

EE449 HW2*

*Note: This homework exaggerated and used LaTeX, I don't know why am I doing this.

Ömer Takkin
No.: 2516987
Electric and Electronic Engineer (EEE)
Middle East Technical University (METU)
Ankara, Turkey
omer.takkin@gmail.com

Abstract—This project explores fundamental reinforcement learning techniques by applying Temporal Difference (TD(0)) learning and Deep Q-Networks (DQN). For this, we implement two tasks: a maze navigation problem and the LunarLander-v3 environment. In the first part, we model a stochastic maze and implement TD(0) algorithm to derive optimal navigation policies. We will analyze performance through utility heatmaps and convergence plots. In the second part, we will build a DQN agent in PyTorch to learn a landing policy for a simulated lunar lander. We will experiment with various hyperparameters and network architectures to evaluate their impact on learning stability and efficiency.

Index Terms—METU, EEE, EE449, HW1, Reinforcement Learning, Help

I. TEMPORAL DIFFERENCE (TD) LEARNING

A. Preliminary

In this part of the experiment, a maze is structured with a probabilistic move agent. The agent tries to reach the goal while avoiding traps.

B. Experiments

1) *Value Iteration For Different Alpha Values*: In this section, we will observe value iteration for different alpha values

Value Function Evolution for $\alpha=0.001$

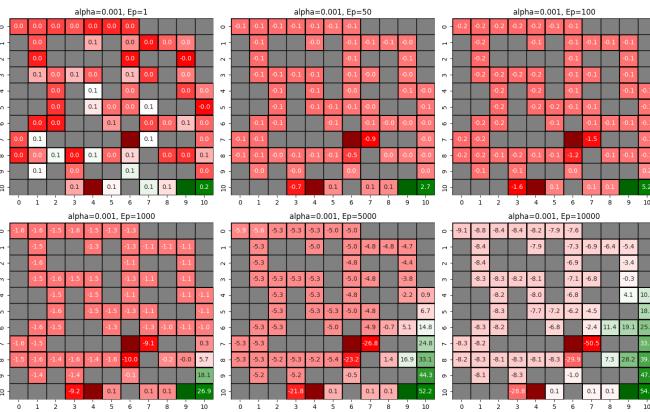


Fig. 1. Value Evaluation for $\alpha = 0.001$

Value Function Evolution for $\alpha=0.01$

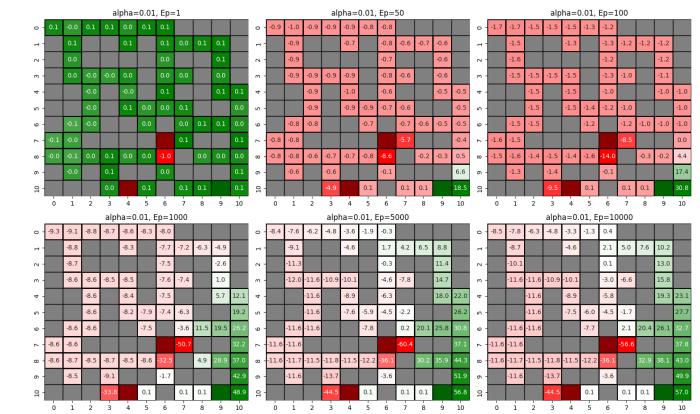


Fig. 2. Value Evaluation for $\alpha = 0.01$

Value Function Evolution for $\alpha=0.1$

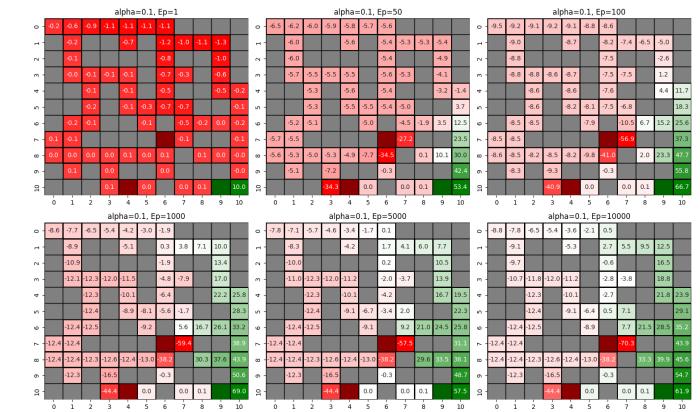


Fig. 3. Value Evaluation for $\alpha = 0.1$

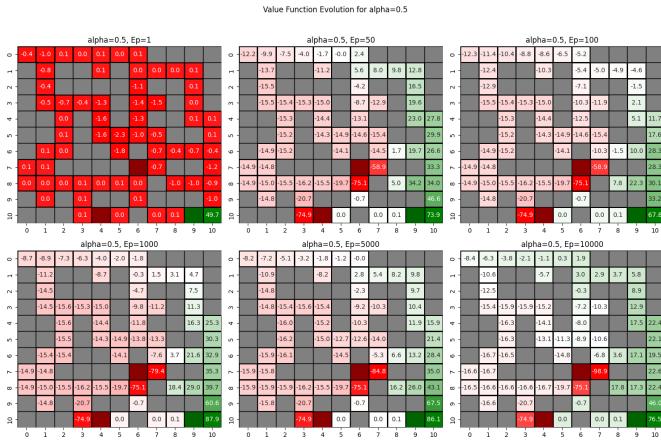


Fig. 4. Value Evaluation for $\alpha = 0.5$

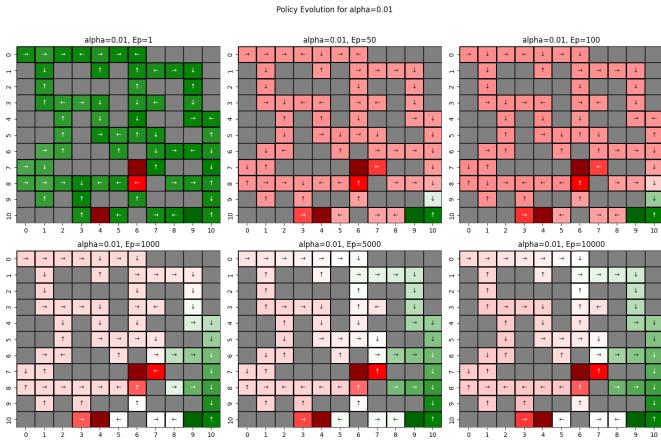


Fig. 7. Policy Evaluation for $\alpha = 0.01$

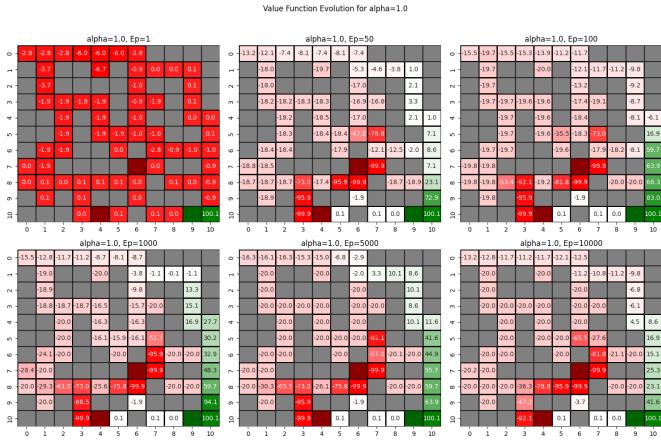


Fig. 5. Value Evaluation for $\alpha = 1.0$

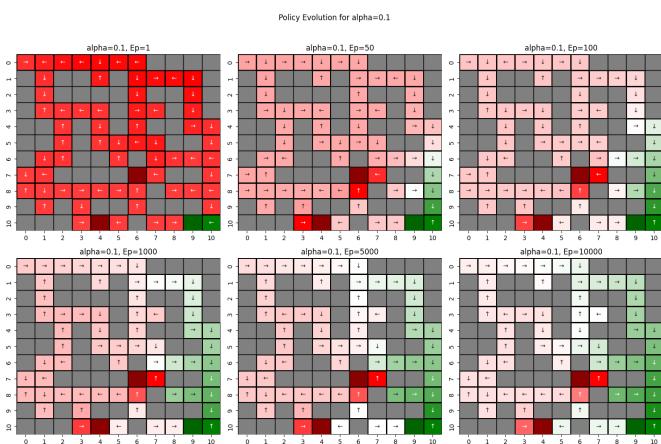


Fig. 8. Policy Evaluation for $\alpha = 0.1$

2) *Policy Iteration For Different Alpha Values:* In this section, we will observe policy iteration for different alpha values

