

**Student ID : 2019310649**

### III. Visualization

Figure 1 displays the performance of the proposed algorithm across different learning rates and discount rates. The figure is organized into three main sections: Maze Navigation, Performance Metrics, and Reward Table.

**Maze Navigation:** The top row shows maze navigation results for four learning rates ( $\gamma = 0.2, 0.4, 0.6, 0.8$ ) and four discount rates ( $\alpha = 0.2, 0.4, 0.6, 0.8$ ). Each maze plot includes a green path and a red path, with a 'Path Length' value in the bottom right corner.

**Performance Metrics:** The middle row shows four line plots of 'Percent' vs 'Discount factor ( $\gamma$ )' for the same learning rates. The plots show that the performance is relatively stable across different discount factors, with a slight increase in performance as the discount factor increases.

**Reward Table:** The bottom row shows a 'Reward Table' heatmap for the same parameters. The color bar on the right indicates the reward values, ranging from -100 (dark blue) to 75 (light blue). The table shows that the reward is generally higher for higher learning rates and discount rates.

## I . Background

According to the wiki, Q-Learning is a model-free reinforcement learning technique used to find the optimal policy for a given finite Markov decision process (MDP). In an MDP, the agent interacts with the environment and learns the optimal policy by receiving feedback in the form of rewards.

An MDP follows the Markov Property, which means that the current state completely characterizes the state of the world, and the future is only concerned with the present, not the past. It is defined by a tuple  $(S, A, R, P, r)$ , where:

- $S$  is the set of states in the environment.
- $A$  is the set of actions that the agent can take.
- $R$  is the reward function that provides the immediate reward when an action is taken in a particular state.
- $P$  is the state transition probability function, which defines the probability of transitioning to a new state when an action is taken in a current state.
- $r$  is the discount factor that determines the importance of future rewards.

An MRP (Markov Reward Process) is a Markov chain augmented with rewards. The goal in an MRP is to maximize the return, which is a measure of the cumulative rewards obtained over time. The return is influenced by the discount rate, denoted as  $\gamma$ .

The discount rate serves two purposes. First, it helps avoid infinite returns in cyclic Markov processes by assigning less weight to future rewards as the time horizon extends. This prevents the return from growing unbounded and ensures convergence. Second, the discount rate reflects the degree of uncertainty about future rewards. A smaller discount rate values immediate rewards more and considers future rewards with less weight, indicating a more myopic decision-making approach. Conversely, a larger discount rate values future rewards more and takes a more forward-looking perspective.

The return in an MRP is computed by summing the discounted rewards over time. The formula for the return is:

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Here, the discount rate determines the relative importance of immediate rewards versus future rewards. It allows for trade-offs between immediate gains and long-term benefits, reflecting the time preference of the decision-maker.

By maximizing the return in an MRP, we aim to find policies that achieve the highest cumulative rewards over time, considering both the immediate rewards and the potential future rewards discounted by  $\gamma$ .

The Bellman equation for an MRP expresses the relationship between the value function of the current state and the value function of the next state.

$$\begin{aligned}
 v(s) &= \mathbb{E}[G_t \mid S_t = s] \\
 &= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s] \\
 &= \mathbb{E}[R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \dots) \mid S_t = s] \\
 &= \mathbb{E}[R_{t+1} + \gamma G_{t+1} \mid S_t = s] \\
 &= \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s] \longrightarrow v(s) = R_s + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'} v(s')
 \end{aligned}$$

In an MDP, which is an extension of MRP, actions to be taken by agents are included. The value function of an MDP is expressed as follows:

- ▶ The state-value function  $v_\pi(s)$  of an MDP is the expected return starting from state  $s$ , and then following policy  $\pi$ .

$$v_\pi(s) = \mathbb{E}_\pi[G_t \mid S_t = s]$$

- ▶ The action-value function  $q_\pi(s, a)$  is the expected return starting from state  $s$ , taking action  $a$ , and then following policy  $\pi$ .

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a]$$

To maximize the return in an MDP, Q-Learning aims to find the optimal policy by maximizing the Q-value function  $Q(s, a)$ , which represents the expected cumulative reward when taking action  $a$  in state  $s$ . The Bellman Optimality Equation provides the theoretical foundation for Q-Learning and is given by:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[ r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right]$$

Here,  $Q^*$  represents the optimal Q-value function,  $R(s, a)$  is the immediate reward,  $\gamma$  is the discount factor,  $P(s' \mid s, a)$  is the transition probability, and  $\max Q^*(s', a')$  represents the maximum Q-value over all possible actions in the next state.

