# Neuroscience of Learning, Memory and Cognition

Instructor: Prof. Karbalaei

## Sharif University of Technology

# Project: Solving Maze Using Q-Learning Algorithm

By Reza Nayeb Habib
Student ID# 401102694

# Contents

# 1   The Maze Generation Algorithm

The Maze Generation Algorithm is coded as below

```python
1  # Adapted from http://code.activestate.com/recipes/578356-random-maze-generator/
2  # Random Maze Generator using Depth-first Search
3  # http://en.wikipedia.org/wiki/Maze_generation_algorithm
4
5
6  import random
7  import matplotlib.pyplot as plt
8  import numpy as np
9
10 random.seed(401102694)
11 mx = 20; my = 20 # width and height of the maze
12
13 maze = [[0 for x in range(mx)] for y in range(my)]
14 dx = [0, 1, 0, -1]; dy = [-1, 0, 1, 0] # 4 directions to move in the maze
15 color = [(0, 0, 0), (255, 255, 255)] # RGB colors of the maze
16
17 # start the maze from a random cell
18 cx = random.randint(0, mx - 1)
19 cy = random.randint(0, my - 1)
20 maze[cy][cx] = 1
21 stack = [(cx, cy, 0)] # stack element: (x, y, direction)
22
23 while len(stack) > 0:
24     (cx, cy, cd) = stack[-1]
25     # to prevent zigzags:
26     # if changed direction in the last move then cannot change again
27     if len(stack) > 2:
28         if cd != stack[-2][2]: dirRange = [cd]
29         else: dirRange = range(4)
30     else: dirRange = range(4)
31
32     # find a new cell to add
33     nlst = [] # list of available neighbors
34     for i in dirRange:
35         nx = cx + dx[i]
36         ny = cy + dy[i]
37         if nx >= 0 and nx < mx and ny >= 0 and ny < my:
38             if maze[ny][nx] == 0:
39                 ctr = 0 # of occupied neighbors must be 1
40                 for j in range(4):
41                     ex = nx + dx[j]; ey = ny + dy[j]
42                     if ex >= 0 and ex < mx and ey >= 0 and ey < my:
43                         if maze[ey][ex] == 1: ctr += 1
44                 if ctr == 1: nlst.append(i)
45
46     # if 1 or more neighbors available then randomly select one and move
47     if len(nlst) > 0:
48         ir = nlst[random.randint(0, len(nlst) - 1)]
49         cx += dx[ir]; cy += dy[ir]; maze[cy][cx] = 1
50         stack.append((cx, cy, ir))
51     else: stack.pop()
52
53 maze = np.array(maze)
54 maze -= 1
55 maze = abs(maze)
56
57 maze[0][0] = 0
58 maze[mx-1][my-1] = 0
```

```
59
60 np.save('maze', np.array(maze))
```

In the project document it is stated that there is always a path between its upper left point (0,0) to its lower right point (my-1,mx-1) this statement is true because the maze Generation Algorithm starts at a random point in the plane then carves out into different paths by seeing the blocks that have more than 2 wall neighbors (0 in the code) and this algorithm always uses already created path points neighbors to create new path points thus the maze is always connected. At the final lines, this algorithm adds (0,0) and (my-1,mx-1) to the path; but how does it make sure that these two points will be also connected by the previous path?

We can ensure these two points are always connected because no cell in this structure will have more than 2 walls beside it so it means every wall will always have at least one path neighbor! Thus when for example (0,0) cell is set to 1 it will always have a neighbor with value one and it will not be an isolated path point and it will be connected to all other path points.