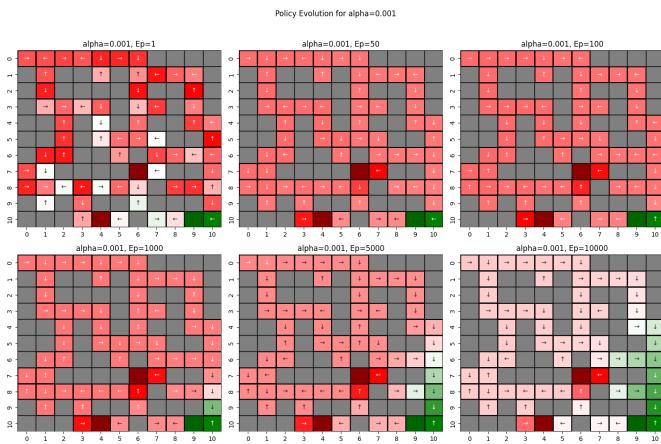


Fig. 6. Policy Evaluation for $\alpha = 0.001$

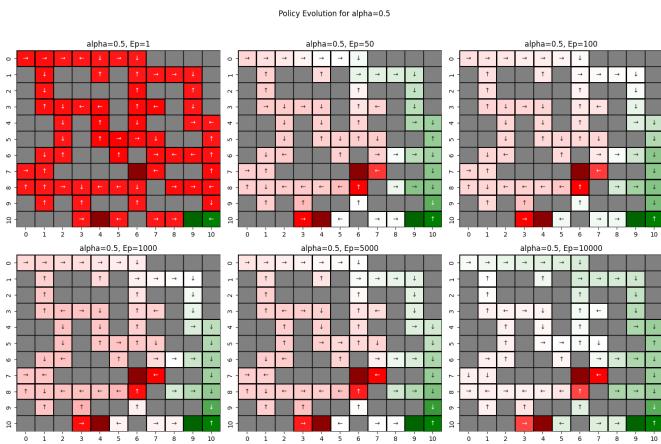


Fig. 9. Policy Evaluation for $\alpha = 0.5$

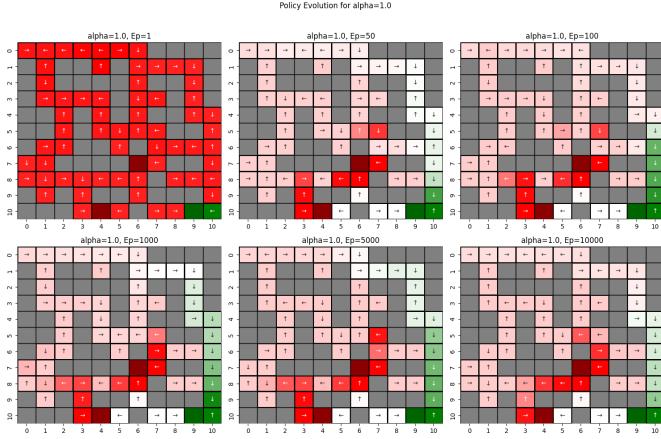


Fig. 10. Policy Evaluation for $\alpha = 1.0$

3) Convergence For Different Alpha Values: In this section, we will observe convergence for different alpha values

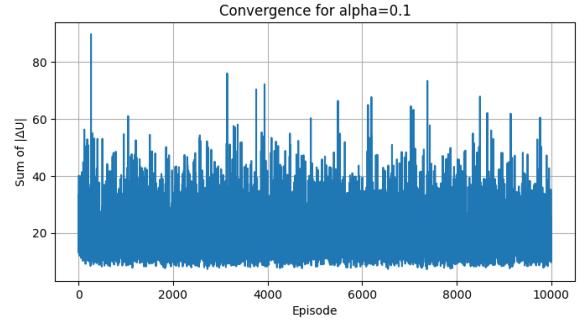


Fig. 13. Convergence for $\alpha = 0.1$

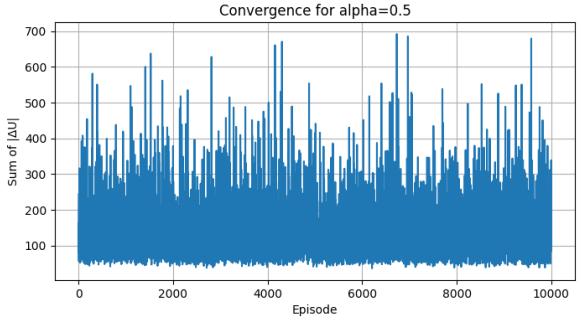


Fig. 14. Convergence for $\alpha = 0.5$

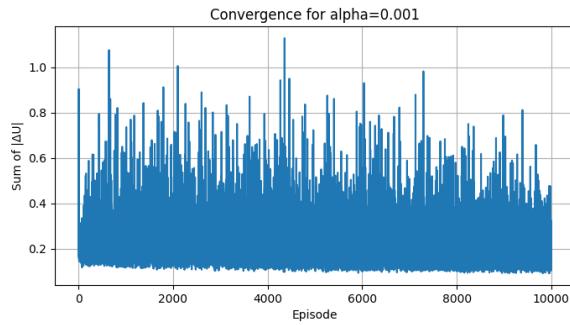


Fig. 11. Convergence for $\alpha = 0.001$

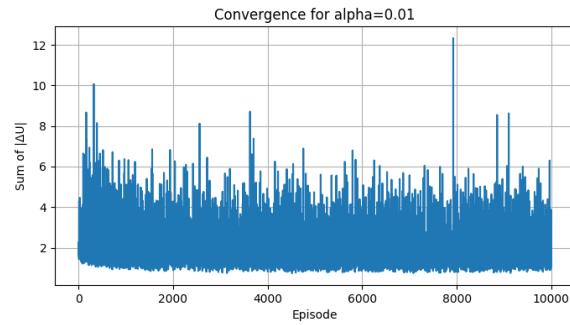


Fig. 12. Convergence for $\alpha = 0.01$

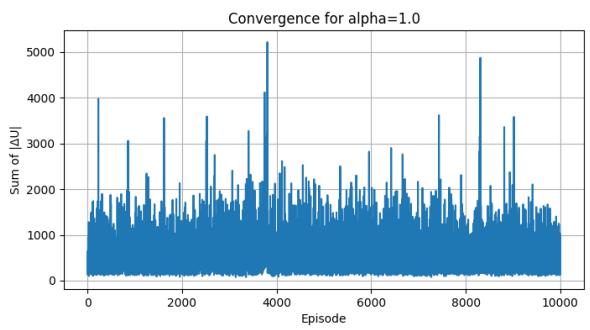


Fig. 15. Convergence for $\alpha = 1.0$

4) Value Iteration For Different Gamma Values: In this section, we will observe value iteration for different gamma values

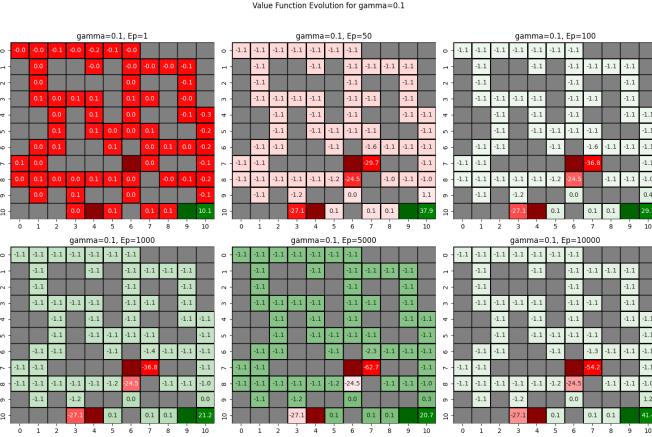


Fig. 16. Value Evaluation for $\gamma = 0.1$

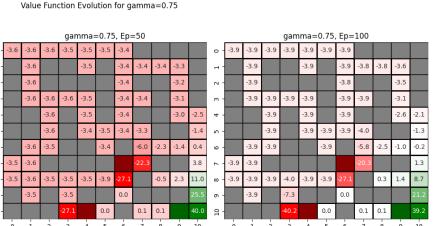


Fig. 19. Value Evaluation for $\gamma = 0.75$

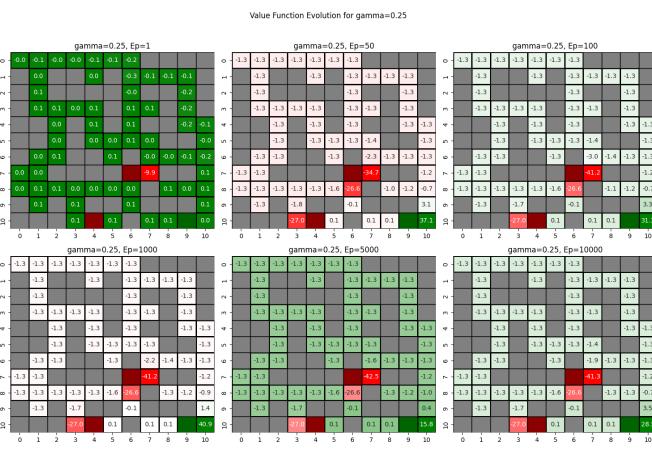


Fig. 17. Value Evaluation for $\gamma = 0.25$

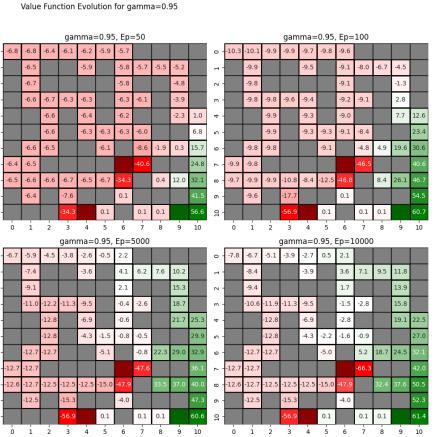


Fig. 20. Value Evaluation for $\gamma = 0.95$

5) Policy Iteration For Different Gamma Values: In this section, we will observe policy iteration for different gamma values

