



## Assignment 4

Mahdi Tabatabaei 400101515  
Github [Repository](#)

**Neuroscience of Learning, Memory, and Cognition**

Dr. Aghajan

January 29, 2025

## — Contents

|                          |           |
|--------------------------|-----------|
| <b>Contents</b>          | <b>1</b>  |
| <b>Maze Generation</b>   | <b>2</b>  |
| <b>Q-Learning</b>        | <b>4</b>  |
| Environment . . . . .    | 4         |
| Agent . . . . .          | 4         |
| <b>Deep Q-Learning</b>   | <b>8</b>  |
| Environment . . . . .    | 9         |
| Agent . . . . .          | 11        |
| Implementation . . . . . | 13        |
| <b>100 × 100 Maze</b>    | <b>18</b> |

## Maze Generation

In this part, we use following code to generate the maze:

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

Our  $20 \times 20$  generated maze is:

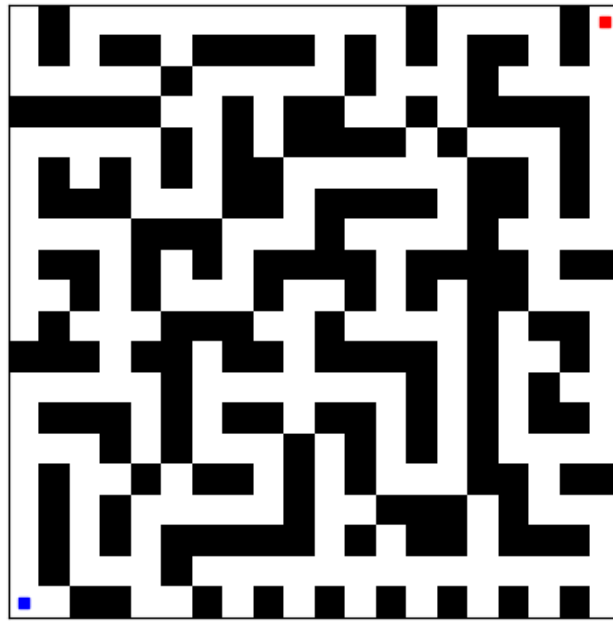


Figure 1:  $20 \times 20$  Generated Maze

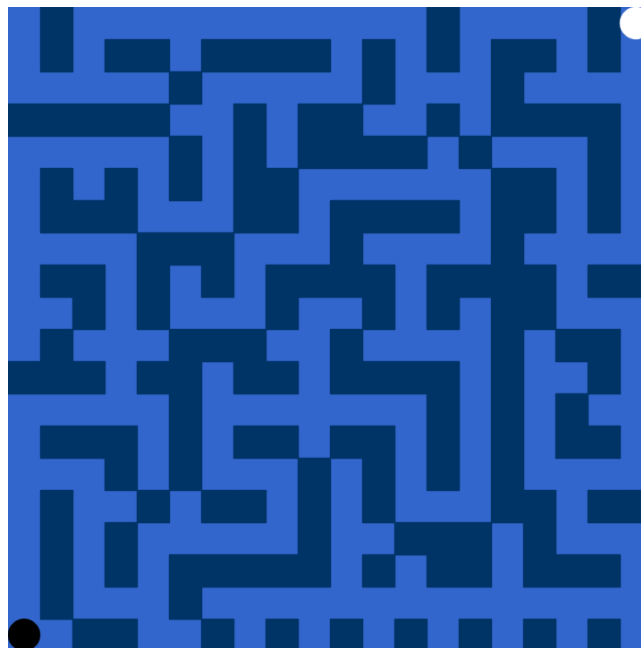


Figure 2: Other representation of  $20 \times 20$  Generated Maze

**Does the Maze\_Generator guarantee that there is one path between the up-right and bottom-left cells?**

The algorithm starts at a random cell and explores as far as possible along branches before backtracking. Each new cell is connected to the existing maze in a way that prevents loops and isolated sections. When adding a new cell, the code verifies that it has exactly one occupied neighbor (`ctr == 1`) and it prevents loops. We can also observe that the last 2 lines of the

code guarantees that the mentioned 2 cells are white and by our explanations, we can discover that there is one path between these 2 cells.

## Q-Learning

Q-learning is a popular model-free reinforcement learning algorithm used in machine learning and artificial intelligence applications. It falls under the category of temporal difference learning techniques, in which an agent picks up new information by observing results, interacting with the environment, and getting feedback in the form of rewards. Q-learning models engage in an iterative process where various components collaborate to train the model. This iterative procedure encompasses the agent exploring the environment and continuously updating the model based on this exploration. The key components of Q-learning include:

- Agents: Entities that operate within an environment, making decisions and taking actions.
- States: Variables that identify an agent's current position in the environment.
- Actions: Operations undertaken by the agent in specific states.
- Rewards: Positive or negative responses provided to the agent based on its actions.
- Episodes: Instances where an agent concludes its actions, marking the end of an episode.
- Q-values: Metrics used to evaluate actions at specific states.

In Q-learning I have implemented two classes: Environment and Agent.

