

DA6400 : Reinforcement Learning  
Programming Assignment #1  
Report

Ritabrata Mandal  
EE24E009

March 30, 2025

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Environments . . . . .	3
1.2	Algorithms . . . . .	3
<b>2</b>	<b>Implementation</b>	<b>4</b>
2.1	<b>CartPole-v1</b> . . . . .	4
2.1.1	Code Snippets . . . . .	4
2.1.2	SARSA Hyper-Parameter Tuning . . . . .	8
2.1.3	SARSA best 3 results . . . . .	8
2.1.4	Q-Learning hyper-parameter tuning . . . . .	10
2.1.5	Q-Learning best 3 results . . . . .	10
2.1.6	Result(SARSA vs Q-Learning) . . . . .	12
2.2	<b>MountainCar-v0</b> . . . . .	12
2.2.1	Code Snippets . . . . .	12
2.2.2	SARSA hyper-parameter tuning . . . . .	15
2.2.3	SARSA best 3 results . . . . .	15
2.2.4	Q-Learning hyper-parameter tuning . . . . .	17
2.2.5	Q-Learning best 3 results . . . . .	17
2.2.6	Result(SARSA vs Q-Learning) . . . . .	19
2.3	<b>MiniGrid-Dynamic-Obstacles-5x5-v0</b> . . . . .	19
2.3.1	Code Snippets . . . . .	19
2.3.2	SARSA hyper-parameter tuning . . . . .	21
2.3.3	SARSA best 3 results . . . . .	22
2.3.4	Q-Learning hyper-parameter tuning . . . . .	24
2.3.5	Q-Learning best 3 results . . . . .	24
2.3.6	Result(SARSA vs Q-Learning) . . . . .	26
<b>3</b>	<b>Github link</b>	<b>26</b>
<b>4</b>	<b>Steps to recreate</b>	<b>26</b>

# 1 Introduction

## 1.1 Environments

In this programming task, we are utilize the following [Gymnasium environments](#) for training and evaluating your policies. The links associated with the environments contain descriptions of each environment.

- [CartPole-v1](#) : A pole is attached by an un-actuated joint to a cart, which moves along a frictionless track. The pendulum is placed upright on the cart and the goal is to balance the pole by applying forces in the left and right direction on the cart.
- [MountainCar-v0](#) : The Mountain Car MDP is a deterministic MDP that consists of a car placed stochastically at the bottom of a sinusoidal valley, with the only possible actions being the accelerations that can be applied to the car in either direction. The goal of the MDP is to strategically accelerate the car to reach the goal state on top of the right hill. There are two versions of the mountain car domain in gymnasium: one with *discrete actions* and one with *continuous*. This version is the one with discrete actions.
- [MiniGrid-Dynamic-Obstacles-5x5-v0](#) : This environment is an empty room with moving obstacles. The goal of the agent is to reach the green goal square without colliding with any obstacle. A large penalty is subtracted if the agent collides with an obstacle and the episode finishes. This environment is useful to test Dynamic Obstacle Avoidance for mobile robots with Reinforcement Learning in Partial Observability.

## 1.2 Algorithms

Training each of the below algorithms and assessing their comparative performance.

- **SARSA**→ with  $\epsilon$ -greedy exploration
- **Q-Learning**→ with Softmax exploration

## 2 Implementation

### 2.1 CartPole-v1

#### 2.1.1 Code Snippets

```
1 def Qtable(state_space, action_space, bin_size=100):
2
3     bins = np.zeros((state_space.shape[0], bin_size))
4
5     bins[0] = np.linspace(-4.8, 4.8, bin_size)
6     bins[1] = np.linspace(-4, 4, bin_size)
7     bins[2] = np.linspace(-0.42, 0.42, bin_size)
8     bins[3] = np.linspace(-4, 4, bin_size)
9
10    q_table = np.zeros((bin_size, bin_size, bin_size, bin_size,
11                        action_space.n))
12
13    return q_table, bins
```

Listing 1: Q-table

```
1 def discretize_state(state_space, bins):
2
3     state_discrete = np.zeros(state_space.shape)
4
5     for i in range(state_space.shape[0]):
6         state_discrete[i] = np.digitize(state_space[i], bins[i])
7
8     return state_discrete.astype(np.int32)
```