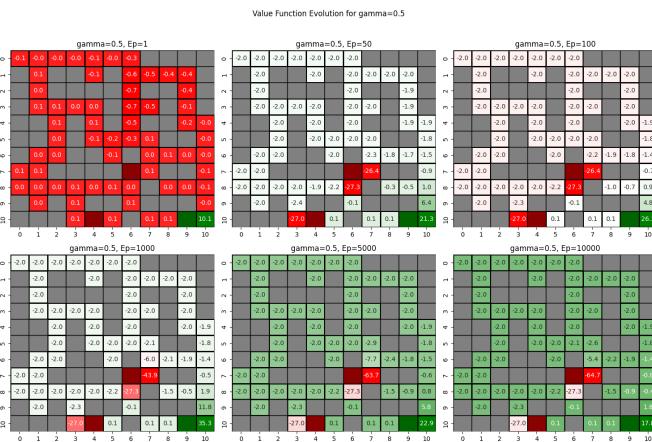


Fig. 18. Value Evaluation for $\gamma = 0.50$

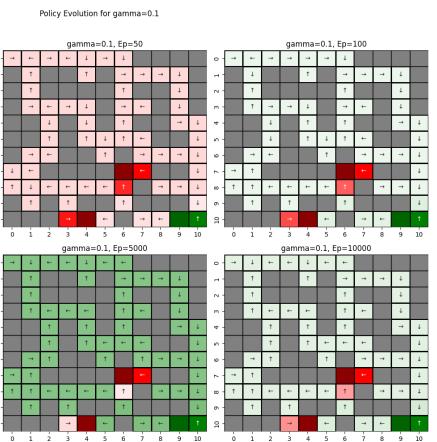


Fig. 21. Policy Evaluation for $\gamma = 0.1$

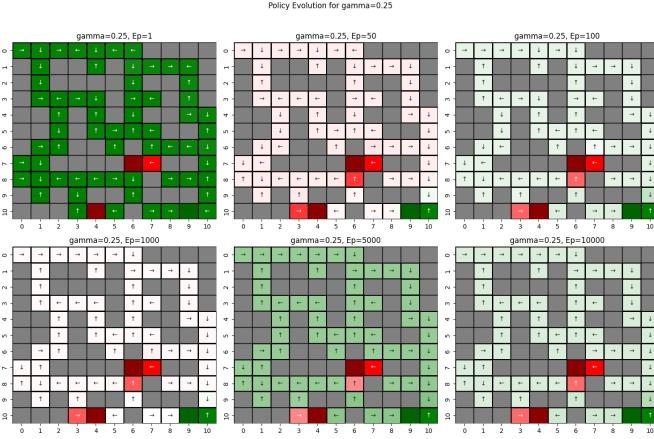


Fig. 22. Policy Evaluation for $\gamma = 0.25$

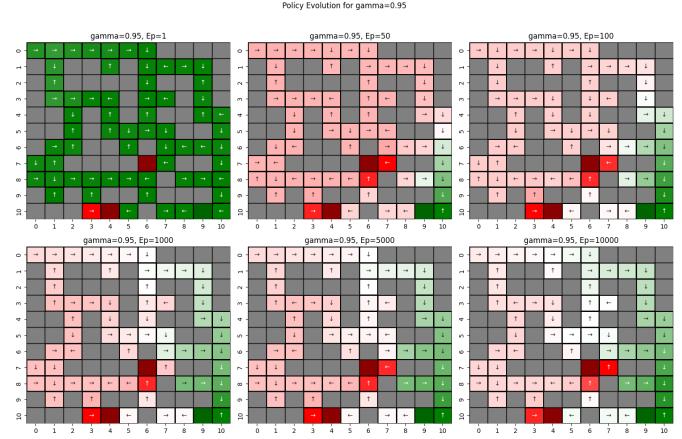


Fig. 25. Policy Evaluation for $\gamma = 0.95$

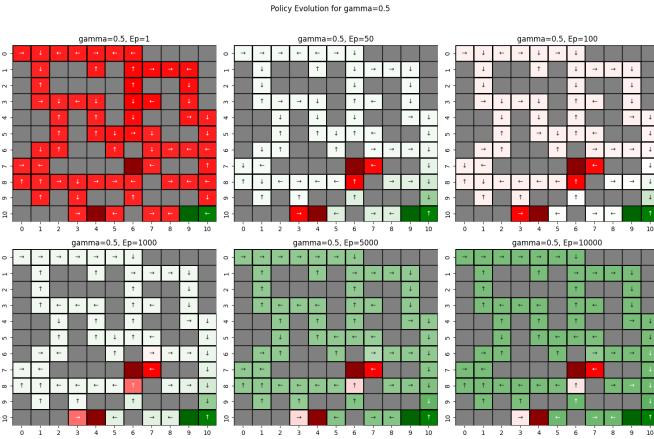


Fig. 23. Policy Evaluation for $\gamma = 0.5$

6) Convergence For Different Gamma Values: In this section, we will observe convergence for different gamma values