The Q-Learning algorithm updates the Q-values iteratively based on the observed rewards and transitions. The Q-Learning update equation is:

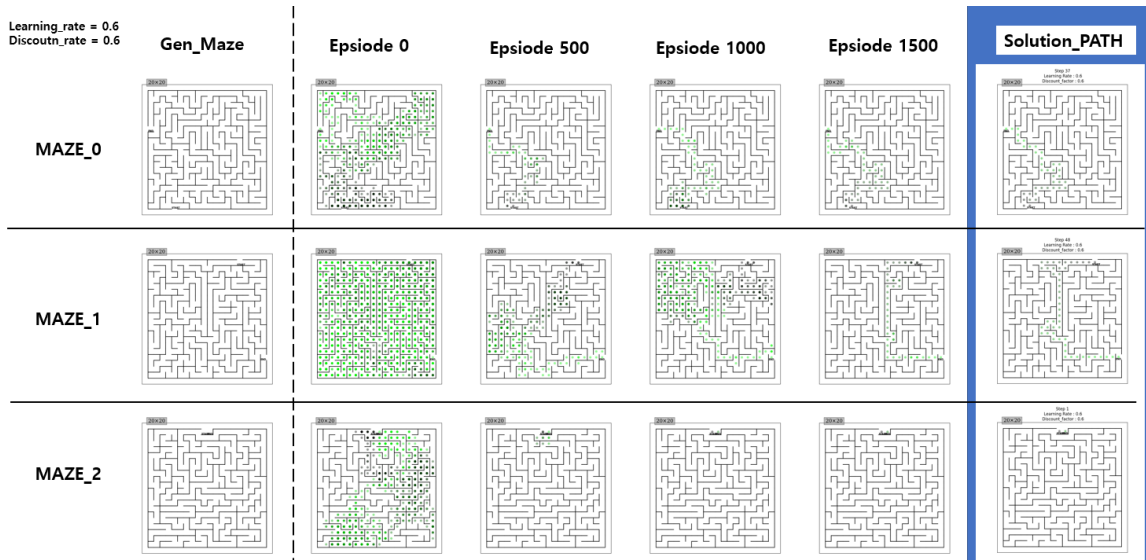
$$Q(s_t, a_t) \leftarrow (1 - \alpha) \cdot \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \left( \underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} \right)$$

Here,  $Q(s, a)$  is the Q-value for state  $s$  and action  $a$ ,  $\alpha$  is the learning rate that determines the weight given to new information,  $R(s, a)$  is the immediate reward,  $\gamma$  is the discount factor, and  $\max Q(s', a')$  represents the maximum Q-value over all possible actions in the next state.

By applying Q-Learning iteratively, the agent can learn the optimal Q-values for each state-action pair and ultimately determine the optimal policy for

## II. Result and Analysis

Let's examine the progression of the Episode by referring to the following image.



Let me analyze the picture above.

First of all, since the Exploration Rate is set high in the early stages of low Episode and all Q\_Table are similar in initialization state, the Agent updates the Q\_Table as it explores most of the space. An episode ends when an agent randomly accidentally arrives at the destination. And as Episode progresses, the agent's Exploration Rate gradually decreases, so it tends to rely on Q\_Table (dependent on experience). For reference, as I will explain later, in the code I made, if there are multiple same Max values when walking on the street based on Q\_Table, it is set to go in a random direction among the directions that are max values. Anyway, as Episode progresses, the agent updates the Q\_Table and becomes increasingly dependent on experience.

There is an important point here. In fact, after a long period of code analysis and research, the agent was observed to move back and forth (Ex, the agent moves left and right indefinitely), and as a result of checking the Q\_Table in this case, the same Left and Right existed at the same value (for example, -0.166 or -0.125). To understand this reason, let's look at the equation of Q-Learning.

$$Q(s_t, a_t) \leftarrow (1 - \alpha) \cdot \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \left( \underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} \right)$$

learned value

Notice the inside of the right parenthesis in the above equation.

As explained in Background, Q-Learning gradually updates the Q\_table of the final destination's Reward in the form of Backpropagation, and the final destination's Reward affects the way to the

final destination, increasing its value. However, if the Discount Rate is very low, it converges to zero of the impact of Reward at the final destination, and most of the values converge to a specific value. The value is the value of zero inside the right bracket in the above equation.

Therefore, if the maze is long from the start position to the end position, the smaller the Discount Rate, the smaller the impact on the final destination, so that most Q\_Table around the start position converges to the same value and repeatedly moves back and forth in the space!

As a way to solve this problem, it was possible to solve the problem by significantly increasing the Discount Rate.

Now let's take a closer look at the picture above. As shown in the figure above, it can be seen that there is less tendency to pierce walls well in the beginning, and the closer to the final destination, the greater the tendency to pierce walls.

