# Maze Simulator and Reinforcement Learning Agent Report

Shivam Chaubey

February 9, 2025

## 1 Introduction

This report details the implementation of a maze simulator and a reinforcement learning (RL) agent designed to navigate the generated maze. The maze is generated using a randomized depth-first search algorithm. The RL agent employs Q-learning to find a path from a designated start point to an end point within the maze. The primary goal is to develop an agent that can find a reasonably short path through various maze configurations.

## 2 Code Overview

The Python code is structured into two main components: the `MazeSimulator` class and the `mazeSolverUsingRL` function.

### 2.1 MazeSimulator Class

The `MazeSimulator` class is responsible for generating and representing the maze.

#### 2.1.1 __init__ Method

This method initializes the maze with a specified width and height, creating a grid filled with walls (represented by 1). It also initializes the start and end points to `None`.

#### 2.1.2 generate_maze Method

This method uses a depth-first search (DFS) algorithm to carve out paths in the maze. It starts from a random cell and recursively visits neighboring cells, marking them as part of the path (represented by 0). The `is_valid_move` helper function ensures that the DFS stays within the maze boundaries and only moves to unvisited cells.

#### 2.1.3 set_start_and_end Method

This method randomly selects two distinct open cells (cells with value 0) in the maze and designates them as the start and end points for the agent.

#### 2.1.4 get_maze_array Method

This method creates a 2D array representation of the maze, where 1 represents walls and 0 represents paths. It also marks the start and end points as 1 for representation purposes.

#### 2.1.5 save_maze_image Method

This method uses Matplotlib to visualize the maze and save it as a PNG image. The start and end points are colored differently for clarity.

### 2.2 mazeSolverUsingRL Function

This function implements a Q-learning algorithm to train an agent to navigate the maze.

### 2.2.1 Q-Learning Implementation

The Q-learning algorithm uses a Q-table to store the expected rewards for taking specific actions in different states. The agent iteratively explores the maze, updating the Q-table based on the rewards it receives.

**Hyperparameters:**

- `alpha` (Learning rate): Determines the extent to which newly acquired information overrides old information.

- `gamma` (Discount factor): Determines the importance of future rewards.

- `epsilon` (Exploration rate): Controls the balance between exploration (trying new actions) and exploitation (using known optimal actions).

- `num_episodes`: The number of training iterations.

**Q-Table Initialization:** The Q-table is a 3D NumPy array that stores the Q-values for each state-action pair. The dimensions correspond to the height, width, and number of possible actions (up, right, down, left).

**Action Mappings:** The agent can take four actions: move up, right, down, or left. These actions are represented as coordinate changes: `(-1, 0), (0, 1), (1, 0), (0, -1)`.

**Reward Function:** The reward function assigns rewards based on the agent's actions:

- Reaching the end point: +100

- Hitting a wall: -10

- Moving to an open cell: -1 (small penalty per step to encourage efficiency)

**Q-Learning Loop:** The Q-learning loop iterates through a specified number of episodes. In each episode, the agent starts from the start point and continues until it reaches the end point or hits a wall. The agent chooses an action based on an epsilon-greedy policy, which balances exploration and exploitation. The Q-table is updated using the Q-learning update rule.

### 2.2.2 Path Extraction

After the Q-table is trained, the function extracts the optimal path by starting at the start point and repeatedly choosing the action with the highest Q-value until it reaches the end point. A check is included to prevent the agent from getting stuck in loops.

# 3 Approach

The approach combines a standard maze generation algorithm (DFS) with a reinforcement learning technique (Q-learning) to solve the maze. This hybrid approach allows for the creation of complex mazes and the training of an agent to efficiently navigate them.

## 3.1 Maze Generation

The DFS algorithm ensures that the maze is fully connected, meaning there is a path from any cell to any other cell. The randomization aspect ensures that each generated maze is unique.