## 2   Drawing The Maze

The following code draws the Maze by Loading the file and using matplotlibś Imshow:

```
1 maze = np.load('maze.npy')
2 plt.imshow(1-maze, cmap='gray')    # black is wall white is path
3 plt.xticks([]); plt.yticks([])
4 plt.show()
```
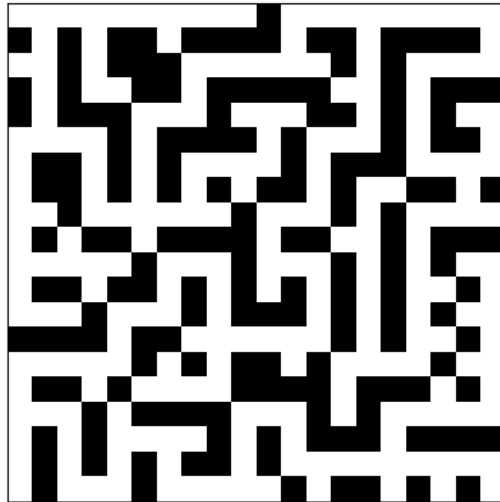
it gives the following result:



Figure 1: the unique maze of my student number seed

## 3   Q-Learning Algorithm

### 3.1   Full Code and Results

```
1 import numpy as np
2 import random
3 import matplotlib.pyplot as plt
4 import imageio
5 import os
```

```python
6
7  class Environment:
8      def __init__(self, maze_file):
9          self.maze = np.load(maze_file)
10         self.ny, self.nx = self.maze.shape
11         self.start = (0, 0)
12         self.goal = (self.nx - 1, self.ny - 1)
13         self.maze_file = maze_file
14         self.output_dir = os.path.join(os.path.dirname(maze_file), "output")
15         os.makedirs(self.output_dir, exist_ok=True)
16             # Create output directory if it doesn't exist
17
18     def get_valid_actions(self, x, y):
19         """Returns a list of valid actions for a given state."""
20         valid_actions = []
21         actions = [(0, -1), (1, 0), (0, 1), (-1, 0)]  # Up, Right, Down, Left
22         for i, (dx, dy) in enumerate(actions):
23             nx_, ny_ = x + dx, y + dy
24             if 0 <= nx_ < self.nx and 0 <= ny_ < self.ny and self.maze[ny_][nx_] == 0:
25                 valid_actions.append(i)
26         return valid_actions
27
28     def get_reward(self, x, y):
29         """Returns the reward for a given state."""
30         if (x, y) == self.goal:
31             return 100  # Big reward for reaching the goal
32         else:
33             return -1  # Small penalty for each step
34
35     def is_goal(self, x, y):
36         """Checks if the given state is the goal."""
37         return (x, y) == self.goal
38
39     def visualize_maze(self):
40         """Visualizes the maze."""
41         plt.imshow(self.maze, cmap='gray')
42         plt.title("Maze")
43         plt.show()
44
45 class Agent:
46     def __init__(self, env, alpha=0.1, gamma=0.9, epsilon=1.0,
47             epsilon_decay=0.995, epsilon_min=0.01):
48         self.env = env
49         self.alpha = alpha  # Learning rate
50         self.gamma = gamma  # Discount factor
51         self.epsilon = epsilon  # Initial exploration rate
52         self.epsilon_decay = epsilon_decay  # Decay per episode
53         self.epsilon_min = epsilon_min  # Minimum epsilon value
54         self.actions = [(0, -1), (1, 0), (0, 1), (-1, 0)]  # Up, Right, Down, Left
55         self.num_actions = len(self.actions)
56         self.Q_table = np.zeros((env.ny, env.nx, self.num_actions))  # Q-table
57
58     def choose_action(self, x, y):
59         """Epsilon-greedy policy."""
60         if random.uniform(0, 1) < self.epsilon:
61             return random.choice(self.env.get_valid_actions(x, y))  # Explore
62         else:
63             valid_actions = self.env.get_valid_actions(x, y)
64             return max(valid_actions, key=lambda a: self.Q_table[y, x, a])  # Exploit
65
66     def train(self, num_episodes):
67         """Trains the agent using Q-learning."""
```

```python
68          for episode in range(num_episodes):
69              x, y = self.env.start  # Start position
70              total_reward = 0
71
72              while not self.env.is_goal(x, y):  # Until reaching the goal
73                  action = self.choose_action(x, y)
74                  dx, dy = self.actions[action]
75                  new_x, new_y = x + dx, y + dy
76
77                  # Reward system
78                  reward = self.env.get_reward(new_x, new_y)
79
80                  # Q-value update
81                  best_next_action = max(self.env.get_valid_actions(new_x, new_y),
82                      key=lambda a: self.Q_table[new_y, new_x, a])