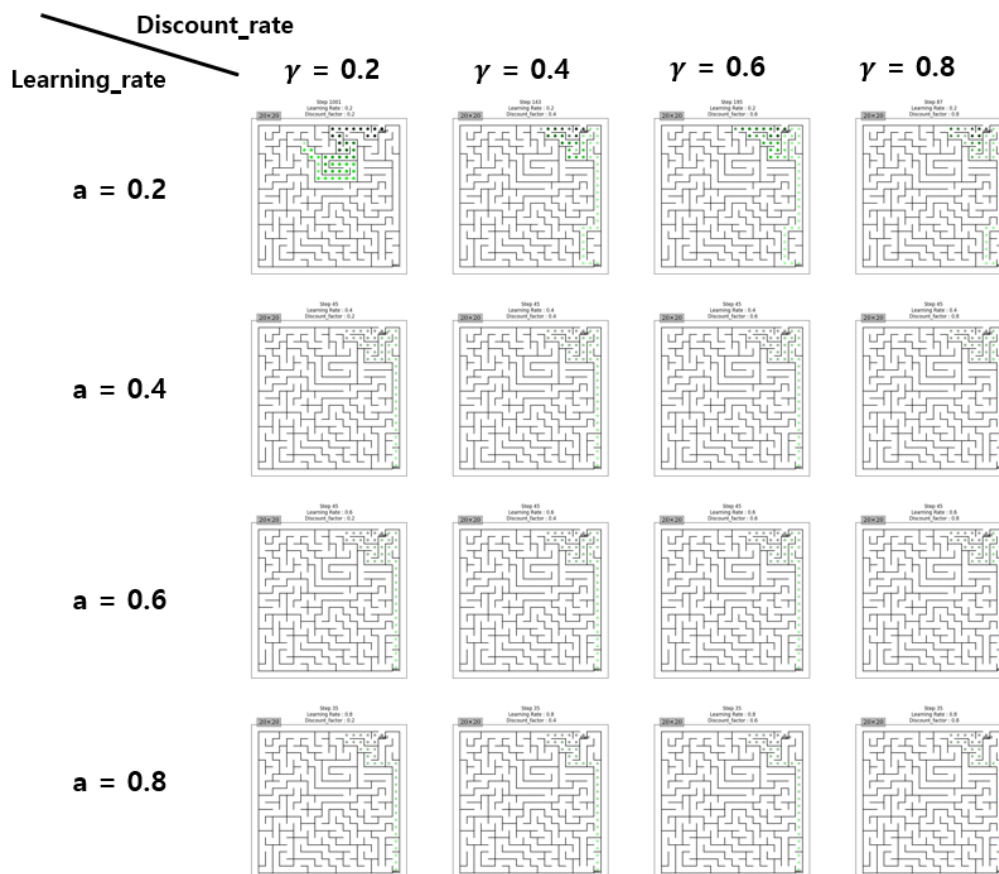
This is due to the fact that when Reward reaches its final destination, +100, -1 when drilling a wall, and -0.1 when just moving, the closer it is to the final destination, the less Discount multiplied, so even if the wall is pierced and reaches the final destination, it has a significant impact on the update of the Q-Table. Therefore, this trend becomes stronger toward the second half, so the tendency is written in the form of a compensation in which the updated value of the Q-Table breaks through a large wall (-1) rather than going back along the road. In the vicinity of the Start position, the effect of the final reward +100 has become very small because the Discount Rate has been multiplied several times, so the Reward (-1) when the wall is drilled is operated larger. The following is part of the Q-Table.

(15, 5 )	-1.1498	-0.2500	-1.1500	-0.2500
(15, 6 )	-1.1499	-0.2500	-0.2500	-1.1499
(15, 7 )	-0.2500	-1.1500	-0.2500	-1.1498
(15, 8 )	-0.2500	-1.1500	-0.2500	-1.1500
(15, 9 )	-0.2500	-1.1500	-1.1499	-0.2500

As you can see in the figure, there is a lot of -1.15. Also, interestingly, if it's not -1.15, you can see a difference of -0.25, i.e. when it's barely pierced the wall ( $-1 + 0.1 = 0.9$ ).

For this reason, it can be seen that less wall drilling is performed around the starting position, and the tendency to wall drilling becomes stronger as we reach the final destination.

For a maze, let's observe what happens when you sweep the Learning Rate and Discount Rate.



As shown in the figure above, the smaller the learning\_rate(a), the smaller the size of the Q\_table is, the smaller the size of the Q\_table's update, so even if the final result is reached and Reward (100), it eventually converges to a specific value as Episode is repeated.

So we can see that 0.2 or 0.4 where a is very small does not reach the final destination.

Because the above maze problem is difficult (because the distance is long from the start position to the target position) it is difficult to learn with these small learning rates.(The update size is small, so the final reward cannot affect and eventually converges all around the starting position, resulting in the Agent repeating only the same position) Note that most cases , it will repeat about 5 spaces or less from the starting position, but in my code, it is designed to behave randomly.

Now let's understand the impact of the Discover Rate. I explained the impact of Discount above, and it is the same. While the Learning Rate governs the overall size of updates, the Discount Rate more directly affects the reduction of final rewards. Smaller Discount Rate converges faster, so you become shortsighted, and the larger the Discount Rate, the more far you can see. Therefore, in order not to face the same converging problem, the Discount Rate must be set according to the tendency of the maze to some extent.