Listing 2: Discretized states

```
1 class EpsilonGreedyPolicy:
2     def __init__(self, epsilon, q_table, env):
3         self.epsilon = epsilon
4         self.q_table = q_table
5         self.env = env
6
7     def get_action(self, state):
8
9         if np.random.rand() < self.epsilon:
10             action = self.env.action_space.sample()
11         else:
12             action = np.argmax(self.q_table[state[0], state[1],
13                                 state[2], state[3]])
14
15         return action
```

Listing 3: epsilon-greedy action selection

```
1 class SarasLearner:
2     def __init__(self, alpha, gamma, epsilon, q_table, bins,
3      env, seed):
4         self.alpha = alpha
5         self.gamma = gamma
6         self.epsilon = epsilon
7         self.q_table = q_table
8         self.env = env
9         self.bins = bins
10        self.seed = seed
11        self.policy = EpsilonGreedyPolicy(self.epsilon, self.
12        q_table, self.env)
13
14    def compute_td_error(self, state, action, next_state,
15    next_action, reward):
16
17        return reward + self.gamma * self.q_table[next_state[0], next_state[1], next_state[2], next_state[3],
18        next_action] - \
19            self.q_table[state[0], state[1], state[2], state[3],
20            action]
21
22    def update_q_table(self, state, action, td_error):
23
24        self.q_table[state[0], state[1], state[2], state[3],
25        action] += self.alpha * td_error
26
27    def learn(self, num_episodes, num_steps):
28        reward_list = []
29        for episode in range(num_episodes):
30            state, _ = self.env.reset(seed=self.seed)
31            state_discrete = discretize_state(state, self.bins)
32            action = self.policy.get_action(state_discrete)
33            total_reward = 0
34
35            for step in range(num_steps):
36                next_state, reward, done = self.env.step(action
37                )[:3]
38                next_state_discrete = discretize_state(
39                next_state, self.bins)
40                next_action = self.policy.get_action(
41                next_state_discrete) # Select next action
42
43                td_error = self.compute_td_error(state_discrete
44                , action, next_state_discrete, next_action, reward)
45                self.update_q_table(state_discrete, action,
46                td_error)
47
48                state_discrete, action = next_state_discrete,
```

```

38     next_action
39         total_reward += reward
40
41         if done:
42             break
43
44         print(f"Seed: {self.seed} Episode: {episode + 1}/{num_episodes}, Total Reward: {total_reward}")
45         reward_list.append(total_reward)
46
47     return reward_list

```

Listing 4: SARSA implementation

```

1 class SoftmaxPolicy:
2     def __init__(self, temperature, q_table, env):
3         self.temperature = temperature
4         self.q_table = q_table
5         self.env = env
6
7     def get_action(self, state):
8         q_values = self.q_table[state[0], state[1], state[2],
9         state[3]]
10        max_q = np.max(q_values)
11        exp_q = np.exp((q_values - max_q) / self.temperature)
12        probabilities = exp_q / np.sum(exp_q)
13        action = np.random.choice(len(q_values), p=
probabilities)
14        return action

```

Listing 5: Softmax action selection

```

1 class QLearner:
2
3     def __init__(self, alpha, gamma, temperature, q_table, bins,
4     , env, seed):
5         self.alpha = alpha
6         self.gamma = gamma
7         self.temperature = temperature
8         self.q_table = q_table
9         self.bins = bins
10        self.env = env
11        self.seed = seed
12        self.policy = SoftmaxPolicy(self.temperature, self.
q_table, self.env)
13
14    def compute_td_error(self, state, action, next_state,
reward):