83                  self.Q_table[y, x, action] += self.alpha * (reward + self.gamma *
84                      self.Q_table[new_y, new_x, best_next_action] - self.Q_table[y, x, action])
85
86                  x, y = new_x, new_y
87                  total_reward += reward
88
89              # Decay epsilon
90              self.epsilon = max(self.epsilon * self.epsilon_decay, self.epsilon_min)
91
92              if episode % 500 == 0:
93                  print(f"Episode {episode}, Total Reward:
94                      {total_reward}, Epsilon: {self.epsilon:.3f}")
95
96      def solve_maze(self):
97          """Uses the trained Q-table to solve the maze."""
98          x, y = self.env.start
99          path = [(x, y)]
100         while not self.env.is_goal(x, y):
101             action = max(self.env.get_valid_actions(x, y),
102                 key=lambda a: self.Q_table[y, x, a])
103             dx, dy = self.actions[action]
104             x, y = x + dx, y + dy
105             path.append((x, y))
106         return path
107
108     def draw_policy(self, episode, save_path=None):
109         """Visualizes the best policy at each cell using arrows.
110         The optimal path (even if it leads to dead ends) is shown in green.
111         """
112         fig, ax = plt.subplots(figsize=(10, 10))
113         ax.imshow(1 - self.env.maze, cmap="gray")  # Show maze
114
115         arrow_map = {
116             0: (0, -1),   # Up
117             1: (1, 0),    # Right
118             2: (0, 1),    # Down
119             3: (-1, 0)    # Left
120         }
121
122         # Compute the best current path, even if not optimal (can hit dead ends)
123         x, y = self.env.start
124         visited = set()
125         best_path = []
126
127         while (x, y) not in visited and not self.env.is_goal(x, y):
128             visited.add((x, y))
129             best_path.append((x, y))
```

```python
130
131              valid_actions = self.env.get_valid_actions(x, y)
132              if not valid_actions:  # If no valid moves, stop
133                  break
134
135              best_action = max(valid_actions, key=lambda a: self.Q_table[y, x, a])
136              dx, dy = arrow_map[best_action]
137              x, y = x + dx, y + dy
138
139          best_path_set = set(best_path)  # Convert to set for quick lookup
140
141          # Draw policy arrows
142          for y in range(self.env.ny):
143              for x in range(self.env.nx):
144                  if self.env.maze[y, x] == 0:  # Ignore walls
145                      valid_actions = self.env.get_valid_actions(x, y)
146                      if valid_actions:  # If there are valid actions
147                          best_action = max(valid_actions, key=lambda a: self.Q_table[y, x, a])
148                          dx, dy = arrow_map[best_action]
149
150                          color = 'green' if (x, y) in best_path_set else 'red'
151                          ax.arrow(x, y, dx * 0.3, dy * 0.3, head_width=0.2,
152                              head_length=0.2, fc=color, ec=color)
153
154          ax.set_title(f"Policy Visualization - Episode {episode}")
155
156          if save_path:
157              plt.savefig(save_path)
158              plt.close(fig)  # Close the figure to avoid displaying it
159          else:
160              plt.show()
161
162
163      def draw_optimal_path(self, path=None, save_path=None):
164          """Draws the maze with the optimal path overlaid and optionally saves it as a PNG."""
165          if path is None:
166              path = self.solve_maze()  # Solve the maze if no path is provided
167
168          fig, ax = plt.subplots(figsize=(10, 10))
169          ax.imshow(self.env.maze, cmap="gray")  # Show maze
170
171          # Extract X and Y coordinates from path
172          x_coords, y_coords = zip(*path)