### III. Code Description

It is a main code that performs TASK 1 to 5 as required in the task. It is important to note that Task 1 to 4 made three Maze and compared each Maze, and Task 5 made another Maze to confirm each tendency as the Learning Rate and Discount Rate were changed.

```

1  from __future__ import absolute_import
2  import os
3  import sys
4  import shutil
5
6  from matplotlib import colors, pyplot as plt
7  import numpy as np
8  sys.path.append(os.path.abspath(os.path.join(os.path.dirname(__file__), '..')))
9  from src.maze import Maze
10 from src.maze_manager import MazeManager
11 from src.debug_viz import DebugViz
12
13 #####
14 # Task 1 & 2 | Create Function about Q-Learning : q_learning & q_learning_path #
15 #####
16 from src.q_learning_manager import QLearningManager
17 print("=====")
18 print("    Task 1,2 Complete")
19 print("=====")
20
21 shutil.rmtree('./temp')
22 shutil.rmtree('./output')
23
24 #####
25 # Task 3 | Train 3 randomly generated 20x20 mazes with Visualization #
26 #####
27 maze_row, maze_col = 20, 20
28 manager = MazeManager()
29 maze_list = []
30 solver_list = []
31
32 for i in range(3):
33     maze = manager.add_maze(maze_row, maze_col)
34     maze_list.append(maze)
35     solver = QLearningManager(maze)
36     debugViz = DebugViz()
37     debugViz.set_filename("maze{}_generation_20x20".format(i))
38     solver_list.append(solver)
39     debugViz.show_maze(maze, display_mode=False)
40
41 # Define the image file names and paths
42 manager.combine_images('temp', 'generation', 'maze_output1.png', 'down')
43
44 print("=====")
45 print("    Task 3 Complete")
46 print("=====")
47
48 #####
49 # Task 4 | Implement the Q-Learning based maze traversal algorithm for the randomly generated three mazes #
50 #####
51 solution_path_list = []
52 solution_cost_list = []
53 solution_q_table = []
54
55 debugViz = DebugViz()
56 for i in range(3):
57     solver_list[i].q_learning(maze_list[i], episodes = 1500, learning_rate = 0.6, discount_factor = 0.6, exploration_rate = 0.95)
58     manager.combine_images('temp', 'episode', 'maze{}_episode_solution.png'.format(i), 'right')
59     solution_q_table.append(solver_list[i].q_table)
60     maze_list[i].solution_path, maze_list[i].solution_cost, maze_list[i].solution_reward = solver_list[i].q_learning_path(maze_list[i], solution_q_table[i])
61     debugViz.show_solution(maze_list[i], display_key = False, media_name="maze{}".format(i), learning_rate = 0.6, discount_factor = 0.6)
62     solution_path_list.append(maze_list[i].solution_path)
63     solution_cost_list.append(maze_list[i].solution_cost)
64
65 manager.combine_images('temp', 'solution', 'maze_q_learning_solution.png', 'down')
66 manager.combine_images('output', '_episode_', 'maze_output1.png', 'down')
67 manager.combine_images('output', '_output_', 'maze_episode_table.png', 'right')
68
69 print("=====")
70 print("    Task 4 Complete")
71 print("=====")
72 # 5. Use Visualization Funct (at Codebase). Demonstrates the algorithm's performance from the 3 randomly gen. mazes
73 print("*****Show Solution Animation*****")
74 # Run the solution animation
75 for i in range(3):
76     debugViz.show_solution_animation(maze_list[i])
77
78 #####
79 # Task 5 | Choose at least two param and change values of them to see the impact of the parameters. #
80 # e.g.,  $\alpha$  (0.2, 0.4, 0.6, 0.8) and  $\gamma$  (0.2, 0.4, 0.6, 0.8) #
81 #####
82
83 shutil.rmtree('./temp')
84 solution2_path_list = []
85 solution2_cost_list = []
86 solution2_q_table = []
87 solution2_reward = []
88 solver_list = []
89
90 a_list = [0.2, 0.4, 0.6, 0.8]
91 r_list = [0.2, 0.4, 0.6, 0.8]
92
93 manager_param = MazeManager()
94 maze_param = manager_param.add_maze(maze_row, maze_col)
95 debugViz_param = DebugViz()
96 manager_param.show_maze(maze_param.id)
97
98 i = 0
99 j = 0
100 for a in a_list:
101     solution2_path_list.append([])
102     solution2_cost_list.append([])
103     solution2_reward.append([])
104     for r in r_list:
105         # Make Solver Manager
106         solver_param = QLearningManager(maze_param)
107         solver_list.append(solver_param)
108
109         # Make Q-Table
110         solver_list[i].q_learning(maze_param, episodes = 1500, learning_rate = a, discount_factor = r, exploration_rate = 0.95)
111         solution2_q_table.append(solver_list[i].q_table)
112
113         # Solve Maze Using Learned Q-Table
114         temp_solution_path, temp_solution_cost, temp_solution_reward = solver_list[i].q_learning_path(maze_param, solution2_q_table[i])
115         solution2_path_list[i].append(temp_solution_path)
116         solution2_cost_list[i].append(temp_solution_cost)
117         solution2_reward[i].append(temp_solution_reward)
118         debugViz_param.show_solution_animation(maze_param, display_key = False, media_name="maze{}_r_{}_out".format(a, r), learning_rate = a, discount_factor = r)
119         j+=1
120
121 manager_param.combine_images('temp', '_out_', 'maze_{}_solution.png'.format(a), 'right')
122 i+=1
123
124 manager_param.combine_images('output', 'maze_sb_', 'maze_a_r_solution_table.png', 'down')