### Environment

The Environment class represents a maze environment where an agent navigates from a starting position to a goal position while avoiding walls. `is_valid_move` Ensures that a move is within maze boundaries and does not land on walls. `get_next_state` Computes the next state by adding the movement vector to the current position.

```
1 class Environment:
2     def __init__(self, maze_file):
3         self.maze = np.load(maze_file)
4         self.actions = [(0, 1), (1, 0), (0, -1), (-1, 0)] # Right, Down, Left
5         , Up
6         self.start_state = (0, self.maze.shape[1] - 1)
7         self.goal_state = (self.maze.shape[0] - 1, 0)
8
9     def is_valid_move(self, x, y):
10         return 0 <= x < self.maze.shape[0] and 0 <= y < self.maze.shape[1] and
11             self.maze[x, y] == 0
12
13     def get_next_state(self, x, y, action):
14         dx, dy = self.actions[action]
15         next_x, next_y = x + dx, y + dy
16         if self.is_valid_move(next_x, next_y):
17             return next_x, next_y
18         return x, y
```

### Agent

The Agent class implements Q-learning. `alpha` is learning rate which shows how much the agent update q-values in each step. `gamma` is discount factor which shows how much future

rewards matter compared to immediate rewards., and **epsilon** is exploration rate which shows the probability of taking a random action instead of the best-known action. **Q-table** stores the estimated rewards for each (state, action) pair with same shape as maze. **train** function runs Q-learning to update the Q-table based on rewards. By reaching the goal, reward is 1, by hitting the wall reward is -0.2, and by valid move reward is -0.01. We update the q-values by using the Bellman equation. By using **generate\_solution\_gif**, we create a gif of how our model is working. Besides, we use **show\_policy** to observe the policy after running the algorithm. By using **\_create\_policy\_gif**, we create the gif for policy evolution after each 100 episodes.

```

1 class Agent:
2     def __init__(self, environment, alpha=0.1, gamma=0.9, epsilon=0.1,
3         num_episodes=5000):
4         self.env = environment
5         self.alpha = alpha
6         self.gamma = gamma
7         self.epsilon = epsilon
8         self.num_episodes = num_episodes
9         self.q_table = np.zeros((*self.env.maze.shape, len(self.env.actions)))
10
11     def train(self, frame_interval=100, gif_name="policy_evolution.gif"):
12         frame_filenames = []
13
14         for episode in range(self.num_episodes):
15             state = self.env.start_state
16             episode_reward = 0
17             while state != self.env.goal_state:
18                 x, y = state
19                 if random.uniform(0, 1) < self.epsilon:
20                     action = random.randint(0, len(self.env.actions) - 1) #
21                     Explore
22                 else:
23                     action = np.argmax(self.q_table[x, y]) # Exploit
24                     next_x, next_y = self.env.get_next_state(x, y, action)
25
26                     if (next_x, next_y) == self.env.goal_state:
27                         reward = 1
28                     elif (next_x, next_y) == (x, y):
29                         reward = -0.2 # High penalty for hitting a wall
30                     else:
31                         reward = -0.01 # Small penalty for valid move
32
33                     max_next_q = np.max(self.q_table[next_x, next_y])
34                     self.q_table[x, y, action] += self.alpha * (
35                         reward + self.gamma * max_next_q - self.q_table[x, y,
36                             action]
37                     )
38
39                     state = (next_x, next_y)
40                     episode_reward += reward
41
42             if episode % frame_interval == 0:
43                 filename = f"temp_policy_frame_{episode}.png"
44                 self._save_policy_frame(episode, filename)
45                 frame_filenames.append(filename)
46
47             if episode % 1000 == 0:

```

```
45         print(f"Episode {episode}, Epsilon: {self.epsilon:.3f}, Reward  
: {episode_reward}")  
46  
47         self._create_policy_gif(frame_filenames, gif_name)  
48  
49         # Cleanup temporary files  
50         for filename in frame_filenames:  
51             os.remove(filename)  
52  
53     def generate_solution_gif(self, output_filename="solution_path.gif"):  
54         frames = []  
55         state = self.env.start_state  
56         solution_path = [state]  
57         while state != self.env.goal_state:  
58             x, y = state  
59             action = np.argmax(self.q_table[x, y])  
60             next_x, next_y = self.env.get_next_state(x, y, action)  
61             solution_path.append((next_x, next_y))  
62             state = (next_x, next_y)  
63             frame = np.copy(self.env.maze)  
64             for sx, sy in solution_path:  
65                 frame[sx, sy] = 0.5 # Path is shown in gray  
66             frame[x, y] = 0.8 # Agent is shown in red  
67             fig, ax = plt.subplots(figsize=(5, 5))  
68             ax.imshow(frame, cmap="gray")  
69             ax.scatter(y, x, c="red", s=100) # Highlight agent with a red dot  
70             ax.axis("off")  
71             plt.tight_layout()  
72  
73             # Save frame to buffer  
74             buf = f"frame_{len(frames)}.png"  
75             plt.savefig(buf, dpi=100, bbox_inches='tight')  
76             frames.append(imageio.imread(buf))  
77             plt.close()  
78             os.remove(buf) # Remove the intermediate PNG file  
79  
80             # Ensure the last frame is included  
81             x, y = solution_path[-1]  
82             fig, ax = plt.subplots(figsize=(5, 5))  
83             ax.imshow(self.env.maze, cmap="gray")  
84             for sx, sy in solution_path:  
85                 ax.scatter(sy, sx, c="gray", s=50)  
86             ax.scatter(y, x, c="red", s=100) # Final position in red  
87             ax.axis("off")  
88             plt.tight_layout()  
89             buf = "final_frame.png"  
90             plt.savefig(buf, dpi=100, bbox_inches='tight')  
91             frames.append(imageio.imread(buf))  
92             plt.close()  
93             os.remove(buf) # Remove the final PNG file  
94  
95             # Save the GIF  
96             imageio.mimsave(output_filename, frames, fps=5)  
97             print(f"Q-learning completed. GIF saved as '{output_filename}'.")  
98  
99     def show_policy(self):  
100         policy_arrows = {  
101             0: "\u2192", # Right arrow
```

```

102         1: "\u2193", # Down arrow
103         2: "\u2190", # Left arrow
104         3: "\u2191"  # Up arrow
105     }
106
107     maze_with_policy = np.copy(self.env.maze)
108     fig, ax = plt.subplots(figsize=(5, 5))
109
110     for i in range(maze_with_policy.shape[0]):
111         for j in range(maze_with_policy.shape[1]):
112             if maze_with_policy[i, j] == 0: # Open cell
113                 if (i, j) == (maze_with_policy.shape[0] - 1, 0): # Bottom
-left corner
114                     ax.plot(j, i, 'bs', markersize = 4)
115                 else:
116                     best_action = np.argmax(self.q_table[i, j])
117                     arrow = policy_arrows[best_action]
118                     ax.text(j, i, arrow, ha='center', va='center', color='
black', fontsize=12)
119     ax = plt.gca()
120     ax.imshow(maze_with_policy, "Greys")
121     plt.xticks([], [])
122     plt.yticks([], [])
123
124     plt.show()
125
126     def _save_policy_frame(self, episode, filename):
127         # Reuse the visualization logic from show_policy() but save to file
128         policy_arrows = {
129             0: "\u2192", 1: "\u2193", 2: "\u2190", 3: "\u2191"
130         }
131         maze = np.copy(self.env.maze)
132         fig, ax = plt.subplots(figsize=(5, 5))
133
134         for i in range(maze.shape[0]):
135             for j in range(maze.shape[1]):
136                 if maze[i, j] == 0: # Open cell
137                     best_action = np.argmax(self.q_table[i, j])
138                     arrow = policy_arrows[best_action]
139                     ax.text(j, i, arrow, ha='center', va='center', color='
black', fontsize=12)
140
141         ax.imshow(maze, "Greys")
142         plt.xticks([], [])
143         plt.yticks([], [])
144         plt.title(f"Episode: {episode}")
145         plt.savefig(filename, bbox_inches='tight', dpi=100)
146         plt.close()
147
148     def _create_policy_gif(self, frame_files, output_filename, fps=2):
149         # Compile frames into a GIF
150         frames = [imageio.imread(filename) for filename in frame_files]
151         imageio.mimsave(output_filename, frames, fps=fps)
152         print(f"Policy evolution GIF saved as '{output_filename}'.")

```



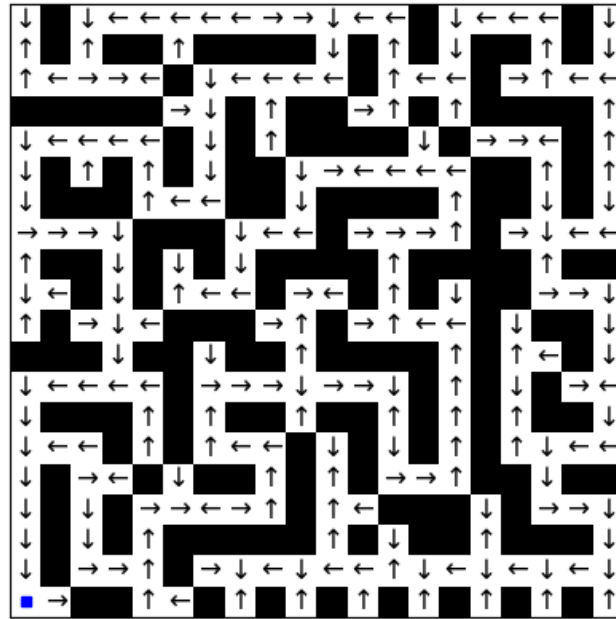


Figure 3: Policy Map Plot After Q-Learning Algorithm

## Deep Q-Learning

Source of this explanations is this [website](#). The Deep Q-Network (DQN) algorithm enhances the traditional Q-learning method by incorporating neural networks to approximate Q-values.

1. Initialize your Main and Target neural networks.
2. Choose an action using the Epsilon-Greedy Exploration Strategy.
3. Update your network weights using the Bellman Equation.