```

```

15     best_next_action_value = np.max(self.q_table[next_state
16         [0], next_state[1], next_state[2], next_state[3], :])
17     return reward + self.gamma * best_next_action_value -
18         self.q_table[state[0], state[1], state[2], state[3], action
19     ]
20
21     def update_q_table(self, state, action, td_error):
22
23
24         self.q_table[state[0], state[1], state[2], state[3],
25             action] += self.alpha * td_error
26
27     def learn(self, num_episodes, num_steps):
28
29         reward_list = []
30
31         for episode in range(num_episodes):
32             state, _ = self.env.reset(seed=self.seed)
33             state_discrete = discretize_state(state, self.bins)
34             total_reward = 0
35
36             for step in range(num_steps):
37                 action = self.policy.get_action(state_discrete)
38                 next_state, reward, done = self.env.step(action
39                     )[:3]
40                 next_state_discrete = discretize_state(
41                     next_state, self.bins)
42
43                 td_error = self.compute_td_error(state_discrete
44                     , action, next_state_discrete, reward)
45                 self.update_q_table(state_discrete, action,
46                     td_error)
47
48                 state_discrete = next_state_discrete
49                 total_reward += reward
50
51                 if done:
52                     break
53
54             print(f"Episode: {episode + 1}/{num_episodes},
55 Total Reward: {total_reward}")
56             reward_list.append(total_reward)
57
58         return reward_list

```

Listing 6: Q-Learning implementation

### 2.1.2 SARSA Hyper-Parameter Tuning

Using Weights & Biases (wandb) with a sweep method based on Bayesian optimization, we identified the best-performing hyper-parameters. Specifically, we set  $\alpha \in [0.1, 0.5]$  and  $\epsilon \in [0.01, 0.15]$ , and ran 4000 episodes while minimizing the regret, defined as  $195 -$  (all-time average return). See the wandb report on this environment here. Additionally, Figure 1 displays the results from 50 sweeps, illustrating the relationship between  $\alpha$ ,  $\epsilon$ , and the regret.

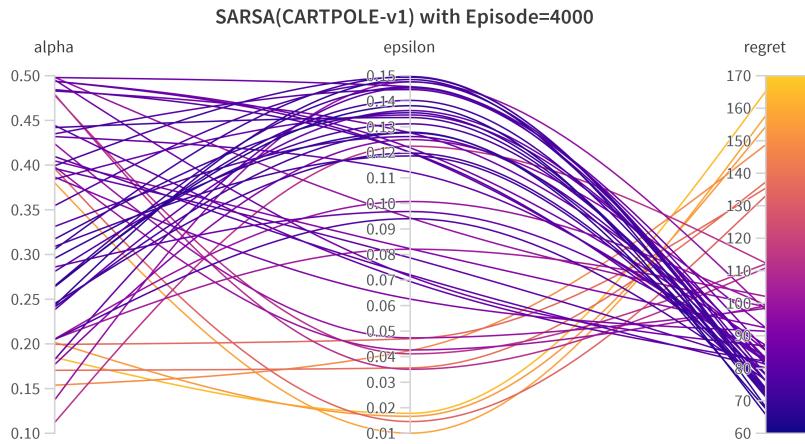
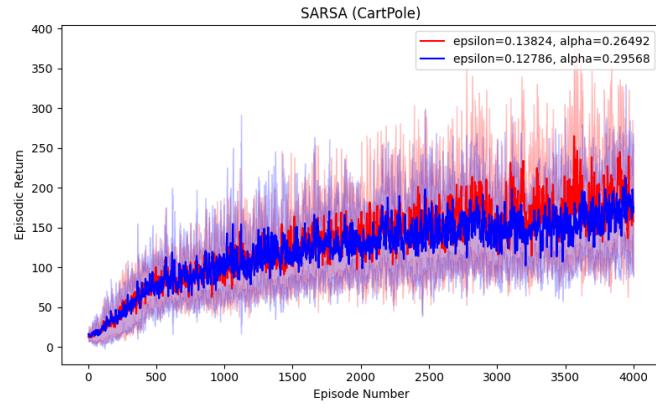


Figure 1: SARSA hyper-parameter sweeps for  $\alpha$  and  $\epsilon$  in CartPole-v1(4000 episodes), with lines color-coded by regret.