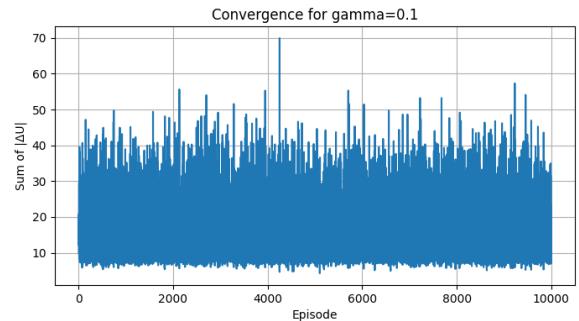


Fig. 26. Convergence for $\gamma = 0.1$

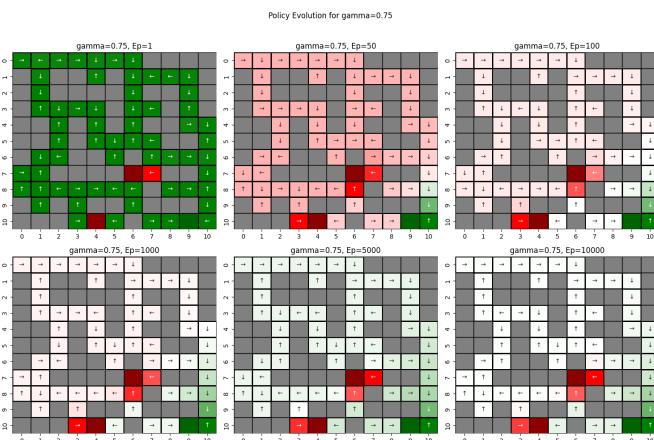


Fig. 24. Policy Evaluation for $\gamma = 0.75$

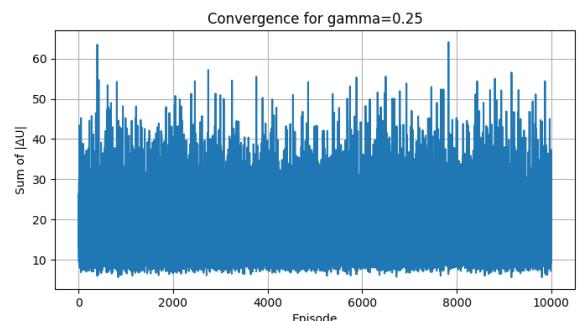


Fig. 27. Convergence for $\gamma = 0.25$

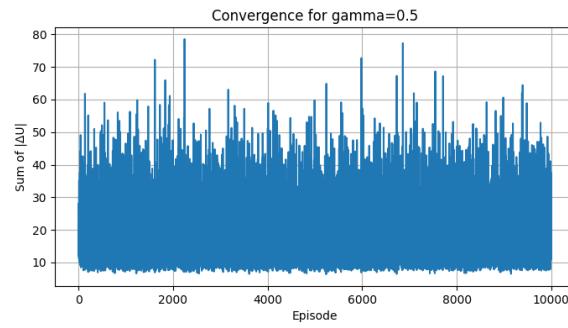


Fig. 28. Convergence for $\gamma = 0.50$

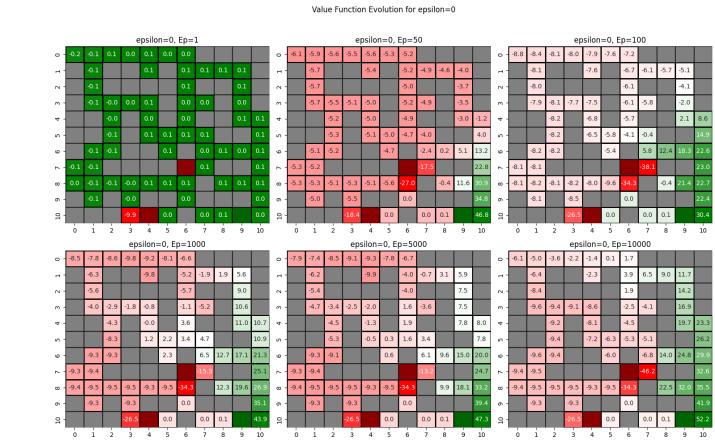


Fig. 31. Value Evaluation for $\epsilon = 0.0$

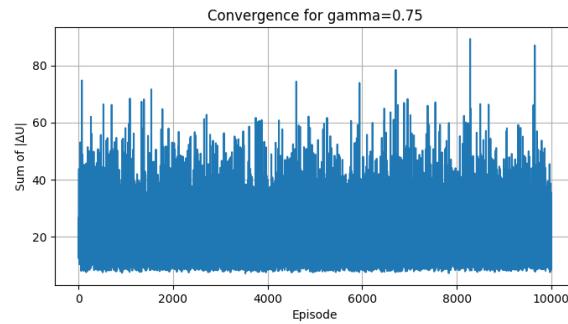


Fig. 29. Convergence for $\gamma = 0.75$

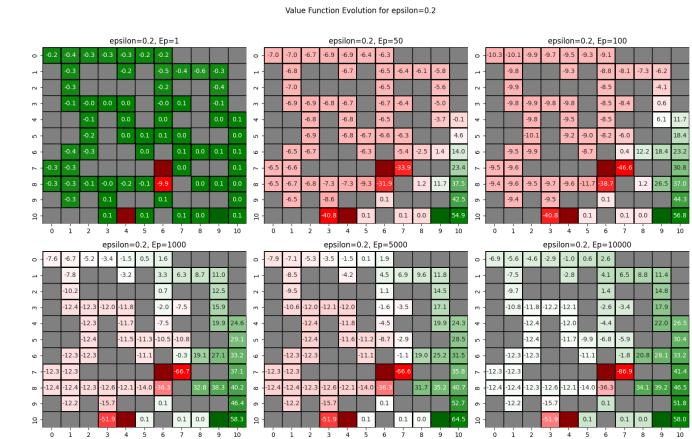


Fig. 32. Value Evaluation for $\epsilon = 0.2$

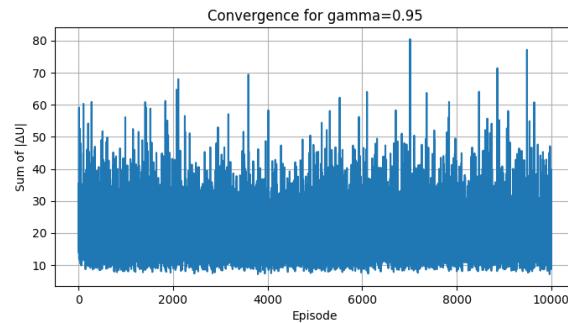


Fig. 30. Convergence for $\gamma = 0.95$

7) Value Iteration For Different Epsilon Values: In this section, we will observe value iteration for different epsilon values

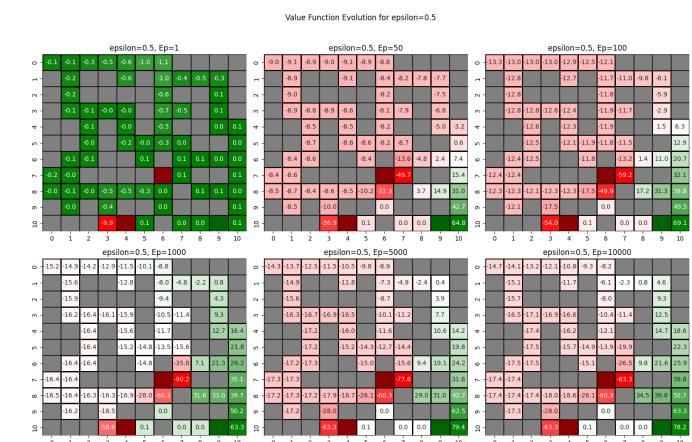


Fig. 33. Value Evaluation for $\epsilon = 0.50$

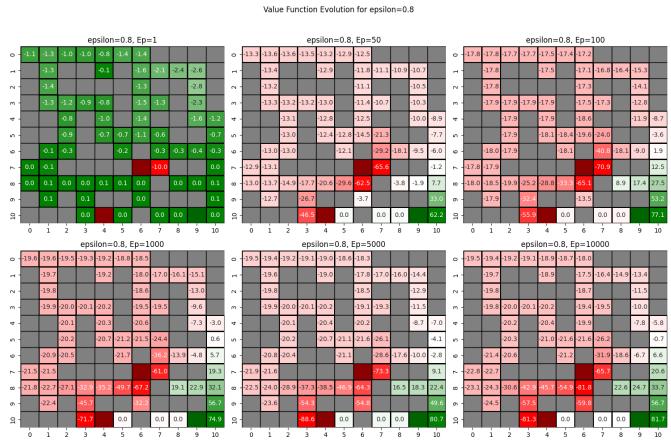


Fig. 34. Value Evaluation for $\epsilon = 0.80$

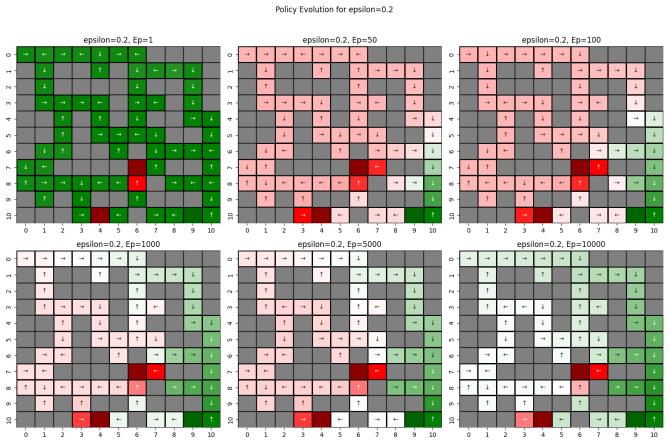


Fig. 37. Policy Evaluation for $\epsilon = 0.20$

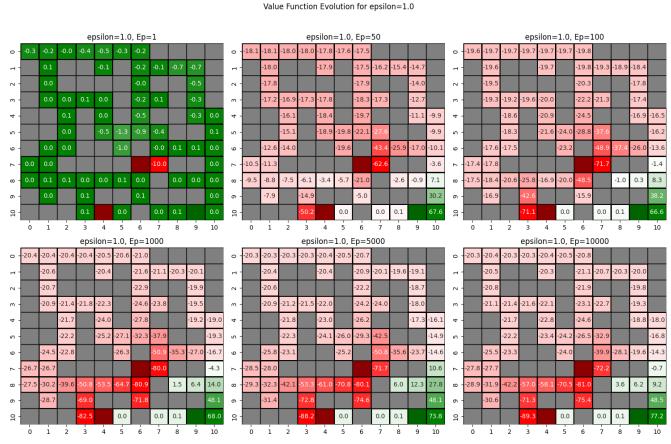


Fig. 35. Value Evaluation for $\epsilon = 1.0$

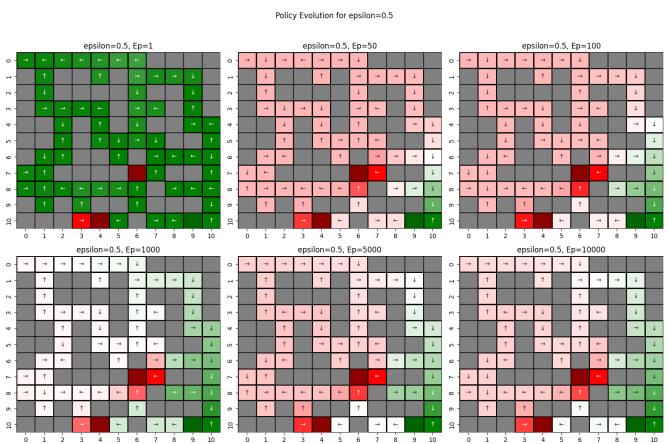


Fig. 38. Policy Evaluation for $\epsilon = 0.50$

8) *Policy Iteration For Different Epsilon Values:* In this section, we will observe policy iteration for different epsilon values

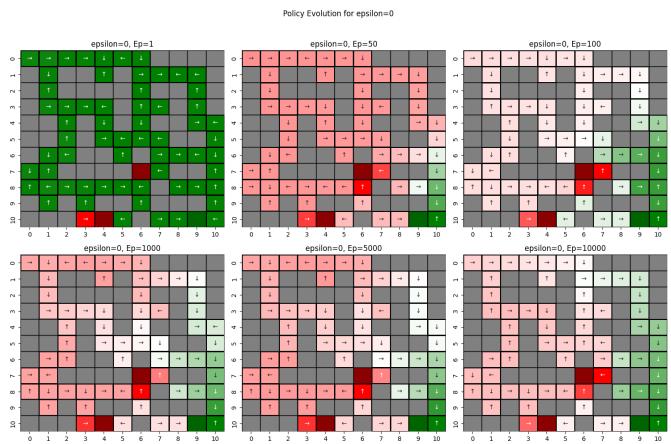


Fig. 36. Policy Evaluation for $\epsilon = 0.0$

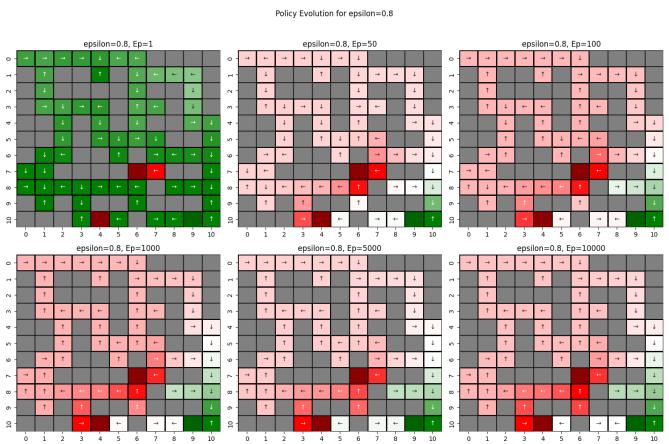


Fig. 39. Policy Evaluation for $\epsilon = 0.80$

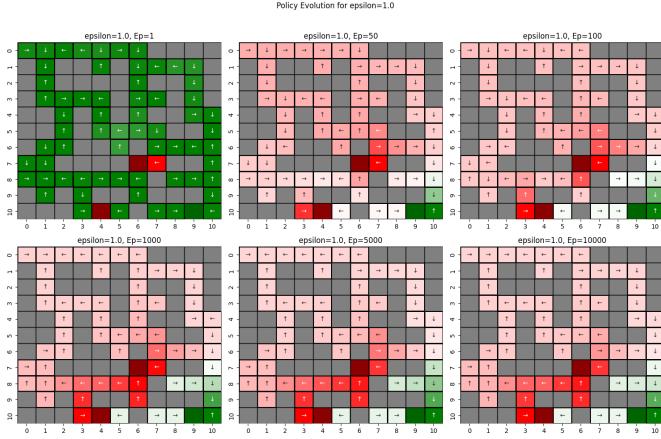


Fig. 40. Policy Evaluation for $\epsilon = 1.0$

9) Convergence For Different Epsilon Values: In this section, we will observe convergence for different epsilon values