### Initialize Your Target and Main Neural Networks

A core difference between Deep Q-Learning and Vanilla Q-Learning is the implementation of the Q-table. Instead of a Q-table, Deep Q-Learning uses a neural network to map input states to  $(action, Q-value)$  pairs. One important feature of Deep Q-Learning is the use of two neural networks: the **Main** and **Target** networks. These networks share the same architecture but have different weights. Every  $N$  steps, the weights of the Main network are copied to the Target network, stabilizing the learning process. In our implementation, we update the Target network every 100 steps.

### Mapping States to (Action, Q-value) Pairs

The Main and Target neural networks map input states to  $(action, Q-value)$  pairs. Each output node represents an action and holds its respective Q-value as a floating-point number. Unlike probabilities, these values do not sum to 1.

### Choosing an Action Using Epsilon-Greedy Strategy

In the Epsilon-Greedy Exploration strategy, the agent chooses a random action with probability  $\epsilon$  and exploits the best-known action with probability  $1 - \epsilon$ . The best-known action corresponds to the action with the highest predicted Q-value from the Main network.

## Updating Network Weights Using the Bellman Equation

Once an action is selected, the agent performs it and updates the Main and Target networks according to the Bellman Equation. Deep Q-Learning employs Experience Replay to improve the learning process.

### Experience Replay

Experience Replay is a mechanism that stores past experiences as tuples (*state, action, reward, next\_state*). This technique allows the agent to learn from past experiences in an offline manner, reducing correlations between consecutive experiences and stabilizing training.

Instead of updating the network after each step, Experience Replay trains the network using small batches every 4 steps. This method improves efficiency and accelerates Deep Q-Learning.

### The Bellman Equation

Similar to Vanilla Q-Learning, Deep Q-Learning updates model weights using the Bellman Equation:

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

The Temporal Difference (TD) target is computed using the Target network rather than the Main network. The Main network is then updated to fit these target Q-values, ensuring stable learning.

Here, we have `environment` and `agent` classes too.

### Environment

```

1 import numpy as np
2 import scipy.special as sp
3 import matplotlib.pyplot as plt
4 import copy
5
6 class MazeEnvironment:
7     def __init__(self, maze, init_position, goal):
8         x = len(maze)
9         y = len(maze)
10
11         self.boundary = np.asarray([x, y])
12         self.init_position = init_position
13         self.current_position = np.asarray(init_position)
14         self.goal = goal
15         self.maze = maze
16
17         self.visited = set()
18         self.visited.add(tuple(self.current_position))
19
20         # initialize the empty cells and the euclidean distance from
21         # the goal (removing the goal cell itself)
22         self.allowed_states = np.asarray(np.where(self.maze == 0)).T.tolist()
23         self.distances = np.sqrt(np.sum((np.array(self.allowed_states) -
24                                         np.asarray(self.goal))**2,
25                                         axis = 1))
26
27         del(self.allowed_states[np.where(self.distances == 0)[0][0]])
28         self.distances = np.delete(self.distances, np.where(self.distances ==
29         0)[0][0])

```

```
30     self.action_map = {0: [0, 1],
31                        1: [0, -1],
32                        2: [1, 0],
33                        3: [-1, 0]}
34
35     self.directions = {0: '→',
36                       1: '←',
37                       2: '↓',
38                       3: '↑'}
39
40     # the agent makes an action from the following:
41     # 1 -> right, 2 -> left
42     # 3 -> down, 4 -> up
43
44     # introduce a reset policy, so that for high epsilon the initial
45     # position is nearer to the goal (useful for large mazes)
46     def reset_policy(self, eps, reg = 7):
47         return sp.softmax(-self.distances/(reg*(1-eps**(2/reg))))**(reg/2)).
squeeze()
48
49     # reset the environment when the game is completed
50     # with probability prand the reset is random, otherwise
51     # the reset policy at the given epsilon is used
52     def reset(self, epsilon, prand = 0):
53         if np.random.rand() < prand:
54             idx = np.random.choice(len(self.allowed_states))
55         else:
56             p = self.reset_policy(epsilon)
57             idx = np.random.choice(len(self.allowed_states), p = p)
58
59         self.current_position = np.asarray(self.allowed_states[idx])
60
61         self.visited = set()
62         self.visited.add(tuple(self.current_position))
63
64         return self.state()
65
66
67     def state_update(self, action):
68         isgameon = True
69
70         # each move costs -0.05
71         reward = -0.05
72
73         move = self.action_map[action]
74         next_position = self.current_position + np.asarray(move)
75
76         # if the goals has been reached, the reward is 1
77         if (self.current_position == self.goal).all():
78             reward = 1
79             isgameon = False
80             return [self.state(), reward, isgameon]
81
82         # if the cell has been visited before, the reward is -0.2
83         else:
84             if tuple(self.current_position) in self.visited:
85                 reward = -0.2
86
```

```

87     # if the moves goes out of the maze or to a wall, the
88     # reward is -1
89     if self.is_state_valid(next_position):
90         self.current_position = next_position
91     else:
92         reward = -1
93
94     self.visited.add(tuple(self.current_position))
95     return [self.state(), reward, isgameon]
96
97     # return the state to be feeded to the network
98     def state(self):
99         state = copy.deepcopy(self.maze)
100         state[tuple(self.current_position)] = 2
101         return state
102
103
104     def check_boundaries(self, position):
105         out = len([num for num in position if num < 0])
106         out += len([num for num in (self.boundary - np.asarray(position)) if
num <= 0])
107         return out > 0
108
109
110     def check_walls(self, position):
111         return self.maze[tuple(position)] == 1
112
113
114     def is_state_valid(self, next_position):
115         if self.check_boundaries(next_position):
116             return False
117         elif self.check_walls(next_position):
118             return False
119         return True
120
121
122     def draw(self, filename):
123         plt.figure()
124         im = plt.imshow(self.maze, interpolation='none', aspect='equal', cmap=
'Greys');
125         ax = plt.gca();
126
127         plt.xticks([], [])
128         plt.yticks([], [])
129
130         ax.plot(self.goal[1], self.goal[0],
131                 'bs', markersize = 4)
132         ax.plot(self.current_position[1], self.current_position[0],
133                 'rs', markersize = 4)
134         plt.savefig(filename, dpi = 300, bbox_inches = 'tight')
135         plt.show()