```

From now on, the Q-Learning Path, which shows the learning code and solution path of Q-Learning, is described as follows.

```
def q_learning(self, maze, episodes, learning_rate, discount_factor, exploration_rate) :
```

```
54 def q_learning(self, maze, episodes, learning_rate, discount_factor, exploration_rate) :
55     # A. Initialize the Q-Table with zeros and set values of parameters by yourself.
56     self.q_table = self.initialize_q_table(maze)
57     maze.solution_path = None
58     maze.solution_cost = None
59     self.solution_reward = None
60     print(" Take some minute ... Please Wait ...")
61     action_list = ["LEFT", "UP", "RIGHT", "DOWN"]
62     for episode in range(episodes):
63         temp_episode_path = list()
64         steps = 0
65         epi_cost = 1
66         total_reward = 0
67         self.before_action = None
68         #predefined Maximum number of iteration
69         max_iter = 10000
70         #print for debug
71         if((episode+1) % (episodes//3) == 0):
72             print("=====")
73             print("Episode : {0:4}, max_iter: {1:5}, a: {2:2.1f} r: {3:2.1f} e: {4:3.2f} | Start ...".format(episode+1, max_iter, learning_rate, discount_factor, exploration_rate))
74             # B. Set the Initial State (the starting position of the agent).
75             current_state = maze.entry_coord
76             # G. Repeat Steps c-f until the agent reaches the goal or a predefined maximum number of iterations is reached.
77             while (current_state != maze.exit_coord):
78                 if(max_iter == 0 and episode > 100):
79                     print("Episode {} is Ended Because of Max_iteration set".format(episode))
80                     print("=====")
81                     break
82                     max_iter -= 1
83                     steps += 1
84                     temp_episode_path.append((current_state, False))
85                     # C. Choose an action based on the current state and the Q-Table, using the e-greedy strategy.
86                     action = self.choose_action(current_state, exploration_rate, self.q_table) #Return ex) LEFT
87                     # D. Perform the action and observe the reward and the next state.
88                     next_state = self.do_action(maze, current_state, action)
89                     reward = self.get_reward(maze, current_state, next_state)
90                     action_idx = self.find_word_idx(action, action_list)
91                     self.q_table[current_state][action_idx] += learning_rate * (reward + discount_factor * max(self.q_table[next_state]) - self.q_table[current_state][action_idx])
92                     total_reward += reward
93                     # F. Set the next state as the current state.
94                     current_state = next_state
95                     epi_cost += 1
96                     temp_episode_path.append((current_state, False))
97                     # Reduce the exploration rate over time.
98                     exploration_rate = exploration_rate * 0.999
99                     #-----this is for Debugging-----
100                     #print log
101                     self.log_iteration_info(episode, episodes, epi_cost, total_reward, max_iter)
102                     if (episode == 0 or (episode + 1) % (episodes//3) == 0):
103                         maze.solution_path = temp_episode_path
104                         debugViz = DebugViz()
105                         debugViz.show_solution(maze, episodes = episode, display_key = False, media_name="solve_maze_episode{episode+1}")
106                         output_folder = os.path.join(os.path.dirname(os.path.abspath(__file__)), '..', 'temp') # output 폴더 생성
107                         if not os.path.exists(output_folder):
108                             os.makedirs(output_folder)
109                         filename = os.path.join(output_folder, f"Q_Table_of_Episode{episode+1}.txt") # 파일 생성 및 저장
110                         debugViz.save_q_table(self.q_table, filename)
111                     print("=====")
```

The above code is the learning code of Q-Learning implemented by the method required by the task. As can be seen in the figure, the method required in the task was satisfied. In particular, the Q\_table update code of Q-Learning is as follows.

```
self.q_table[current_state][action_idx]
+= learning_rate * (reward + discount_factor * max(self.q_table[next_state]) - self.q_table[current_state][action_idx])
```

It can be seen that the requirements of the assignment were met.

Debugging messages are printed through Print to make it easier for users to understand, and Choose Action, Do Action, and Get Reward functions are called and used. This will be explained later. In the case of the Exploration Rate, it was gradually reduced as the Episode progressed, and through this, as the Episode progressed, it was searched based on experience. Max\_iteration was set to a very large number of 10000, which was set in this way to reach the final destination because of the nature of Q-Learning, the Reward of the final destination is transmitted to the starting position. Therefore, in the case of the initial 100 Episode, the agent continued to move unless the final destination was reached.



```

29 # =====
30 def get_reward(self, maze, current_state, next_state):
31     if next_state == maze.exit_coord:
32         reward = 100
33     elif maze.is_wall(next_state, current_state):
34         reward = -1
35     else:
36         reward = -0.1
37     return reward
38

```

As shown on the left, the compensation function returns 100 when reaching the final destination, -1 when hit by a wall, and -0.1 when just moving.

```

208 def is_wall(self, next_state, current_state):
209
210     # If Agent go Out side Maze
211     if(next_state == current_state):
212         return True
213
214     next_k, next_i = next_state
215     curr_k, curr_i = current_state
216
217     # action == "LEFT":
218     if(next_i == curr_i - 1):
219         action = "left"
220     # action == "UP":
221     if(next_k == curr_k + 1):
222         action = "bottom"
223     # action == "RIGHT":
224     if(next_i == curr_i + 1):
225         action = "right"
226     # action == "DOWN":
227     if(next_k == curr_k - 1):
228         action = "top"
229
230     # print(current_state)
231     # print("action: " + action)
232     # print("is wall: " + str(self.grid[curr_k][curr_i].walls))
233     # print(self.grid[curr_k][curr_i].walls[action])
234     # input("")
235
236     return self.grid[curr_k][curr_i].walls[action]
237

```

Earlier, the compensation function used the is\_wall function to determine whether it was a wall, and this was determined through the internal walls variable of the grid class. There was a problem in which Top and Bottom were matched in reverse in a given code, and the is\_wall function solved the problem by writing it in reverse. In the is\_wall function, the action is determined based on the current position and the next position to determine whether it is a wall.