173
174          # Plot the path with red line
175          ax.plot(x_coords, y_coords, color='red', linewidth=2, marker='o'
176              , markersize=4, label="Optimal Path")
177
178          # Mark start and goal positions
179          ax.scatter([self.env.start[0]], [self.env.start[1]], color='green'
180              , s=100, label="Start (0,0)")
181          ax.scatter([self.env.goal[0]], [self.env.goal[1]], color='blue'
182              , s=100, label=f"Goal ({self.env.goal[0]},{self.env.goal[1]})")
183
184          ax.legend()
185          plt.title("Optimal Path in the Maze")
186
187          if save_path:
188              plt.savefig(save_path)
189              plt.close(fig)  # Close the figure to avoid displaying it
190          else:
191              plt.show()
```

```python
192
193     def train_with_visualization(self, num_episodes, frame_time, d_episode):
194         """Trains the agent and saves a GIF of the policy evolution."""
195         frames = []  # List to store frames
196
197         for episode in range(num_episodes):
198             x, y = self.env.start  # Start position
199             total_reward = 0
200
201             while not self.env.is_goal(x, y):  # Until reaching the goal
202                 action = self.choose_action(x, y)
203                 dx, dy = self.actions[action]
204                 new_x, new_y = x + dx, y + dy
205
206                 # Reward system
207                 reward = self.env.get_reward(new_x, new_y)
208
209                 # Q-value update
210                 best_next_action = max(self.env.get_valid_actions(new_x, new_y),
211                     key=lambda a: self.Q_table[new_y, new_x, a])
212                 self.Q_table[y, x, action] += self.alpha * (reward + self.gamma *
213                     self.Q_table[new_y, new_x, best_next_action] - self.Q_table[y, x, action])
214
215                 x, y = new_x, new_y
216                 total_reward += reward
217
218             # Decay epsilon
219             self.epsilon = max(self.epsilon * self.epsilon_decay, self.epsilon_min)
220
221             # Save policy visualization at key episodes
222             if episode % d_episode == 0 or episode == num_episodes - 1:
223                 save_path = os.path.join(self.env.output_dir, f"policy_episode_{episode}.png")
224                 self.draw_policy(episode, save_path)
225                 frames.append(imageio.imread(save_path))  # Store frame
226
227             if episode % d_episode == 0:
228                 print(f"Episode {episode}, Total Reward: {total_reward},
229                     Epsilon: {self.epsilon:.3f}")
230
231         # Save GIF with longer duration between frames
232         gif_path = os.path.join(self.env.output_dir, "policy_evolution.gif")
233         imageio.mimsave(gif_path, frames, duration=frame_time*num_episodes/d_episode)
234         print(f"GIF saved as {gif_path}")
235
236         # Clean up temporary files
237         for episode in range(0, num_episodes, d_episode):
238             os.remove(os.path.join(self.env.output_dir, f"policy_episode_{episode}.png"))
239         os.remove(os.path.join(self.env.output_dir, f"policy_episode_{num_episodes - 1}.png"))
240
241
242 # Main execution
243 if __name__ == "__main__":
244     env = Environment('.\\maze.npy')
245     agent = Agent(env)
246
247     # Train the agent
248     agent.train_with_visualization(num_episodes=1000, frame_time=5.0, d_episode=25)
249
250     # Solve the maze and visualize the optimal path
251     solution_path = agent.solve_maze()
252
253     optimal_path_image_path = os.path.join(env.output_dir, "optimal_path.png")
```

```
254
255     agent.draw_optimal_path(solution_path, save_path=optimal_path_image_path)
256     print(f"Optimal path visualization saved as {optimal_path_image_path}")
257     agent.draw_optimal_path(solution_path)
```