## 3.2 Reinforcement Learning

Q-learning is a model-free reinforcement learning algorithm that learns an optimal policy by iteratively updating a Q-table. The Q-table stores the expected rewards for taking specific actions in different states. The epsilon-greedy policy allows the agent to explore the maze initially and then gradually exploit its knowledge to find the shortest path.

## 4 Code Listing

```python
1  import random
2  import matplotlib.pyplot as plt
3  import numpy as np
4
5  # Do not change the maze generation
6
7  class MazeSimulator:
8      def __init__(self, width, height):
9          self.width = width
10         self.height = height
11         self.maze = [[1 for _ in range(width)] for _ in range(height)]  # 1 = wall, 0 =
   path
12         self.start = None
13         self.end = None
14
15     def generate_maze(self):
16         def is_valid_move(x, y):
17             return 0 <= x < self.height and 0 <= y < self.width and self.maze[x][y] == 1
18
19         def dfs(x, y):
20             self.maze[x][y] = 0
21             directions = [(0, 2), (2, 0), (0, -2), (-2, 0)]
22             random.shuffle(directions)
23
24             for dx, dy in directions:
25                 nx, ny = x + dx, y + dy
26                 if is_valid_move(nx, ny):
27                     self.maze[x + dx // 2][y + dy // 2] = 0  # Remove the wall
28                     dfs(nx, ny)
29
30         start_x, start_y = random.randrange(0, self.height, 2), random.randrange(0, self
   .width, 2)
31         dfs(start_x, start_y)
32
33     def set_start_and_end(self):
34         open_cells = [(x, y) for x in range(self.height) for y in range(self.width) if
   self.maze[x][y] == 0]
35         self.start, self.end = random.sample(open_cells, 2)
36
37     def get_maze_array(self):
38         maze_array = [[1 if self.maze[x][y] == 0 or (x, y) in [self.start, self.end]
   else 0 for y in range(self.width)] for x in range(self.height)]
39         return maze_array
40
41     def save_maze_image(self, filename="maze.png"):
42         maze_copy = [[2 if (x, y) == self.start else 3 if (x, y) == self.end else self.
   maze[x][y] for y in range(self.width)] for x in range(self.height)]
43         plt.figure(figsize=(10, 10))
44         plt.imshow(maze_copy, cmap="viridis", origin="upper")
45         plt.axis("off")
46         plt.savefig(filename)
47         plt.close()
48
49
50     def mazeSolverUsingRL(self):
51         """
52         Solve the maze using Q-learning.
53
54         The RL agent learns a policy for moving from the start to the end cell.
55         Actions: 0 = Up, 1 = Right, 2 = Down, 3 = Left.
56         Rewards:
57         - +100 for reaching the end.
58         - -10 for hitting a wall (invalid move).
59         - -1 per step to encourage shorter paths.
60
61         Returns:
62             path (list of tuple): The sequence of (row, col) states representing the
   optimal path.
63         """
64
65         # Hyperparameters
```

```
66          alpha = 0.1            # Learning rate
67          gamma = 0.9            # Discount factor
68          epsilon = 0.1         # Exploration rate
69          num_episodes = 100000  # Number of training episodes
70          max_steps = self.width * self.height * 4  # Maximum steps per episode (safeguard
    )
71
72          # Initialize the Q-table with zeros.
73          # Dimensions: [height][width][number_of_actions]
74          q_table = np.zeros((self.height, self.width, 4))
75
76          # Define the four possible actions as (dx, dy) movements.
77          actions = [(-1, 0),  # Up
78                     (0, 1),    # Right
79                     (1, 0),    # Down
80                     (0, -1)]   # Left
81
82          def is_valid_state(x, y):
83              """Return True if (x, y) is within bounds and is an open cell (path)."""
84              return 0 <= x < self.height and 0 <= y < self.width and self.maze[x][y] == 0
85
86          def choose_action(state):
87              """Choose an action using an epsilon-greedy strategy."""
88              if random.uniform(0, 1) < epsilon:
89                  return random.choice(range(4))  # Explore: choose a random action
90              else:
91                  return np.argmax(q_table[state[0], state[1]])  # Exploit: choose the
    best-known action
92
93          def get_reward(state):
94              """