What is important here is that in the Do\_action function, the code was designed so that the next position of the agent would return to its current position if the agent acted to go out of the labyrinth, which prevented the agent from exiting the labyrinth. Here, we can see that the is\_wall function returns True if the next position and the current position are the same, so when Get Reward, it returns the action of trying to go out of the maze as the Reward of the action of hitting the wall.

```

362 def initialize_q_table(self, maze):
363     q_table = {}
364     for row in range(maze.num_rows):
365         for col in range(maze.num_cols):
366             state = (row, col)
367             q_table[state] = {0, 0, 0, 0} # Actions : (LEFT, RIGHT, UP, DOWN)
368     return q_table
369
370 def find_argmax(self, list):
371     if not list:
372         print("No list ERROR at find_argmax")
373         return None # handle empty list
374     max_value = list[0]
375     max_index = 0
376     for i, value in enumerate(list):
377         if value > max_value:
378             max_value = value
379             max_index = i
380     return max_index
381
382 def find_word_idx(self, word, list):
383     for idx in range(len(list)):
384         if word == list[idx]:
385             return idx
386     print("Can't find word IDX ERROR")
387     return None
388
389 def _random_argmax_word(self, action_q_list):
390     max_value = max(action_q_list)
391     max_indices = [i for i, value in enumerate(action_q_list) if value == max_value]
392     selected_index = random.choice(max_indices)
393     action_list = ["LEFT", "UP", "RIGHT", "DOWN"]
394     return action_list[selected_index]
395

```

The left code is the base code used by default.

The Find\_Argmax function finds the index of the Max value in the list.

Find\_word\_idx is used when you want to find the index number of the Action, "Left", in a list containing the Action.

The \_random\_argmax\_word is a function that returns a random action among the same values if the MAX value is the same among the Q values for the action at any position in the Q\_Table. Based on this, when converging to the same value, it goes back and forth from side to side or up and down It is designed to avoid moving situations. If you are lucky to move randomly like this and get to where the final reward affects you, you can follow the MAX Q and move to the final position.

```

397 def _choose_action(self, current_coors, exploration_rate, q_table):
398     ==
399     Actions : (LEFT, UP, RIGHT, DOWN)
400     max
401     # Using e-greedy strategy
402     # Can't go before state
403
404     action_q_list = q_table[current_coors]
405     action_list = ["LEFT", "UP", "RIGHT", "DOWN"]
406     if random.uniform(0, 1) < exploration_rate:
407         random.shuffle(action_list)
408         explora_act = action_list[0]
409     else:
410         explora_act = self._random_argmax_word(action_q_list)
411
412     return explora_act # (LEFT, UP, RIGHT, DOWN)
413
414 def do_action(self, maze, current_state, action):
415     curr_k, curr_i = current_state
416
417     if action == "LEFT":
418         next_k = curr_k
419         next_i = curr_i - 1
420     elif action == "UP":
421         next_k = curr_k + 1
422         next_i = curr_i
423     elif action == "RIGHT":
424         next_k = curr_k
425         next_i = curr_i + 1
426     elif action == "DOWN":
427         next_k = curr_k - 1
428         next_i = curr_i
429     else:
430         raise ValueError("Invalid action value: {}".format(action))
431
432     if(next_k < 0 or next_i < 0 or next_k >= maze.num_rows or next_i >= maze.num_cols):
433         next_k = curr_k
434         next_i = curr_i
435         return("Shouldn't act : Can't go Outside Maze")
436
437     next_state = (next_k, next_i)
438
439     return next_state
440

```

The \_choose\_action function implements the E-Greedy strategy required in the task. Based on this, the agent experiences a direction other than the direction in which the Q value is maximum, creating a better Q\_Table, and finding the optimal movement path while enjoying several places. As you can see from the code, the reciprocating problem was solved by using the above-described \_random\_argmax\_word when it was not in an exploration situation.