This code gives the following Results (the Gif file is included in the zip sent with this pdf file):
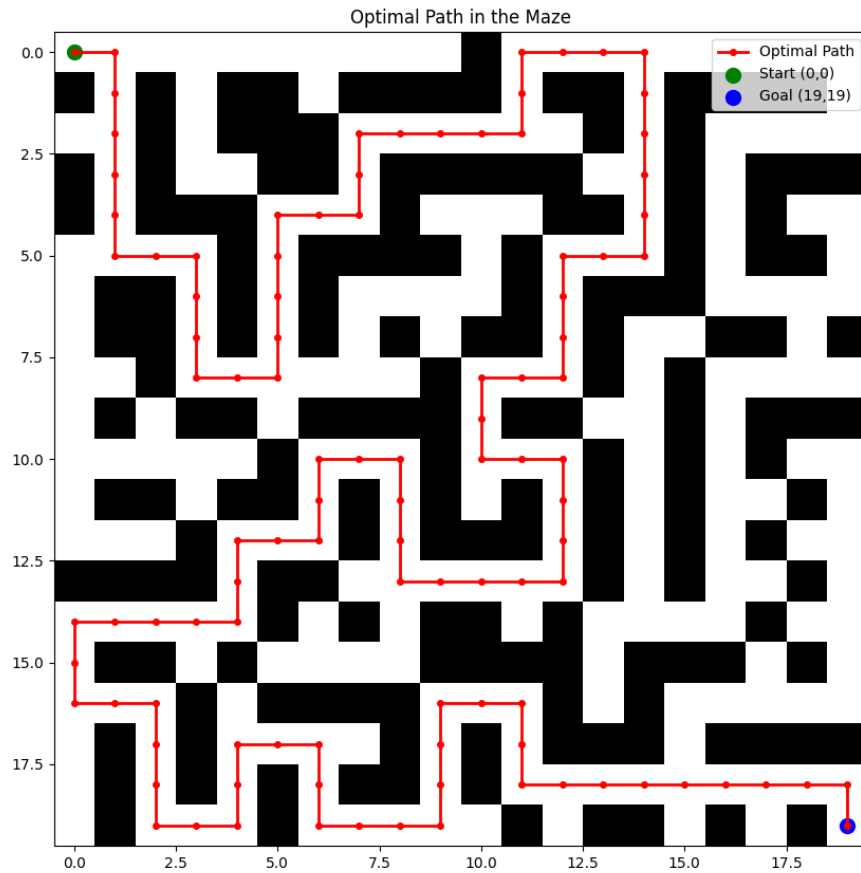