95              Return the reward for moving into a given state.
96              - 100 if the state is the end.
97              - -10 if the state is invalid (e.g. a wall or out-of-bounds).
98              - -1 for any other (valid) step.
99              """
100             if state == self.end:
101                 return 100
102             if not is_valid_state(state[0], state[1]):
103                 return -10
104             return -1
105
106         # --------------------------
107         # Q-learning Training Loop
108         # --------------------------
109         for episode in range(num_episodes):
110             state = self.start
111             steps = 0
112             while steps < max_steps:
113                 # Choose an action from the current state
114                 action = choose_action(state)
115                 x, y = state
116                 dx, dy = actions[action]
117                 new_state = (x + dx, y + dy)
118                 reward = get_reward(new_state)
119
120                 # Determine the maximum Q-value for the new state (if valid)
121                 if is_valid_state(new_state[0], new_state[1]):
122                     best_next_q = np.max(q_table[new_state[0], new_state[1]])
123                 else:
124                     best_next_q = 0
125
126                 # Update the Q-value for the current state and action
127                 q_table[x, y, action] += alpha * (reward + gamma * best_next_q - q_table
    [x, y, action])
128
129                 # Terminate the episode if:
130                 # - The agent reaches the goal.
131                 # - The agent takes an invalid move (hits a wall).
132                 if new_state == self.end or not is_valid_state(new_state[0], new_state
    [1]):
133                     break
134
```

```
135              # Otherwise, move to the new state and continue
136              state = new_state
137              steps += 1
138
139         # -------------------------------
140         # Extract the Optimal Path
141         # -------------------------------
142         path = [self.start]
143         current_state = self.start
144         visited = set([self.start])  # To help prevent loops
145
146         # Follow the greedy policy derived from the Q-table until the end is reached.
147         while current_state != self.end:
148             action = np.argmax(q_table[current_state[0], current_state[1]])
149             dx, dy = actions[action]
150             new_state = (current_state[0] + dx, current_state[1] + dy)
151
152             # If the new state is invalid or we've already visited it, break out to
    avoid an infinite loop.
153             if not is_valid_state(new_state[0], new_state[1]) or new_state in visited:
154                 break
155
156             path.append(new_state)
157             visited.add(new_state)
158             current_state = new_state
159
160         return path
161
162
163 # Game flow
164 if __name__ == "__main__":
165     width, height = 21, 21
166
167     # Initialize simulator
168     simulator = MazeSimulator(width, height)
169     simulator.generate_maze()
170     simulator.set_start_and_end()
171
172     # Generate and display maze
173     maze_array = simulator.get_maze_array()
174     print("Maze array:")
175     for row in maze_array:
176         print(row)
177     print(f"Start point: {simulator.start}")
178     print(f"End point: {simulator.end}")
179
180     # Solve maze using RL
181     path = simulator.mazeSolverUsingRL()  # Call the solver
182     print(f"\nRL Agent Path ({len(path)} steps):")
183     print(path)
184
185     # Save visualization with path
186     simulator.save_maze_image("maze.png")
187     print("\nMaze image saved as 'maze.png'")
```

Listing 1: Maze Simulator Code

# 5 Conclusion

The implemented maze simulator and Q-learning agent provide a functional solution for generating and navigating mazes. The agent learns to find paths through the maze, demonstrating the effectiveness of reinforcement learning techniques in solving pathfinding problems. The performance of the agent can be further improved by fine-tuning the hyperparameters and exploring other reinforcement learning algorithms. Future work could also involve visualizing the agent's learning process and comparing the performance of different maze-solving algorithms.