```

## Agent

```

1 import numpy as np
2 import scipy.special as sp
3 import matplotlib.pyplot as plt
4 import copy

```

```
5 import torch
6 import torch.nn as nn
7 import torch.optim as optim
8 import collections
9
10 Transition = collections.namedtuple('Experience',
11                                     field_names=['state', 'action',
12                                                  'next_state', 'reward',
13                                                  'is_game_on'])
14
15
16 class Agent:
17     def __init__(self, maze, memory_buffer, use_softmax = True):
18         self.env = maze
19         self.buffer = memory_buffer # this is actually a reference
20         self.num_act = 4
21         self.use_softmax = use_softmax
22         self.total_reward = 0
23         self.min_reward = -self.env.maze.size
24         self.isgameon = True
25
26
27     def make_a_move(self, net, epsilon, device = 'cuda'):
28         action = self.select_action(net, epsilon, device)
29         current_state = self.env.state()
30         next_state, reward, self.isgameon = self.env.state_update(action)
31         self.total_reward += reward
32
33         if self.total_reward < self.min_reward:
34             self.isgameon = False
35         if not self.isgameon:
36             self.total_reward = 0
37
38         transition = Transition(current_state, action,
39                                next_state, reward,
40                                self.isgameon)
41
42         self.buffer.push(transition)
43
44
45     def select_action(self, net, epsilon, device = 'cuda'):
46         state = torch.Tensor(self.env.state()).to(device).view(1,-1)
47         qvalues = net(state).cpu().detach().numpy().squeeze()
48
49         # softmax sampling of the qvalues
50         if self.use_softmax:
51             p = sp.softmax(qvalues/epsilon).squeeze()
52             p /= np.sum(p)
53             action = np.random.choice(self.num_act, p = p)
54
55         # else choose the best action with probability 1-epsilon
56         # and with probability epsilon choose at random
57         else:
58             if np.random.random() < epsilon:
59                 action = np.random.randint(self.num_act, size=1)[0]
60             else:
61                 action = np.argmax(qvalues, axis=0)
62                 action = int(action)
```

```

63
64         return action
65
66
67     def plot_policy_map(self, net, filename, offset):
68         net.eval()
69         with torch.no_grad():
70             fig, ax = plt.subplots()
71             ax.imshow(self.env.maze, 'Greys')
72
73             for free_cell in self.env.allowed_states:
74                 self.env.current_position = np.asarray(free_cell)
75                 qvalues = net(torch.Tensor(self.env.state()).view(1,-1).to('
cuda'))
76                 action = int(torch.argmax(qvalues).detach().cpu().numpy())
77                 policy = self.env.directions[action]
78
79                 ax.text(free_cell[1]-offset[0], free_cell[0]-offset[1], policy
)
80             ax = plt.gca();
81
82             plt.xticks([], [])
83             plt.yticks([], [])
84
85             ax.plot(self.env.goal[1], self.env.goal[0],
86                     'bs', markersize = 4)
87             plt.savefig(filename, dpi = 300, bbox_inches = 'tight')
88             plt.show()

```

## Implementation

```

1 class ExperienceReplay:
2     def __init__(self, capacity):
3         self.capacity = capacity
4         self.memory = collections.deque(maxlen=capacity)
5
6     def __len__(self):
7         return len(self.memory)
8
9     def push(self, transition):
10        self.memory.append(transition)
11
12    def sample(self, batch_size, device = 'cuda'):
13        indices = np.random.choice(len(self.memory), batch_size, replace = False)
14
15        states, actions, next_states, rewards, isgameon = zip(*[self.memory[idx]
16                                                                for idx in indices])
17
18        return torch.Tensor(states).type(torch.float).to(device), \
19               torch.Tensor(actions).type(torch.long).to(device), \
20               torch.Tensor(next_states).to(device), \
21               torch.Tensor(rewards).to(device), torch.tensor(isgameon).to(device)
22
23 class fc_nn(nn.Module):
24     def __init__(self, Ni, Nh1, Nh2, No = 4):
25         super().__init__()
26
27         self.fc1 = nn.Linear(Ni, Nh1)