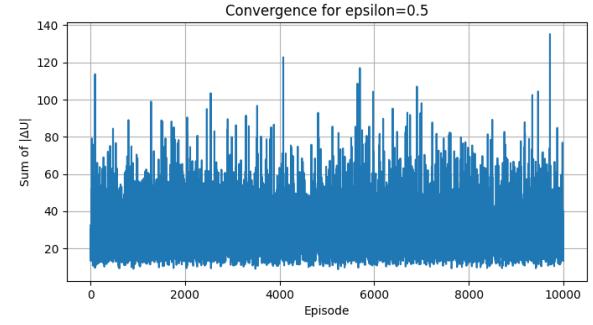


Fig. 43. Convergence for $\epsilon = 0.50$

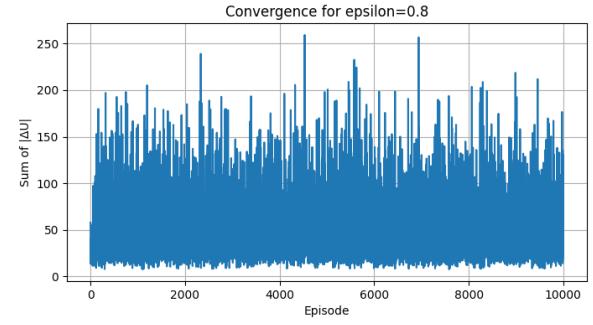


Fig. 44. Convergence for $\epsilon = 0.80$

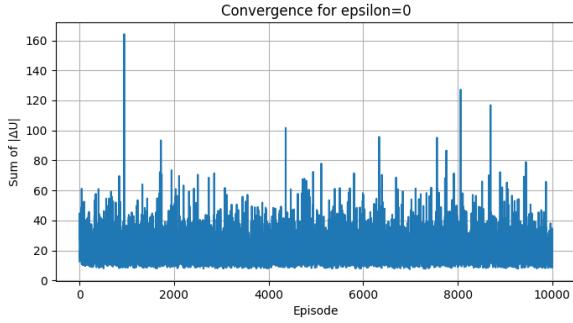


Fig. 41. Convergence for $\epsilon = 0.0$

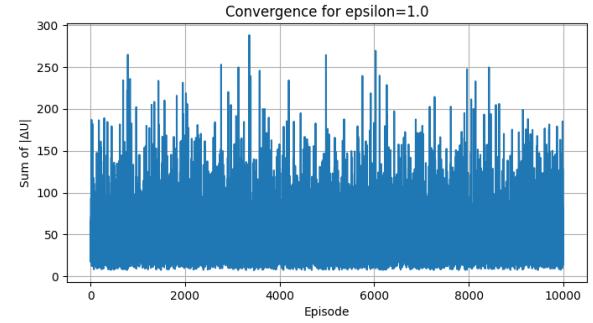


Fig. 45. Convergence for $\epsilon = 1.0$

C. Discussion

The transition probabilities affect stochasticity in the agent's movements. It simulates real-world uncertainty. The reward function penalizes each move slightly so that the agent is encouraged for efficiency. Also, we heavily penalize traps to discourage the risk of getting trapped and highly reward reaching the goal. In that way we motivate the agent to find a safe and optimal path while learning to avoid traps and redundant steps.

Early in training, utility values are random. However, over time, these values begin to show the agent's learned expectations, higher values appear near the goal and lower near traps. Heatmaps show how the agent's understanding of the environment sharpens with experience, especially after 1000+ episodes.

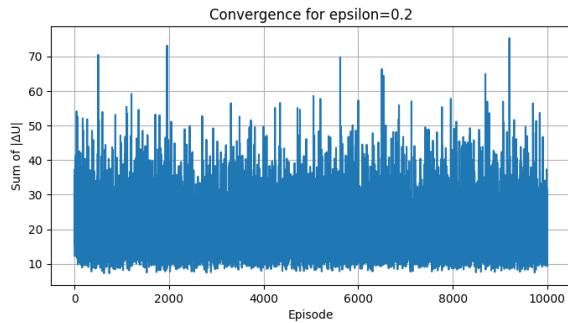


Fig. 42. Convergence for $\epsilon = 0.20$

For our cases, the utility values generally converged after around 5000 to 7000 episodes. This value depends on the parameter configuration. Convergence was calculated by observing minimal changes in utility values and a flattening in the convergence plots.

A very low learning rate, which is α is 0.001 slowed convergence. On the contrary, a very high one, which is α is 1.0 caused instability. Similarly, low discount factors (prioritized immediate rewards and discouraged planning, while high values encouraged long-term strategies but required more iterations to stabilize).

One of the challenges encountered is handling stochastic transitions and ensuring valid moves, especially near walls. I had to check in every move to ensure my move is valid.

Low ϵ causes to premature exploitation and poor exploration. Therefore, it causes suboptimal policies. On the other hand, high ϵ (like 1.0) causes exploration but delayed convergence. A balanced ϵ (e.g., 0.2) provided a good mix, enabling the agent to explore early and exploit later.

There are several ways to make the problem harder. There are Making the maze larger, increasing trap density, or introducing dynamic elements (e.g., moving obstacles) could make the problem more challenging. There are several ways to make the problem easier. There are simplifying the layout or increasing reward granularity could speed up learning.

Incorporating eligibility traces ($TD(\lambda)$) or switching to model-based methods could improve learning. Combining TD with supervised imitation learning or using function approximation (like neural networks) would also enhance performance and generalization in larger or more complex mazes.