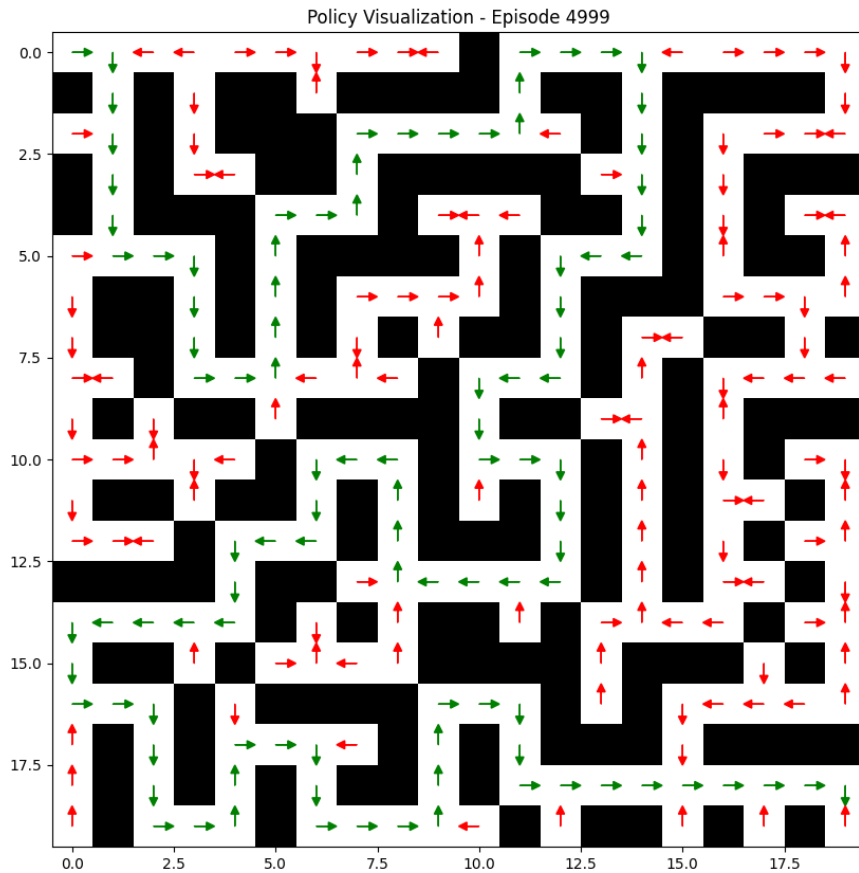Figure 2: the Optimal path in the maze found by Q-learning