```

```

28     self.fc2 = nn.Linear(Nh1, Nh2)
29     self.fc3 = nn.Linear(Nh2, No)
30
31     self.act = nn.ReLU()
32
33     def forward(self, x, classification = False, additional_out=False):
34         x = self.act(self.fc1(x))
35         x = self.act(self.fc2(x))
36         out = self.fc3(x)
37
38         return out
39
40 class conv_nn(nn.Module):
41     channels = [16, 32, 64]
42     kernels = [3, 3, 3]
43     strides = [1, 1, 1]
44     in_channels = 1
45
46     def __init__(self, rows, cols, n_act):
47         super().__init__()
48         self.rows = rows
49         self.cols = cols
50
51         self.conv = nn.Sequential(nn.Conv2d(in_channels = self.in_channels,
52                                             out_channels = self.channels[0],
53                                             kernel_size = self.kernels[0],
54                                             stride = self.strides[0]),
55                                   nn.ReLU(),
56                                   nn.Conv2d(in_channels = self.channels[0],
57                                             out_channels = self.channels[1],
58                                             kernel_size = self.kernels[1],
59                                             stride = self.strides[1]),
60                                   nn.ReLU()
61                                   )
62
63         size_out_conv = self.get_conv_size(rows, cols)
64
65         self.linear = nn.Sequential(nn.Linear(size_out_conv, rows*cols*2),
66                                   nn.ReLU(),
67                                   nn.Linear(rows*cols*2, int(rows*cols/2)),
68                                   nn.ReLU(),
69                                   nn.Linear(int(rows*cols/2), n_act),
70                                   )
71
72     def forward(self, x):
73         x = x.view(len(x), self.in_channels, self.rows, self.cols)
74         out_conv = self.conv(x).view(len(x),-1)
75         out_lin = self.linear(out_conv)
76         return out_lin
77
78     def get_conv_size(self, x, y):
79         out_conv = self.conv(torch.zeros(1,self.in_channels, x, y))
80         return int(np.prod(out_conv.size()))
81
82 def Qloss(batch, net, gamma=0.99, device="cuda"):
83     states, actions, next_states, rewards, _ = batch
84     lbatch = len(states)
85     state_action_values = net(states.view(lbatch,-1))

```

```
86 state_action_values = state_action_values.gather(1, actions.unsqueeze(-1))
87 state_action_values = state_action_values.squeeze(-1)
88
89 next_state_values = net(next_states.view(1batch, -1))
90 next_state_values = next_state_values.max(1)[0]
91
92 next_state_values = next_state_values.detach()
93 expected_state_action_values = next_state_values * gamma + rewards
94
95 return nn.MSELoss()(state_action_values, expected_state_action_values)
```

```
1 from environment import MazeEnvironment
2 from agent import Agent
3
4 maze = np.load('maze.npy')
5 initial_position = [0, len(maze)-1]
6 goal = [len(maze)-1, 0]
7 maze_env = MazeEnvironment(maze, initial_position, goal)
8
9 buffer_capacity = 10000
10 buffer_start_size = 1000
11 memory_buffer = ExperienceReplay(buffer_capacity)
12
13 agent = Agent(maze = maze_env,
14               memory_buffer = memory_buffer,
15               use_softmax = True
16               )
17
18 net = fc_nn(maze.size, maze.size, maze.size, 4)
19 optimizer = optim.Adam(net.parameters(), lr=1e-4)
20 device = 'cuda'
21 net.to(device)
22 batch_size = 32
23 gamma = 0.9
24
25 num_epochs = 4000
26 cutoff = 3000
27 epsilon = np.exp(-np.arange(num_epochs)/(cutoff))
28 epsilon[epsilon > epsilon[100*int(num_epochs/cutoff)]] = epsilon[100*int(
num_epochs/cutoff)]
29
30 loss_log = []
31 best_loss = 1e5
32
33 running_loss = 0
34
35 for epoch in range(num_epochs):
36     loss = 0
37     counter = 0
38     eps = epsilon[epoch]
39
40     agent.isgameon = True
41     _ = agent.env.reset(eps)
42
43     while agent.isgameon:
44         agent.make_a_move(net, eps)
45         counter += 1
46
```



```
47         if len(agent.buffer) < buffer_start_size:
48             continue
49
50         optimizer.zero_grad()
51         batch = agent.buffer.sample(batch_size, device = device)
52         loss_t = Qloss(batch, net, gamma = gamma, device = device)
53         loss_t.backward()
54         optimizer.step()
55
56         loss += loss_t.item()
57
58         if (agent.env.current_position == agent.env.goal).all():
59             result = 'won'
60         else:
61             result = 'lost'
62
63         if epoch%1000 == 0:
64             agent.plot_policy_map(net, 'sol_epoch_'+str(epoch)+'.pdf',
[0.35,-0.3])
65
66         loss_log.append(loss)
67
68         if (epoch > 2000):
69             running_loss = np.mean(loss_log[-50:])
70             if running_loss < best_loss:
71                 best_loss = running_loss
72                 torch.save(net.state_dict(), "best.torch")
73                 estop = epoch
74
75         print('Epoch', epoch, '(number of moves ' + str(counter) + ')')
76         print('Game', result)
77         print([' + '#*(100-int(100*(1 - epoch/num_epochs))) +
78             ' '*int(100*(1 - epoch/num_epochs)) + ']')
79         print('\t Average loss: ' + f'{loss:.5f}')
80         if (epoch > 2000):
81             print('\t Best average loss of the last 50 epochs: ' + f'{
best_loss:.5f}' + ', achieved at epoch', estop)
82         clear_output(wait = True)
83
84         torch.save(net.state_dict(), "net.torch")
85         agent.plot_policy_map(net, 'solution.pdf', [0.35,-0.3])
```