II. DEEP Q-LEARNING

A. Experiments

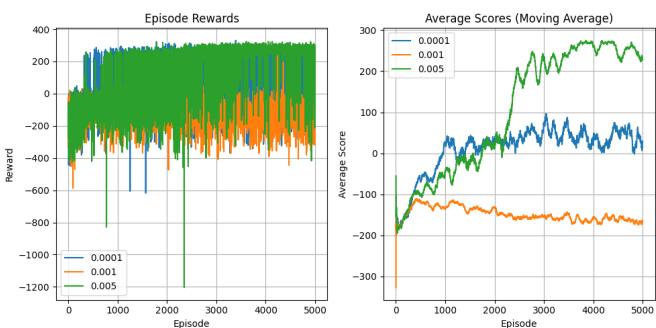


Fig. 46. Learning curves for learning rate

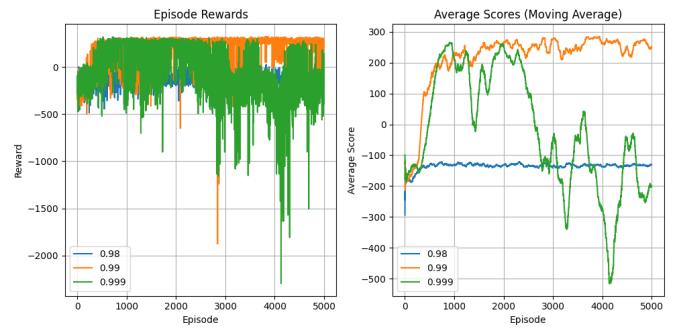


Fig. 47. Learning curves for discount factor

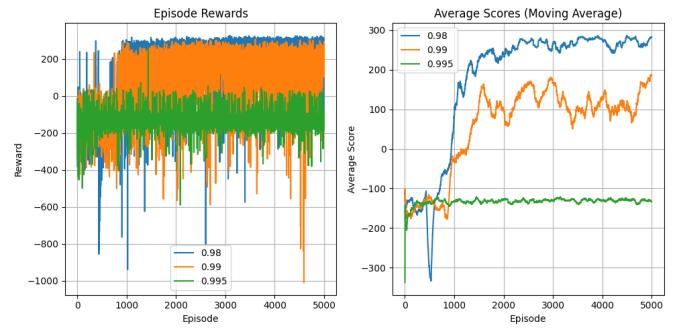


Fig. 48. Learning curves for epsilon decay

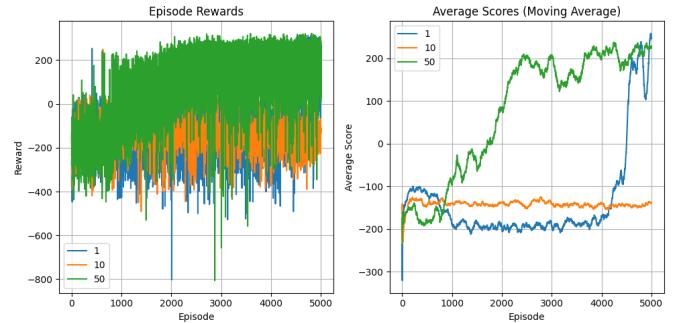


Fig. 49. Learning curves for target update frequency

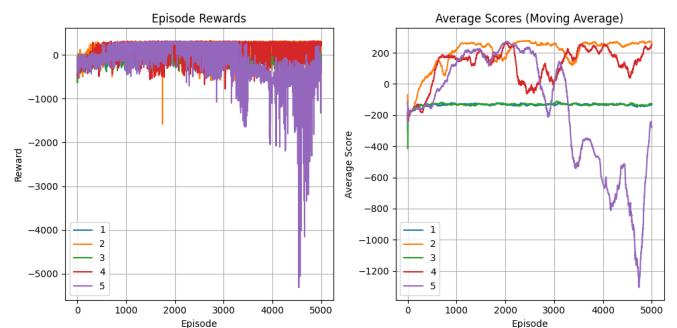


Fig. 50. Learning curves for net no

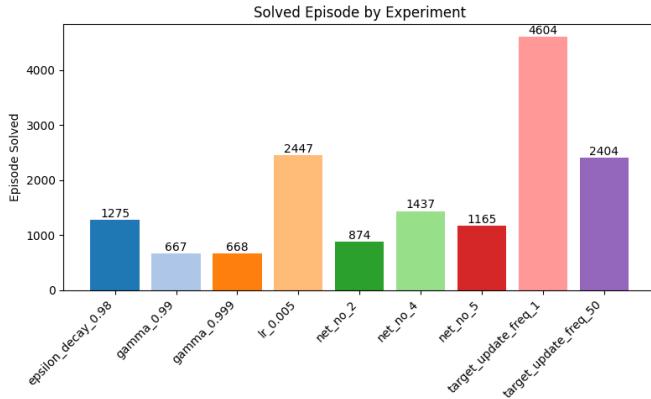


Fig. 51. Solved episodes for experiments

B. Discussion

The choice of learning rate (α) affects the training dynamics of the DQN agent. A lower learning rate generally leads to stable but slow convergence. This may be very beneficial in preventing divergence, but it may also prevent the agent from adapting quickly to new environmental patterns. On the other hand, a higher learning rate can lead to faster initial learning. However, most of the time it can introduce instability to the model due to overshooting optimal values. Among the tested values, 1e-3 offered the most balanced outcome. It achieved a better final performance. Moreover, it maintained a very good learning speed.

The discount factor (γ) affects the agent's preference for long-term versus short-term rewards. A smaller lower value prioritizes immediate rewards, which causes faster early improvements but may be bad for capturing the long-term strategies. On the other hand, a high value prioritizes future rewards, which causes better long-term planning. However, it may delay convergence due to the slower propagation of value updates. The best value for γ is 0.99. It allows for focusing on long-term goals without overriding or delaying learning progress.

The exploration decay rate (ϵ -decay) is very important for balancing exploration and exploitation. For small values, randomness reduces rapidly, which can cause the agent to settle prematurely on suboptimal policies. On the other hand, a slower decay causes the agent to sample a wider range of state-action pairs, and it may discover more optimal strategies. However, it may also slow down convergence. The experiments showed that slower decay rates tended to yield higher final performance due to more thorough policy refinement, although they required more episodes to reach the same level of stability.

The frequency of target network updates affects the stability and reliability of Q-value estimates. Frequent updates tend to increase instability because of shifting the target values rapidly. It can lead to oscillating behavior. On the other hand, infrequent updates cause the target network to lag behind the policy network. It means it results in outdated targets that slow learning progress. The default frequency of $f = 10$ is the most

effective. It is stable in the learning process by periodically refreshing the target values without causing instability.

The model of the Q-network plays a very important role in the agent's ability to approximate the action-value function effectively. Networks with fewer parameters tend to converge quickly, but capacity can be shallow. It may result in suboptimal policies. Deeper or wider architectures offer greater function approximation capability and can capture more complex patterns. However, they may overfit. Among the tested models, Net 3 delivered the best performance by balancing expressiveness and stability. It converged efficiently and was able to reach the solution criterion faster than both smaller and larger networks in most cases.

The learning rate had the most noticeable impact on the agents' overall performance out of all the hyperparameters that were investigated. Compared to both lower and higher values, a carefully selected learning rate like 1e-3 allowed for faster and more stable convergence. Furthermore, the neural network architecture significantly affected the agent's ability to generalize from experiences. Additionally, there seemed to be interplay among the hyperparameters. For example, deeper networks were more susceptible to rapid epsilon decay, while a higher learning rate necessitated more robust network architectures to preserve stability. Despite the independent testing of the parameters, these dependencies imply that tuning hyperparameters together rather than separately may be able to yield the best results.

Important information about the training dynamics under various configurations was revealed by the learning curves. Sharp increases and decreases in episode rewards, which signify instability in the learning process, are characteristic of volatile curves that are usually produced by configurations with high learning rates or very frequent target network updates. Smoother learning curves, on the other hand, were typically linked to a gradual epsilon decay schedule, steady update frequencies, and moderate learning rates. The more steady improvement in policy quality over time was represented by these smoother curves. The significance of choosing hyperparameters that strike a balance between learning speed and stability to guarantee stable performance and dependable convergence is highlighted by the variation in learning trajectories.

APPENDIX

A. Part 1

1) main.py:

```
1 import sys
2 import os
3 sys.path.append(os.path.abspath('_given'))
4
5 import numpy as np
6 import matplotlib.pyplot as plt
7 from utils import plot_value_function, plot_policy, plot_convergence # type: ignore
8
9 import gc
10
11 # Ensure results directory exists
12 os.makedirs('part1/results', exist_ok=True)
13
14 class MazeEnvironment:
15     def __init__(self):
16         # Define the maze layout, rewards, and action space
17         self.start_pos = (0, 0)
18         self.current_pos = self.start_pos
19         self.state_penalty = -1
20         self.trap_penalty = -100
21         self.goal_reward = 100
22         self.actions = {0: (-1, 0), 1: (1, 0), 2: (0, -1), 3: (0, 1)}
23         # Transition probabilities
24         self.prob_intended = 0.75
25         self.prob_opposite = 0.05
26         self.prob_perpendicular = 0.10
27
28     def reset(self):
29         self.current_pos = self.start_pos
30         return self.current_pos
31
32     def step(self, action):
33         r = np.random.rand()
34         if r < self.prob_intended:
35             chosen = action
36         elif r < self.prob_intended + self.prob_opposite:
37             chosen = {0:1, 1:0, 2:3, 3:2}[action]
38         elif r < self.prob_intended + self.prob_opposite + self.prob_perpendicular:
39             chosen = {0:2, 1:3, 2:1, 3:0}[action]
40         else:
41             chosen = {0:3, 1:2, 2:0, 3:1}[action]
42
43         move = self.actions[chosen]
44         next_state = (self.current_pos[0] + move[0], self.current_pos[1] + move[1])
45         if (0 <= next_state[0] < self.maze.shape[0] and 0 <= next_state[1] < self.maze.shape[1]
46             and self.maze[next_state] != 1):
47             self.current_pos = next_state
48
49         cell_value = self.maze[self.current_pos]
50         if cell_value == 2:
51             return self.current_pos, self.trap_penalty, True
52         elif cell_value == 3:
53             return self.current_pos, self.goal_reward, True
54         else:
55             return self.current_pos, self.state_penalty, False
56
57 class MazeTD0(MazeEnvironment):
58     def __init__(self, maze, alpha=0.001, gamma=0.95, epsilon=0.2, episodes=10000):
59         super().__init__()
60         self.maze = maze
61         self.alpha = alpha
62         self.gamma = gamma
63         self.epsilon = epsilon
64         self.episodes = episodes
65         # Initialize utilities
66         self.utility = np.random.rand(*self.maze.shape) * 0.1
67         self.utility[maze == 1] = np.nan
68
69     def choose_action(self, state):
70         valid = []
```

```

71     for a, move in self.actions.items():
72         ns = (state[0] + move[0], state[1] + move[1])
73         if 0 <= ns[0] < self.maze.shape[0] and 0 <= ns[1] < self.maze.shape[1] and self.maze[ns] != 1:
74             valid.append(a)
75     if not valid:
76         return None
77     if np.random.rand() < self.epsilon:
78         return np.random.choice(valid)
79     best, best_val = None, -np.inf
80     for a in valid:
81         move = self.actions[a]
82         ns = (state[0] + move[0], state[1] + move[1])
83         val = self.utility[ns]
84         if val > best_val:
85             best, best_val = a, val
86     return best
87
88 def update_utility_value(self, s, r, ns):
89     curr = self.utility[s]
90     new = self.utility[ns]
91     self.utility[s] = curr + self.alpha * (r + self.gamma * new - curr)
92
93 def run_episodes(self, snapshot_episodes=None, max_steps=5000000):
94     if snapshot_episodes is None:
95         snapshot_episodes = []
96     snapshots = {}
97     convergence = np.zeros(self.episodes)
98
99     reached_step = 0 # checking if we hit max_steps
100
101    for ep in range(1, self.episodes + 1):
102        state = self.reset()
103        diff = 0
104        steps = 0
105        done = False
106
107        while not done and steps < max_steps:
108            act = self.choose_action(state)
109            if act is None:
110                break
111            old = self.utility[state]
112            next_state, reward, done = self.step(act)
113            self.update_utility_value(state, reward, next_state)
114            diff += abs(self.utility[state] - old)
115            state = next_state
116            steps += 1
117
118            reached_step = max(reached_step , steps)
119            convergence[ep - 1] = diff
120            if ep in snapshot_episodes:
121                snapshots[ep] = self.utility.copy()
122            print(f"Reached maximum step is: {reached_step}")
123    return snapshots, convergence
124
125
126
127 def save_results(param, val, snapshots, convergence):
128     base = f"part1/results/{param}_{val}"
129     os.makedirs(base, exist_ok=True)
130
131     for ep, u in snapshots.items():
132         np.save(f"{base}/utility_ep{ep}.npy", u)
133
134     np.save(f"{base}/convergence.npy", np.array(convergence))
135
136
137
138 if __name__ == '__main__':
139     # Define maze
140     maze = np.array([
141         [0,0,0,0,0,0,1,1,1,1],
142         [1,0,1,1,0,1,0,0,0,1],
143         [1,0,1,1,1,1,0,1,1,0,1],
144         [1,0,0,0,0,1,0,0,1,0,1],