Figure 3: the Optimal policy in the maze found by Q-learning

## 3.2  Code Explanation

`Environment`:
This class contains the maze cells mapping the reward function of the maze and the path the maze file is loaded from. Itś variables contain the goal point of the maze, the start point of the maze, the maze shape and ...

`get_valid_action(x,y)`:
This method of the `Environment` class gives the valid actions that can be taken at position (x,y) without hitting a wall or going outside of the maze.

`get_reward(x,y)`:
This method gives the reward of any given cell (x,y). all cells have a default reward of -1 to avoid the agent from wandering and punish him for taking longer paths. Also the goal point has a reward of 100 for motivating the agent to reach the end point.

`is_goal(x,y)`:
Checks if the cell is the goal cell; useful in functions needing to check if they have reached the end of the maze or not. Returns True if (nx,ny) is given.

`visualize_maze()`:
Uses `Imshow` to draw the maze.

`Agent`:
This class creates an agent that uses Q-Learning (Explained Later) to find the best policy for traversing in the maze. This class has the Environment, learning rate, Discount Factor, exploration rate and its decay factor, and the Q-table as its variables. It creates Q-table from the size of the maze and the number of actions that can be taken (number of directions that can be traversed) and Initializes them as zeros.

`choose_action(x,y)`:
This method of the Agent class is essentially the actor in actor-critic algorithm. It uses the Q-table to choose the best policy and act on that policy, this action is called exploitation; It also does exploration by generating a random value between 0 and 1 and if the value is smaller than the exploration probability (epsilon) it chooses its action randomly from valid actions.

`train(num_episodes)` and `train_with_visualization(num_episodes, frame_time, d_episode)`:
These functions implement the formula for Q-learning as stated below:

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left[ r + \gamma \max_{a'} Q(s',a') - Q(s,a) \right] \tag{1}$$

by this line of code:

```
1    self.Q_table[y, x, action] += self.alpha * (reward + self.gamma
2        * self.Q_table[new_y, new_x, best_next_action] - self.Q_table[y, x, action])
```

In every iteration it chooses an action by `choose_action(x,y)` and then updates the Q-table according the action taken. This process is repeated until the goal state is reached, then the episode is done and a new one starts it iterates this until the number of given episodes of learning is reached. The function with visualize is the same it only Visualizes the best policy using `draw_policy()` and saves it periodically on specific period of episodes; after that it creates a gif from the given files and saves it in the output folder.

`draw_policy(episode, save_path=None)`
This function first draws the maze with `Imshow` then uses `scatter` to draw the best policy at each cell with an arrow (it draws the arrows on the optimal path with green and others with red).

`draw_optimal_path (path=None, save_path =None)`:
This function finds the optimal path by using `solve_maze()` and getting the path that will be taken on exploitation mode and again uses scatter for drawing the optimal path given by `solve_maze`.