The do\_action function is a function that, given an action, simply moves the agent in that direction and returns information about the next location.

```
def q_learning_path(self, maze, learned_q_table):
```

		Solution PATH Q-TABLE				
		State	Left	UP	Right	Down
124	def q_learning_path(self, maze, learned_q_table):	(2, 19)	-0.4846	-0.4982	-1.3877	-1.3760
125	# Set the start and goal positions	(2, 18)	-1.3268	-1.3950	-0.4877	-0.4888
126	start_pos = maze.entry_coord	(1, 18)	-1.3414	-0.4846	-0.4760	-1.3625
127	goal_pos = maze.exit_coord	(1, 19)	-0.4888	-1.3877	-1.3760	-0.4780
128	# Set the initial position	(0, 19)	-0.4625	-0.4760	-1.3780	-1.3780
129	curr_pos = start_pos	(0, 18)	-0.4531	-1.3888	-0.4780	-1.3625
130	self.before_action = None	(0, 17)	-0.4625	-0.4414	-0.4625	-1.3531
131	solution_path = []	(1, 17)	-1.3888	-0.4268	-1.3888	-0.4531
132	solution_cost = 0	(2, 16)	-0.4888	-1.3885	-1.3846	-0.4414
133	solution_reward = 0	(2, 15)	-1.2388	-0.3570	-0.4885	-1.3760
134	action_list = ["LEFT", "UP", "RIGHT", "DOWN"]	(3, 15)	-1.1985	-0.3213	-1.2856	-0.3856
135	max_cnt = 0	(4, 15)	-0.2766	-1.2198	-0.3570	-0.3570
136	print("=====")	(4, 14)	-0.2287	-1.2558	-0.3213	-1.1985
137	print("Solution PATH Q-TABLE")	(4, 13)	-0.1580	-1.3262	-0.2766	-1.1481
138	print("State   Left   UP   Right   Down  ")	(4, 12)	-0.6313	-0.0637	-0.2287	-0.9080
139	print("=====")	(5, 12)	0.0454	-0.8546	-1.3262	-0.1580
140	# Move until reaching the goal	(5, 11)	-0.1989	0.1818	-0.0637	-0.6313
141	while curr_pos != goal_pos:	(6, 11)	0.1813	0.3522	0.0454	0.0454
142	if max_cnt > 1000:	(7, 11)	1.3767	1.6903	-0.7182	-0.7182
143	print("Path fail")	(8, 11)	2.2379	1.2227	0.2785	1.2522
144	break	(8, 10)	2.0223	2.9223	1.6903	0.4761
145	max_cnt += 1	(9, 10)	3.7779	3.5979	1.2227	2.2379
146	# Move to the next position	(9, 9)	4.8474	4.8474	2.9223	2.9223
147	action = self.choose_action(curr_pos, 0, learned_q_table)	(10, 9)	7.3893	6.2295	4.4979	2.8779
148	# Format the Q-table values and highlight the selected action	(10, 8)	8.5856	9.2616	5.7474	5.7474
149	formatted_values = [	(11, 8)	11.7020	11.4212	6.2295	7.3893
150	f"{(3531+value/10.46)}{033[0]} if action_list[index] == action else f'{(value/10.46)}' for index, value in enumerate(learned_q_table[curr_pos])	(11, 7)	15.8774	15.3814	8.3616	8.5856
151	]	(11, 6)	19.6317	19.9718	12.6020	10.9820
152	print(f"({0:2})   {1:2})   {2:2})   {3:2})   {4:2})   {5:2})  ")	(12, 6)	23.7647	25.0898	14.4014	15.8774
153	curr_pos[0], curr_pos[1], formatted_values[0], formatted_values[1], formatted_values[2], formatted_values[3]]	(13, 6)	31.4872	19.3077	18.3518	19.9718
154	solution_path.append((curr_pos[0], curr_pos[1]), False) # Append current cell to total search path	(13, 5)	39.4840	24.4847	25.0898	23.7647
155	next_pos = self.do_action(maze, curr_pos, action)	(13, 4)	49.4880	30.9558	31.4872	31.8558
156	solution_reward += self.get_reward(maze, curr_pos, next_pos)	(13, 3)	63.1080	39.0448	38.5840	39.0448
157	before_action = action	(13, 2)	79.0080	50.0560	50.3800	50.9560
158	curr_pos = next_pos	(13, 1)	100.0000	63.8200	62.2000	63.8200
159	solution_cost += 1					
160	print("=====")					
161	solution_path.append((curr_pos[0], curr_pos[1]), False)					
162	#=====Base Solution path=====					
163	temp_episode_path = solution_path					
164	maze.solution_path = temp_episode_path					
165	return solution_path, solution_cost, solution_reward					

The above figure is a function to find solution\_Path based on the Q\_Table required in the task. As you can see in the figure, Print was added for debugging, which is printed as shown in the left picture. Based on this, it is easy to understand what value of the Q\_Table, in which direction the Agent moves, and what the next state has become, and debugging is easy.

The important thing is to find the solution path based only on Q\_TABLE, so by putting 0 in the Exploration Factor of the \_Choose\_Action function, the agent always relied on Q\_TABLE.

It's not important for the rest of the visualization code, so I'll just attach the code and skip it.