```

```

145 [1,1,0,1,0,1,0,1,1,0,0],
146 [1,1,0,1,0,0,0,0,1,1,0],
147 [1,0,0,1,1,0,1,0,0,0,0],
148 [0,0,1,1,1,1,2,0,1,1,0],
149 [0,0,0,0,0,0,0,1,0,0,0],
150 [1,0,1,0,1,1,0,1,1,1,0],
151 [1,1,1,0,2,0,1,0,0,3,0]
152 ])
153
154 # Hyperparameter experiments
155 experiments = {
156     'alpha': {'values': [0.001, 0.01, 0.1, 0.5, 1.0], 'default': 0.1},
157     'gamma': {'values': [0.10, 0.25, 0.50, 0.75, 0.95], 'default': 0.95},
158     'epsilon': {'values': [0, 0.2, 0.5, 0.8, 1.0], 'default': 0.2}
159 }
160 episodes = 10000
161 snapshots_to_keep = [1, 50, 100, 1000, 5000, 10000]
162
163 for param, config in experiments.items():
164     for val in config['values']:
165         # Set parameters
166         alpha = experiments['alpha']['default'] if param != 'alpha' else val
167         gamma = experiments['gamma']['default'] if param != 'gamma' else val
168         epsilon = experiments['epsilon']['default'] if param != 'epsilon' else val
169
170         print(f"Selected values are: alpha={alpha}, gamma={gamma}, epsilon={epsilon}")
171
172         # Run TD(0)
173         agent = MazeTD0(maze, alpha=alpha, gamma=gamma, epsilon=epsilon, episodes=episodes)
174         snapshots, convergence = agent.run_episodes(snapshots_to_keep)
175
176         # Plot and save value function snapshots
177         fig, axes = plt.subplots(2, 3, figsize=(15, 10))
178         axes = axes.flatten()
179         for ax, ep in zip(axes, snapshots_to_keep):
180             plot_value_function(snapshots[ep], maze, ax=ax, title=f"{param}={val}, Ep={ep}")
181             fig.suptitle(f"Value Function Evolution for {param}={val}")
182             plt.tight_layout(rect=[0, 0, 1, 0.96])
183             fig.savefig(f"part1/results/{param}_{val}_value_evolution.png")
184             plt.close(fig)
185
186         # Plot and save policy snapshots
187         fig, axes = plt.subplots(2, 3, figsize=(15, 10))
188         axes = axes.flatten()
189         for ax, ep in zip(axes, snapshots_to_keep):
190             plot_policy(snapshots[ep], maze, ax=ax, title=f"{param}={val}, Ep={ep}")
191             fig.suptitle(f"Policy Evolution for {param}={val}")
192             plt.tight_layout(rect=[0, 0, 1, 0.96])
193             fig.savefig(f"part1/results/{param}_{val}_policy_evolution.png")
194             plt.close(fig)
195
196         # Plot and save convergence
197         fig, ax = plt.subplots(figsize=(8, 4))
198         plot_convergence(convergence, ax=ax, title=f"Convergence for {param}={val}")
199         fig.savefig(f"part1/results/{param}_{val}_convergence.png")
200         plt.close(fig)

```

B. Part 2

1) main.py:

```
1 import gymnasium as gym
2 import numpy as np
3 import torch
4 from collections import deque
5 import random
6 import json
7
8 from models import QNetwork_1, QNetwork_2, QNetwork_3, QNetwork_4, QNetwork_5
9
10 # Define a replay memory to store experiences:
11 class ReplayMemory:
12     def __init__(self, capacity):
13         self.capacity = capacity
14         self.buffer = []
15         self.position = 0
16     def push(self, state, action, reward, next_state, done):
17         if len(self.buffer) < self.capacity:
18             self.buffer.append((state, action, reward, next_state, done))
19         else:
20             self.buffer[self.position] = (state, action, reward, next_state, done)
21         self.position = (self.position + 1) % self.capacity
22     def sample(self, batch_size):
23         return random.sample(self.buffer, batch_size)
24     def __len__(self):
25         return len(self.buffer)
26
27 # Define the DQN Agent encapsulating the networks, memory, and training procedure:
28 class DQNAgent:
29     def __init__(self,
30                  state_dim,
31                  action_dim,
32                  Net_No=3,
33                  memory_size=50000,
34                  batch_size=64,
35                  gamma=0.99,
36                  alpha=1e-3,
37                  epsilon_start=1.0,
38                  epsilon_min=0.01,
39                  epsilon_decay=0.995,
40                  target_update_freq=10):
41
42         self.state_dim = state_dim
43         self.action_dim = action_dim
44         self.gamma = gamma
45         self.alpha = alpha
46         self.batch_size = batch_size
47         self.epsilon = epsilon_start
48         self.epsilon_min = epsilon_min
49         self.epsilon_decay = epsilon_decay
50         self.memory = ReplayMemory(memory_size)
51
52     match Net_No:
53         case 1:
54             self.policy_net = QNetwork_1(state_dim, action_dim)
55             self.target_net = QNetwork_1(state_dim, action_dim)
56         case 2:
57             self.policy_net = QNetwork_2(state_dim, action_dim)
58             self.target_net = QNetwork_2(state_dim, action_dim)
59         case 3:
60             self.policy_net = QNetwork_3(state_dim, action_dim)
61             self.target_net = QNetwork_3(state_dim, action_dim)
62         case 4:
63             self.policy_net = QNetwork_4(state_dim, action_dim)
64             self.target_net = QNetwork_4(state_dim, action_dim)
65         case 5:
66             self.policy_net = QNetwork_5(state_dim, action_dim)
67             self.target_net = QNetwork_5(state_dim, action_dim)
68         case _:
69             print("Error on Net_No: missing")
70
71     self.target_net.load_state_dict(self.policy_net.state_dict())
72
```

```

73     self.target_net.eval()
74     self.optimizer = torch.optim.Adam(self.policy_net.parameters(), lr=alpha)
75     self.target_update_freq = target_update_freq
76     self.solved_score = 200.0          #target average reward for solved
77     self.solved_window = 100          #number of episodes to average for solved check
78
79     def get_action(self, state):
80         # Select an action for the given state using epsilon-greedy policy.
81         if random.random() < self.epsilon:
82             return random.randrange(self.action_dim)
83         else:
84             state_v = torch.tensor(state, dtype=torch.float32).unsqueeze(0)
85             with torch.no_grad():
86                 q_values = self.policy_net(state_v)
87             return int(q_values.argmax(dim=1)[0])
88
89     def train_step(self):
90         #Perform one training step.
91         if len(self.memory) < self.batch_size:
92             return # Not enough samples to train yet.
93         batch = self.memory.sample(self.batch_size)
94         states = torch.tensor(np.array([exp[0] for exp in batch]), dtype=torch.float32)
95         actions = torch.tensor(np.array([exp[1] for exp in batch]), dtype=torch.int64).unsqueeze(1)
96         rewards = torch.tensor(np.array([exp[2] for exp in batch]), dtype=torch.float32).unsqueeze(1)
97         next_states = torch.tensor(np.array([exp[3] for exp in batch]), dtype=torch.float32)
98         dones = torch.tensor(np.array([exp[4] for exp in batch]), dtype=torch.float32).unsqueeze(1)
99         curr_Q = self.policy_net(states).gather(1, actions)
100        next_Q = self.target_net(next_states).max(1)[0].detach().unsqueeze(1)
101
102        # Compute target Q
103        target_Q = rewards + (1 - dones) * (self.gamma * next_Q)
104        #target_Q = self.alpha * (rewards + self.gamma * next_Q - curr_Q)
105
106        # Compute loss
107        loss = torch.nn.functional.mse_loss(curr_Q, target_Q)
108
109        # Optimize
110        self.optimizer.zero_grad()
111        loss.backward()
112        self.optimizer.step()
113
114    def update_target(self):
115        # Update target network to match policy network.
116        self.target_net.load_state_dict(self.policy_net.state_dict())
117
118    def decay_epsilon(self):
119        # Decay exploration rate after each episode.
120        self.epsilon = max(self.epsilon_min, self.epsilon * self.epsilon_decay)
121
122
123    def run_experiment(config, output_path):
124        # Main training loop for a single experiment
125        env = gym.make("LunarLander-v3")
126        state_dim = env.observation_space.shape[0]
127        action_dim = env.action_space.n
128
129        #print(f"lr={config['lr']} | g={config['gamma']} | ed={config['epsilon_decay']} | f={config['target_update_freq']} | NN={config['net_no']}")
130
131        agent = DQNAgent(state_dim,
132                           action_dim,
133                           Net_No=config['net_no'],
134                           gamma=config['gamma'],
135                           epsilon_decay=config['epsilon_decay'],
136                           target_update_freq=config['target_update_freq'],
137                           )
138        num_episodes = config['num_episodes']
139
140        rewards_history = []
141        avg_scores = []
142        solved_episode = None
143        scores_window = deque(maxlen=agent.solved_window)
144
145        for episode in range(1, num_episodes + 1):