`solve_maze()`:
This fucntion follows the path given by Q-table until it reaches the goal point. It should not be used before the model is completely trained to reach the goal point otherwise it will enter an infinite loop.

## 3.3 Deep Q-Learning Implementation

The Deep Q-Learning algorithm is very similar to normal Q-Learning. The difference is that the state variables (here x and y) are fed into a deep neural network and it gives the Q value for all possible actions. Here we show how we implemented this with pytorch:

```
1 import numpy as np
2 import random
```

```python
 3  import torch
 4  import torch.nn as nn
 5  import torch.optim as optim
 6  import torch.nn.functional as F
 7  import matplotlib.pyplot as plt
 8  import imageio
 9  import os
10  from collections import deque
11
12  # Environment Class
13  class Environment:
14      def __init__(self, maze_file):
15          self.maze = np.load(maze_file)
16          self.ny, self.nx = self.maze.shape
17          self.start = (0, 0)
18          self.goal = (self.nx - 1, self.ny - 1)
19          self.actions = [(0, -1), (1, 0), (0, 1), (-1, 0)]  # Up, Right, Down, Left
20          self.output_dir = os.path.join(os.path.dirname(maze_file), "output_dqn")
21          os.makedirs(self.output_dir, exist_ok=True)
22
23      def get_valid_actions(self, x, y):
24          """Returns a list of valid actions for a given state."""
25          valid_actions = []
26          for i, (dx, dy) in enumerate(self.actions):
27              nx_, ny_ = x + dx, y + dy
28              if 0 <= nx_ < self.nx and 0 <= ny_ < self.ny and self.maze[ny_, nx_] == 0:
29                  valid_actions.append(i)
30          return valid_actions
31
32      def get_reward(self, x, y):
33          """Returns the reward for a given state."""
34          return 100 if (x, y) == self.goal else -1
35
36      def is_goal(self, x, y):
37          """Checks if the given state is the goal."""
38          return (x, y) == self.goal
39
40  # Deep Q-Network (DQN) Model
41  class DQN(nn.Module):
42      def __init__(self, input_dim, output_dim):
43          super(DQN, self).__init__()
44          self.fc1 = nn.Linear(input_dim, 128)
45          self.fc2 = nn.Linear(128, 128)
46          self.fc3 = nn.Linear(128, output_dim)
47
48      def forward(self, x):
49          x = F.relu(self.fc1(x))
50          x = F.relu(self.fc2(x))
51          return self.fc3(x)
52
53  # DQN Agent
54  class DQNAgent:
55      def __init__(self, env, gamma=0.9, lr=0.001, batch_size=64, memory_size=50000):
56          self.env = env
57          self.gamma = gamma
58          self.batch_size = batch_size
59          self.epsilon = 1.0
60          self.epsilon_decay = 0.995
61          self.epsilon_min = 0.01
62          self.memory = deque(maxlen=memory_size)
63
64          self.model = DQN(2, len(env.actions))
```

```python
65          self.target_model = DQN(2, len(env.actions))
66          self.target_model.load_state_dict(self.model.state_dict())
67          self.optimizer = optim.Adam(self.model.parameters(), lr=lr)
68
69      def get_state(self, x, y):
70          """Encodes state as a tensor."""
71          return torch.tensor([x / self.env.nx, y / self.env.ny], dtype=torch.float32)
72
73      def choose_action(self, x, y):
74          """Epsilon-greedy policy."""
75          if random.random() < self.epsilon:
76              return random.choice(self.env.get_valid_actions(x, y))
77          else:
78              state = self.get_state(x, y).unsqueeze(0)
79              q_values = self.model(state)
80              valid_actions = self.env.get_valid_actions(x, y)
81              return max(valid_actions, key=lambda a: q_values[0, a].item())
82
83      def store_experience(self, state, action, reward, next_state, done):
84          """Stores experience in replay memory."""
85          self.memory.append((state, action, reward, next_state, done))
86
87      def train(self):
88          """Trains the model using replay memory."""
89          if len(self.memory) < self.batch_size:
90              return
91
92          batch = random.sample(self.memory, self.batch_size)
93          states, actions, rewards, next_states, dones = zip(*batch)
94
95          states = torch.stack(states)
96          next_states = torch.stack(next_states)
97          actions = torch.tensor(actions, dtype=torch.long)
98          rewards = torch.tensor(rewards, dtype=torch.float32)
99          dones = torch.tensor(dones, dtype=torch.float32)
100
101         q_values = self.model(states).gather(1, actions.unsqueeze(1)).squeeze(1)
102         next_q_values = self.target_model(next_states).max(1)[0].detach()
103         target_q_values = rewards + self.gamma * next_q_values * (1 - dones)
104
105         loss = F.mse_loss(q_values, target_q_values)
106         self.optimizer.zero_grad()
107         loss.backward()
108         self.optimizer.step()
109
110     def update_target_network(self):
111         """Updates target network weights."""
112         self.target_model.load_state_dict(self.model.state_dict())
113
114     def train_agent(self, num_episodes, update_target_every=50, frame_time=5.0, d_episode=50):
115         """Trains the DQN agent and creates a GIF."""
116         frames = []
117         for episode in range(num_episodes):
118             x, y = self.env.start
119             total_reward = 0
120             state = self.get_state(x, y)
121
122             while not self.env.is_goal(x, y):
123                 action = self.choose_action(x, y)
124                 dx, dy = self.env.actions[action]
125                 new_x, new_y = x + dx, y + dy
126
```

```
127                  reward = self.env.get_reward(new_x, new_y)
128                  next_state = self.get_state(new_x, new_y)
129                  done = self.env.is_goal(new_x, new_y)
130
131                  self.store_experience(state, action, reward, next_state, done)
132                  self.train()
133
134                  x, y = new_x, new_y
135                  state = next_state
136                  total_reward += reward
137
138              self.epsilon = max(self.epsilon * self.epsilon_decay, self.epsilon_min)
139
140              if episode % update_target_every == 0:
141                  self.update_target_network()
142
143              if episode % d_episode == 0:
144                  save_path = os.path.join(self.env.output_dir, f"policy_{episode}.png")
145                  self.visualize_policy(save_path)
146                  frames.append(imageio.imread(save_path))
147
148              print(f"Episode {episode}, Total Reward: {total_reward}, Epsilon: {self.epsilon:.3f}")
149
150          gif_path = os.path.join(self.env.output_dir, "policy_evolution.gif")
151          imageio.mimsave(gif_path, frames, duration=frame_time*num_episodes/d_episode)
152          print(f"GIF saved as {gif_path}")
153
154      def visualize_policy(self, save_path=None):
155          """Visualizes the current policy."""
156          fig, ax = plt.subplots(figsize=(10, 10))
157          ax.imshow(1 - self.env.maze, cmap="gray")
158
159          for y in range(self.env.ny):
160              for x in range(self.env.nx):
161                  if self.env.maze[y, x] == 0:
162                      best_action = max(self.env.get_valid_actions(x, y), key=lambda a: self.model(self.ge
163                      dx, dy = self.env.actions[best_action]
164                      ax.arrow(x, y, dx * 0.3, dy * 0.3, head_width=0.2, head_length=0.2, fc="red", ec="re
165
166          if save_path:
167              plt.savefig(save_path)
168              plt.close()
169          else:
170              plt.show()
171
172  # Main Execution
173  if __name__ == "__main__":
174      env = Environment("maze100.npy")
175      agent = DQNAgent(env)
176      agent.train_agent(num_episodes=5000)
177
178      # Save final policy visualization
179      final_policy_path = os.path.join(env.output_dir, "final_policy.png")
180      agent.visualize_policy(final_policy_path)
```