# CSE - 584 HomeWork 2

**Name :** Raghavendra Jagirdar

**PSU Id :** 911938232

## 1. Abstract

This code implements a Q-Learning algorithm to solve a simple problem in a grid environment. The grid environment consists of a 2D grid where an agent starts at one location and tries to reach a goal location while avoiding obstacles. The Q-Learning algorithm helps the agent learn the best policy through interactions with the environment. At each step, the agent receives feedback in the form of rewards based on its actions, and it updates a Q-value table to improve its decisions. The main goal is to maximize the cumulative reward by finding an optimal policy that specifies the best action to take in each state.

## 2. Annotated Code

```
1   import numpy as np
2
3   # Define the environment parameters
4   num_states = 16  # Number of states in the grid (4x4 grid, for instance)
5   num_actions = 4  # Number of actions (up, down, left, right)
6   q_table = np.zeros((num_states, num_actions))  # Initialize Q-table with zeros
7
8   # Hyperparameters for the Q-learning algorithm
9   alpha = 0.1        # Learning rate (how much we accept new information)
10  gamma = 0.9        # Discount factor (how much we consider future rewards)
11  epsilon = 0.1      # Exploration factor (probability of choosing random action)
12
13  # Function to choose an action using the epsilon-greedy strategy
14  def choose_action(state):
15      # With probability epsilon, choose a random action (exploration)
16      if np.random.uniform(0, 1) < epsilon:
17          return np.random.randint(0, num_actions)
18      else:
19          # Choose the action with the highest Q-value (exploitation)
20          return np.argmax(q_table[state, :])
21
22  # Function to update the Q-value table
23  def update_q_table(state, action, reward, next_state):
24      # Get the maximum Q-value for the next state
25      max_q_next = np.max(q_table[next_state, :])
26      # Calculate the new Q-value using the Q-learning formula
27      q_table[state, action] = q_table[state, action] + \
28          alpha * (reward + gamma * max_q_next - q_table[state, action])
29
30  # Simulate the environment
31  def simulate_environment(state, action):
32      # Example environment dynamics
33      # Assume the agent moves to a new state based on action taken
34      next_state = state  # Placeholder for the next state
35      reward = -1  # Default reward is -1 (penalty for each step)
36
37      # Example transitions (these should be customized for a real problem)
```

```python
25        max_q_next = np.max(q_table[next_state, :])
26        # Calculate the new Q-value using the Q-learning formula
27        q_table[state, action] = q_table[state, action] + \
28            alpha * (reward + gamma * max_q_next - q_table[state, action])
29
30    # Simulate the environment
31    def simulate_environment(state, action):
32        # Example environment dynamics
33        # Assume the agent moves to a new state based on action taken
34        next_state = state  # Placeholder for the next state
35        reward = -1  # Default reward is -1 (penalty for each step)
36
37        # Example transitions (these should be customized for a real problem)
38        if state == 0 and action == 1:  # If at state 0 and taking action 'down'
39            next_state = 4
40            reward = 0  # Neutral reward
41        elif state == 4 and action == 1:  # If at state 4 and taking action 'down'
42            next_state = 8
43            reward = 10  # Reward for reaching goal state
44
45        # Return the next state and reward
46        return next_state, reward
47
48    # Main Q-Learning loop
49    num_episodes = 1000  # Total episodes for training
50    for episode in range(num_episodes):
51        state = 0  # Starting state for each episode (e.g., top-left corner of grid)
52
53        while state != 8:  # Continue until goal state (e.g., bottom-right corner)
54            action = choose_action(state)  # Choose action based on epsilon-greedy
55            next_state, reward = simulate_environment(state, action)  # Simulate env
56            update_q_table(state, action, reward, next_state)  # Update Q-table
57            state = next_state  # Move to the next state
58
59    # Display the learned Q-table
60    print("Learned Q-Table:")
61    print(q_table)
```

**Output :**

```
Learned Q-Table:
[[ 6.35806493  9.           6.72638768  6.47073138]
 [ 0.          0.           0.          0.         ]
 [ 0.          0.           0.          0.         ]
 [ 0.          0.           0.          0.         ]
 [ 7.4763163  10.           7.72345135  7.28964404]
 [ 0.          0.           0.          0.         ]
 [ 0.          0.           0.          0.         ]
 [ 0.          0.           0.          0.         ]
 [ 0.          0.           0.          0.         ]
 [ 0.          0.           0.          0.         ]
 [ 0.          0.           0.          0.         ]
 [ 0.          0.           0.          0.         ]
 [ 0.          0.           0.          0.         ]
 [ 0.          0.           0.          0.         ]
 [ 0.          0.           0.          0.         ]
 [ 0.          0.           0.          0.         ]]
```

**Explanation of the Core Sections**

Here are the core functions with detailed comments added to each line:

1) `choose_action(state)`
    a) Implements the epsilon-greedy strategy, balancing exploration and exploitation.
    b) With probability `epsilon`, it selects a random action to explore new possibilities.
    c) Otherwise, it chooses the action with the highest Q-value to exploit known information.
2) `update_q_table(state, action, reward, next_state)`
    a) Calculates the maximum Q-value for the next state, representing the expected future reward.
    b) Updates the Q-value of the current state-action pair using the Q-Learning update formula:
    $Q(s,a) \leftarrow Q(s,a) + \alpha(r + \gamma \max_{a'} Q(s',a') - Q(s,a))$ $Q(s, a) \leftarrow Q(s, a) + \alpha \left( r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right)$ $Q(s,a) \leftarrow Q(s,a) + \alpha(r + \gamma \max_{a'} Q(s',a') - Q(s,a))$
    c) Here, $Q(s,a)$ $Q(s, a)$ $Q(s,a)$ is the Q-value for the current state-action pair, $\alpha$ $\alpha$ $\alpha$ is the learning rate, $r$ $r$ $r$ is the immediate reward, $\gamma$ $\gamma$ $\gamma$ is the discount factor, and $\max_{a'} Q(s',a')$ $\max_{a'} Q(s', a')$ $\max_{a'} Q(s',a')$ is the maximum Q-value for the next state.
3) `simulate_environment(state, action)`
    a) Simulates environment dynamics based on the current state and action.
    b) Returns the next state and reward to be used for updating the Q-table.
    c) In this simple implementation, the dynamics are manually defined; in more complex scenarios, this function would interact with an actual environment (e.g., OpenAI Gym).