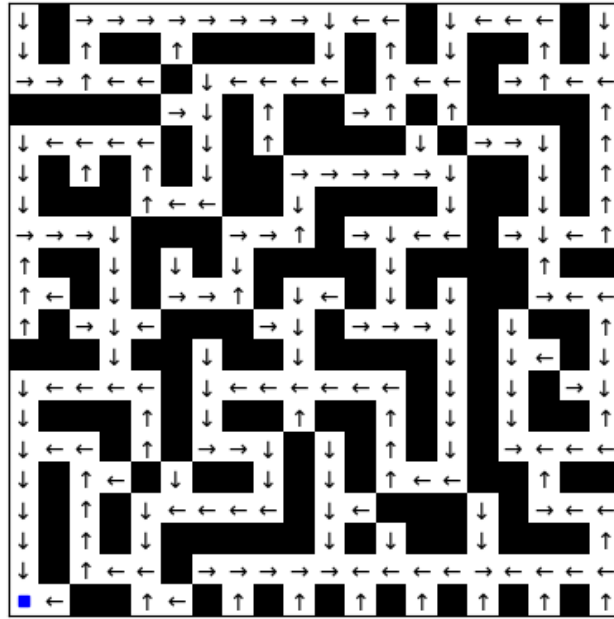


Figure 4: Policy Map Plot After 4000 Epochs of Training with Deep Q-Learning Algorithm

By comparing q-learning and deep q-learning, we understand that our deep method works better. But the training of the deep method takes much time rather than the normal method. I have used 5000 episodes for q-learning and 4000 epochs for deep q-learning. Running of q-learning method takes 20 seconds and running of deep method takes 4 hours.

## 100 × 100 Maze

In this part, first, we create the maze.

