environment and Q-table for the Q-learning algorithm
2 env = MazeEnv() # Instantiate the MazeEnv class to create the maze environment
3 q_table = np.zeros((4, 5, 4)) # Initialize a Q-table with dimensions corresponding to
   the maze grid (4 rows, 5 columns) and 4 possible actions
4
5 # Define Q-Learning parameters
6 alpha = 0.1 # Learning rate: determines the extent to which newly acquired
   information overrides old information
7 gamma = 0.9 # Discount factor: measures the importance of future rewards versus
   immediate rewards
8 epsilon = 0.1 # Exploration rate: probability of choosing a random action instead of
   the best-known action
9 episodes = 500 # Total number of training episodes
10
11 # Q-Learning algorithm: training the agent over multiple episodes
12 for episode in range(episodes): # Loop through each episode
13     state = env.reset() # Reset the environment to the starting state and obtain the
   initial state
14     done = False # Flag to indicate whether the goal has been reached
15
16     while not done: # Continue the episode until the goal is reached
17         if random.uniform(0, 1) < epsilon: # Epsilon-greedy strategy: decide whether to
   explore or exploit
18             action = random.choice(range(4)) # Exploration: randomly select an action
   (0-3 corresponding to up, down, left, right)
19         else:
20             action = np.argmax(q_table[state[0], state[1]]) # Exploitation: choose the
   action with the highest Q-value for the current state
21
22         # Execute the chosen action in the environment
23         next_state, reward, done = env.step(action) # Perform the action and observe
   the next state, reward, and whether the goal is reached
24
25         # Q-value update using the Q-Learning update rule
26         old_value = q_table[state[0], state[1], action] # Retrieve the current Q-value
   for the state-action pair

```

```

27     next_max = np.max(q_table[next_state[0], next_state[1]]) # Find the maximum Q-
    value for the next state across all possible actions
28
29     # Calculate the new Q-value using the Bellman equation
30     new_value = (1 - alpha) * old_value + alpha * (reward + gamma * next_max)
31     q_table[state[0], state[1], action] = new_value # Update the Q-table with the
    new Q-value
32
33     state = next_state # Transition to the next state for the subsequent step

```

Listing 2: Core Q-Learning Algorithm