```

```

146     state, _ = env.reset()
147     total_reward = 0
148     done = False
149
150     for t in range(1000): # limit max steps per episode to avoid long runs
151         action = agent.get_action(state)
152         next_state, reward, terminated, truncated, _ = env.step(action)
153         done = terminated or truncated
154         agent.memory.push(state, action, reward, next_state, done)
155         agent.train_step()
156         state = next_state
157         total_reward += reward
158         if done:
159             break
160
161     # Episode end
162     rewards_history.append(total_reward)
163     scores_window.append(total_reward)
164     # Compute moving average
165     avg = np.mean(scores_window)
166     avg_scores.append(avg)
167     # Check solved
168     if solved_episode is None and len(scores_window) == agent.solved_window and avg >= agent.
169         solved_score:
170             solved_episode = episode
171     # Epsilon decay
172     agent.decay_epsilon()
173     # Update target network
174     if episode % agent.target_update_freq == 0:
175         agent.update_target()
176     print(f"Episode {episode} | Reward: {total_reward:.2f} | Avg: {avg:.2f} | Epsilon: {agent.epsilon
177         :.3f}")
178
179     env.close()
180     # Save results
181     results = {
182         'episode_rewards': rewards_history,
183         'average_scores': avg_scores,
184         'hyperparameters': config,
185         'solved_episode': solved_episode
186     }
187     with open(output_path, 'w') as f:
188         json.dump(results, f)
189     print(f"Results saved to {output_path}")
190
191     # Example usage
192     if __name__ == "__main__":
193         base_config = {
194             'lr': 1e-3,
195             'gamma': 0.99,
196             'epsilon_decay': 0.995,
197             'target_update_freq': 10,
198             'net_no': 3,
199             'num_episodes': 5000
200         }
201
202         param_grid = {
203             'lr': [1e-4, 1e-3, 5e-3],
204             'gamma': [0.98, 0.99, 0.999],
205             'epsilon_decay': [0.98, 0.99, 0.995],
206             'target_update_freq': [1, 10, 50],
207             'net_no': [1, 2, 3, 4, 5]
208         }
209
210         for param, values in param_grid.items():
211             for val in values:
212                 config = base_config.copy()
213                 config[param] = val
214                 save_path = f'part2/results/{param}_{val}.json'
215                 print(f"Running experiment with {param}={val}")
216                 run_experiment(config, save_path)

```

2) main.py:

```

1 import torch
2
3 # Q-Network FC-128, ReLU
4 class QNetwork_1(torch.nn.Module):
5     def __init__(self, state_dim, action_dim, hidden_dim=128):
6         super(QNetwork_1, self).__init__()
7         self.fc1 = torch.nn.Linear(state_dim, hidden_dim)
8         self.fc2 = torch.nn.Linear(hidden_dim, action_dim)
9
10    def forward(self, x):
11        x = torch.relu(self.fc1(x))
12        x = torch.relu(self.fc2(x))
13        return x
14
15
16    # Q-Network FC-64, ReLU and FC-64, ReLU
17    class QNetwork_2(torch.nn.Module):
18        def __init__(self, state_dim, action_dim, hidden_dim1=64, hidden_dim2=64):
19            super(QNetwork_2, self).__init__()
20            self.fc1 = torch.nn.Linear(state_dim, hidden_dim1)
21            self.fc2 = torch.nn.Linear(hidden_dim1, hidden_dim2)
22            self.fc3 = torch.nn.Linear(hidden_dim2, action_dim)
23
24        def forward(self, x):
25            x = torch.relu(self.fc1(x))
26            x = torch.relu(self.fc2(x))
27            x = torch.relu(self.fc3(x))
28            return x
29
30
31    # Q-Network FC-128, ReLU and FC-128, ReLU
32    class QNetwork_3(torch.nn.Module):
33        def __init__(self, state_dim, action_dim, hidden_dim1=128, hidden_dim2=128):
34            super(QNetwork_3, self).__init__()
35            self.fc1 = torch.nn.Linear(state_dim, hidden_dim1)
36            self.fc2 = torch.nn.Linear(hidden_dim1, hidden_dim2)
37            self.fc3 = torch.nn.Linear(hidden_dim2, action_dim)
38
39        def forward(self, x):
40            x = torch.relu(self.fc1(x))
41            x = torch.relu(self.fc2(x))
42            x = torch.relu(self.fc3(x))
43            return x
44
45
46    # Q-Network FC-128, ReLU and FC-128, ReLU and FC-128, ReLU
47    class QNetwork_4(torch.nn.Module):
48        def __init__(self, state_dim, action_dim, hidden_dim1=128, hidden_dim2=128, hidden_dim3=128):
49            super(QNetwork_4, self).__init__()
50            self.fc1 = torch.nn.Linear(state_dim, hidden_dim1)
51            self.fc2 = torch.nn.Linear(hidden_dim1, hidden_dim2)
52            self.fc3 = torch.nn.Linear(hidden_dim2, hidden_dim3)
53            self.fc4 = torch.nn.Linear(hidden_dim3, action_dim)
54
55        def forward(self, x):
56            x = torch.relu(self.fc1(x))
57            x = torch.relu(self.fc2(x))
58            x = torch.relu(self.fc3(x))
59            x = torch.relu(self.fc4(x))
60            return x
61
62
63    # Q-Network FC-256, ReLU and FC-256, ReLU
64    class QNetwork_5(torch.nn.Module):
65        def __init__(self, state_dim, action_dim, hidden_dim1=256, hidden_dim2=256):
66            super(QNetwork_5, self).__init__()
67            self.fc1 = torch.nn.Linear(state_dim, hidden_dim1)
68            self.fc2 = torch.nn.Linear(hidden_dim1, hidden_dim2)
69            self.fc3 = torch.nn.Linear(hidden_dim2, action_dim)
70
71        def forward(self, x):
72            x = torch.relu(self.fc1(x))
73            x = torch.relu(self.fc2(x))

```

```
74     x = torch.relu(self.fc3(x))
75     return x
```

```

3) main.py: -----
1 import sys
2 import os
3 sys.path.append(os.path.abspath('_given'))
4 from utils import plot_learning_curves, plot_solved_episodes # type: ignore
5
6
7 def main():
8     # Directory containing experiment JSON results
9     base_dir = os.path.join(os.path.dirname(__file__), 'results')
10
11    # Hyperparameter experiments and their values (must match JSON filenames)
12    experiments = {
13        'lr': [0.0001, 0.001, 0.005],
14        'gamma': [0.98, 0.99, 0.999],
15        'epsilon_decay': [0.98, 0.99, 0.995],
16        'target_update_freq': [1, 10, 50],
17        'net_no': [1, 2, 3, 4, 5]
18    }
19
20    # Plot learning curves for each hyperparameter
21    for param, values in experiments.items():
22        # Construct list of JSON paths for this parameter
23        json_paths = []
24        labels = []
25        for val in values:
26            # Format value to match filename
27            val_str = str(val)
28            json_file = f'{param}_{val_str}.json'
29            path = os.path.join(base_dir, json_file)
30            if os.path.isfile(path):
31                json_paths.append(path)
32                labels.append(str(val))
33            else:
34                print(f"Warning: missing result file {path}")
35
36        if not json_paths:
37            print(f"No data for parameter '{param}', skipping.")
38            continue
39
40        # Output image filename
41        out_png = f"part2/results_fig/learning_curves_{param}.png"
42        # Generate and save figure
43        plot_learning_curves(
44            json_paths=json_paths,
45            labels=labels,
46            output_file=out_png
47        )
48        print(f"Saved learning curves for '{param}' at: {out_png}")
49
50    # Plot bar chart of solved episodes for all experiments
51    # Gather all JSON result files
52    all_json = []
53    for fname in os.listdir(base_dir):
54        if fname.endswith('.json'):
55            all_json.append(os.path.join(base_dir, fname))
56
57    if not all_json:
58        print("No JSON result files found, skipping solved-episode bar chart.")
59        return
60
61    all_labels = [
62        os.path.splitext(os.path.basename(path))[0]
63        for path in all_json
64    ]
65
66    # Output bar chart filename
67    out_bar = 'part2/results_fig/solved_episodes.png'
68
69    # Generate and save bar chart
70    plot_solved_episodes(
71        json_paths=all_json,
72        labels = all_labels,
73        output_file=out_bar

```

```
74     )
75     print(f"Saved_solved_episodes_bar_chart_at:{out_bar}")
76
77
78 if __name__ == '__main__':
79     main()
```