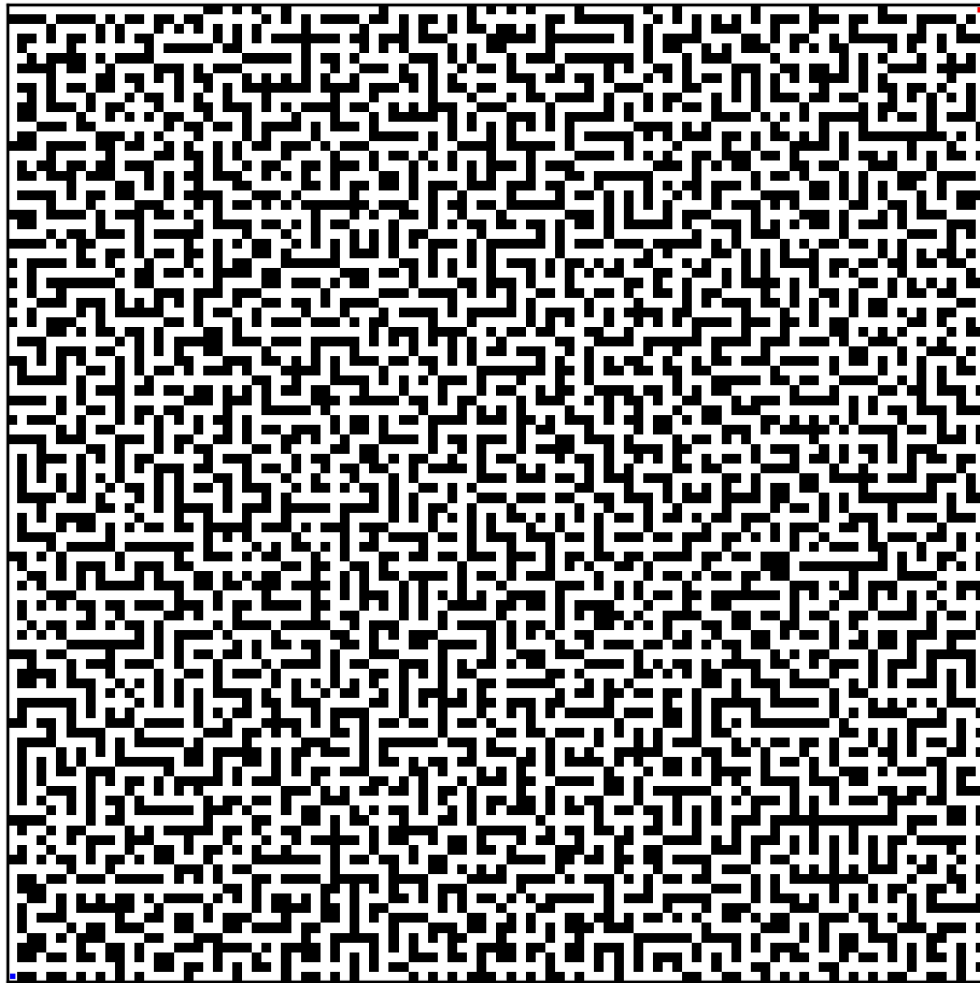


Figure 5: 100×100 Generated Maze

We use all scripts that we've used in previous parts. Now, We compare the results. It is to say that the running for a  $100 \times 100$  maze takes more times in comparison with a  $20 \times 20$  maze. for my computational limitations, the running of the q-learning method for  $100 \times 100$  maze was for 500 episodes that could not learn well. the time for 1000 episodes was about 1 hours that it was not good to and because of that, I prefer to run the script for 500 episodes. The comparison of q-learning and deep q-learning for  $100 \times 100$  maze was not required. I know the time for running its could based on my gpu and colab's limitations was not achievable. but the time and the computational cost for the larger maze would be more.

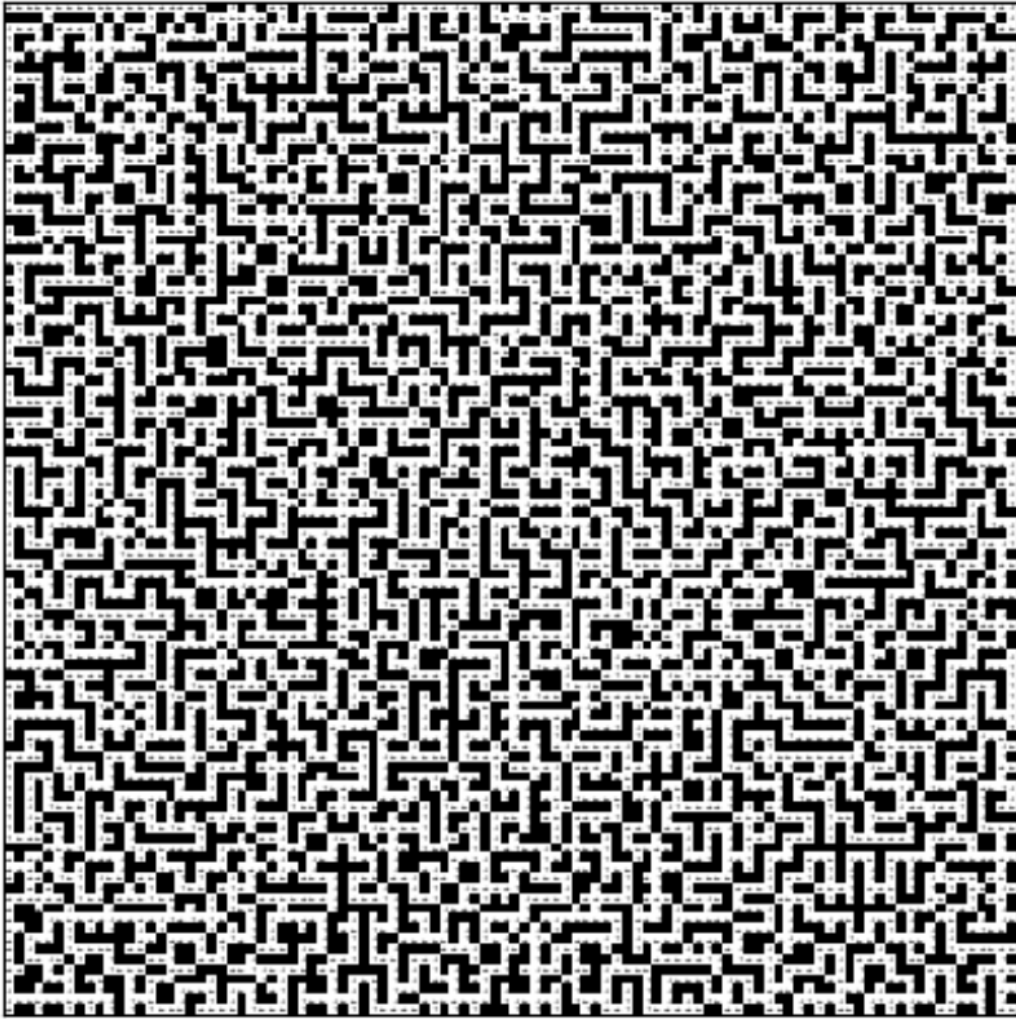


Figure 6: Policy Map Plot After q-learning algorithm for  $100 \times 100$  Generated Maze