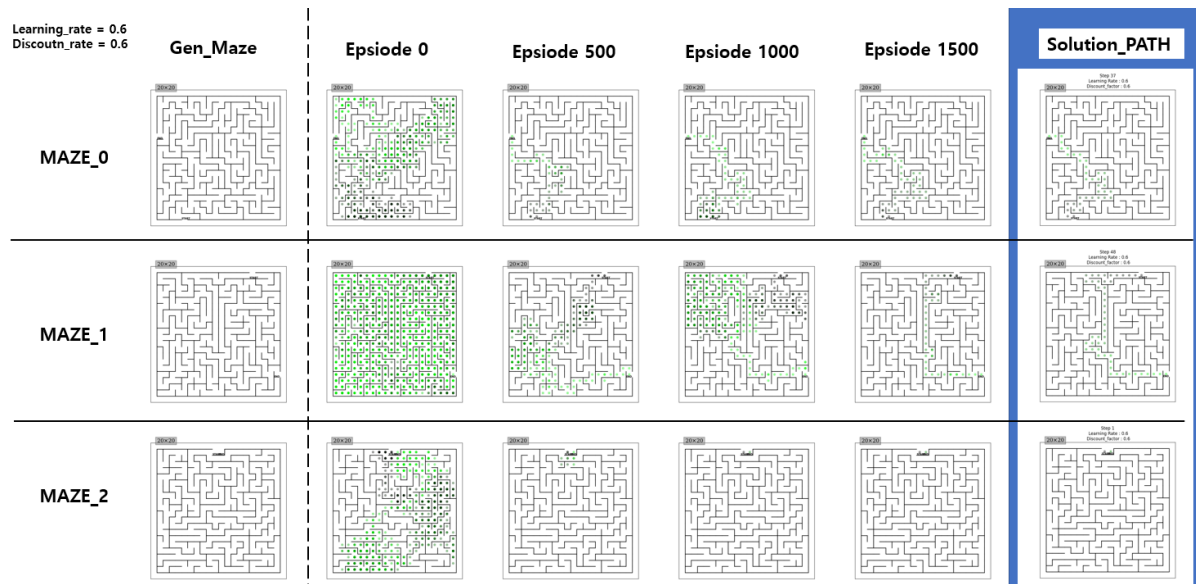
```

1 # Importing Libraries
2 import sys
3 import os
4 from gym import Env
5 from gym.spaces import Discrete, Box
6 from gym.wrappers import Monitor
7 from gym.utils import seeding
8 from gym.envs.registration import register
9
10 # Registering the environment
11 register(
12     id='Maze-v0',
13     entry_point='MazeEnv:MazeEnv',
14     kwargs={'maze': 'Maze', 'seed': 0},
15 )
16
17 # Importing the environment
18 from MazeEnv import MazeEnv
19
20 # Creating the environment
21 env = MazeEnv('Maze', 'Maze', 'seed')
22
23 # Defining the Q-table
24 def create_q_table(self, filename):
25     """
26     Create a Q-table for the maze environment.
27     """
28     with open(filename, 'w') as f:
29         f.write("Q-table\n")
30         f.write("-----\n")
31         f.write("State | Left | UP | Right | Down |")
32         f.write("-----\n")
33         for state, action_values in self.q_table.items():
34             f.write(f"({0:2}) | {1:2}) | {2:2}) | {3:2}) | {4:2}) | {5:2}) |")
35             f.write(f"{(3531+value/10.46)}{033[0]} if action_list[index] == action else f'{(value/10.46)}' for index, value in enumerate(action_values)}")
36             f.write("\n")
37
38 # Defining the Q-table
39 def print_q_table(self, q_table):
40     """
41     Print the Q-table.
42     """
43     print("Q-table")
44     print("-----")
45     print("State | Left | UP | Right | Down |")
46     print("-----")
47     for state, action_values in q_table.items():
48         print(f"({0:2}) | {1:2}) | {2:2}) | {3:2}) | {4:2}) | {5:2}) |")
49         print(f"{(3531+value/10.46)}{033[0]} if action_list[index] == action else f'{(value/10.46)}' for index, value in enumerate(action_values)}")
50         print("\n")
51
52 # Defining the solution
53 def show_solution(self, maze, episode, display_key = False, media_name=None, learning_rate = None, discount_factor = None, cell_size=1):
54     """
55     Show the solution of the maze.
56     """
57     viz = VizDumper(maze, cell_size, self.media_name)
58     viz.show_solution(episode, display_key, learning_rate, discount_factor)
59
60 # Defining the solution
61 def show_solution_animation(self, maze, media_name = None, cell_size=1):
62     """
63     Show the animation of the path that the solver took.
64     """
65     viz = VizDumper(maze, cell_size, self.media_name)
66     viz.show_solution_animation()
67
68 # Defining the solution
69 def set_filename(self, filename):
70     """
71     Set the filename for saving the solution and images.
72     """
73     self.filename = filename
74
75 # Defining the solution
76 def main():
77     """
78     Main function.
79     """
80     env = MazeEnv('Maze', 'Maze', 'seed')
81     create_q_table(env, 'q_table.txt')
82     print_q_table(env, env.q_table)
83     show_solution(env, 1, display_key = False, media_name=None, learning_rate = None, discount_factor = None, cell_size=1)
84     show_solution_animation(env, 'Maze', media_name = None, cell_size=1)
85     set_filename(env, 'Maze')
86
87 if __name__ == '__main__':
88     main()

```

## IV. Visualization

### TASK 1~4



### TASK 5