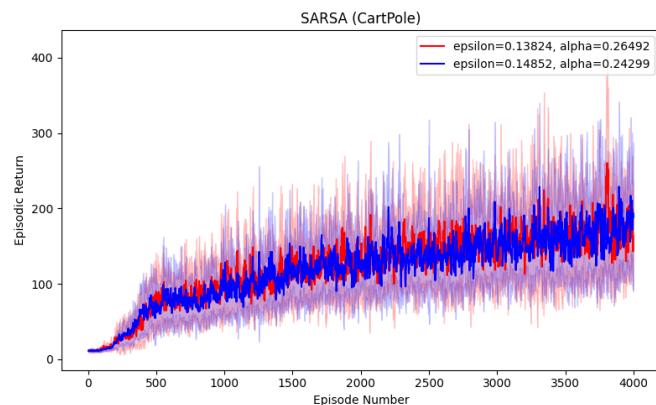
### 2.1.3 SARSA best 3 results

	$\alpha$	$\epsilon$	$\gamma$	regret
1	0.26492	0.13824	0.99	65.7379
2	0.29568	0.12786	0.99	67.5396
3	0.24299	0.14852	0.99	68.09645

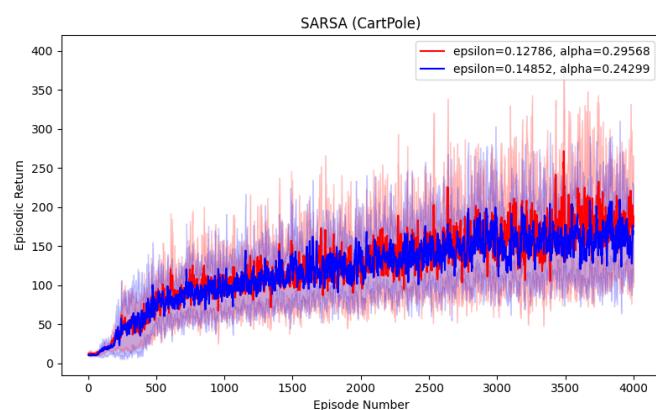
Figure 2 include the episode vs reward plot for best three hyper-parameter combination.



(a) 1 vs 2



(b) 1 vs 3



(c) 2 vs 3

Figure 2: Performance comparison of the three best SARSA parameter configurations: (a) 1 vs 2, (b) 1 vs 3, and (c) 2 vs 3.

#### 2.1.4 Q-Learning hyper-parameter tuning

Using Weights & Biases (wandb) with a sweep method based on Bayesian optimization, we identified the best-performing hyper-parameters. Specifically, we set  $\alpha \in [0.1, 0.5]$  and  $\tau \in [0.1, 1.5]$ , and ran 4000 episodes while minimizing the regret, defined as 195 – (all-time average return). See the wandb report on this environment here. Additionally, Figure 3 displays the results from 50 sweeps, illustrating the relationship between  $\alpha$ ,  $\tau$ , and the regret.

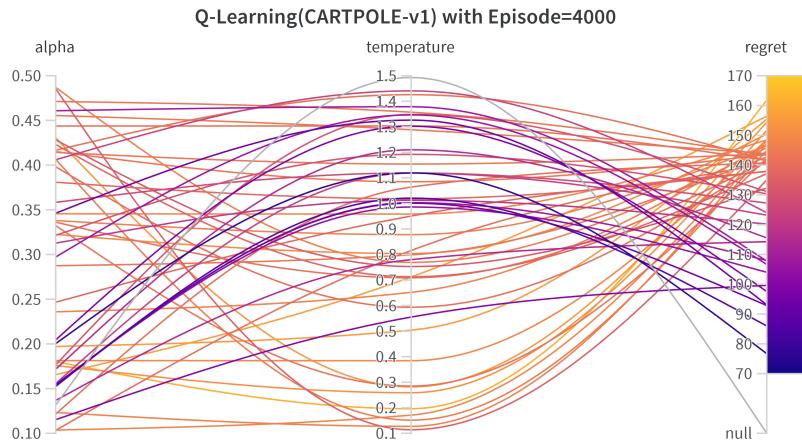


Figure 3: Q-Learning hyper-parameter sweeps for  $\alpha$  and  $\tau$  in Cartpole-v1(4000 episodes), with lines color-coded by regret.

#### 2.1.5 Q-Learning best 3 results

	$\alpha$	$\tau$	$\gamma$	regret
1	0.20034	1.11949	0.99	76.7096
2	0.15243	1.02074	0.9	86.03815
3	0.15404	1.0012	0.9	92.8918

Figure 4 include the episode vs reward plot for best three hyper-parameter combination.

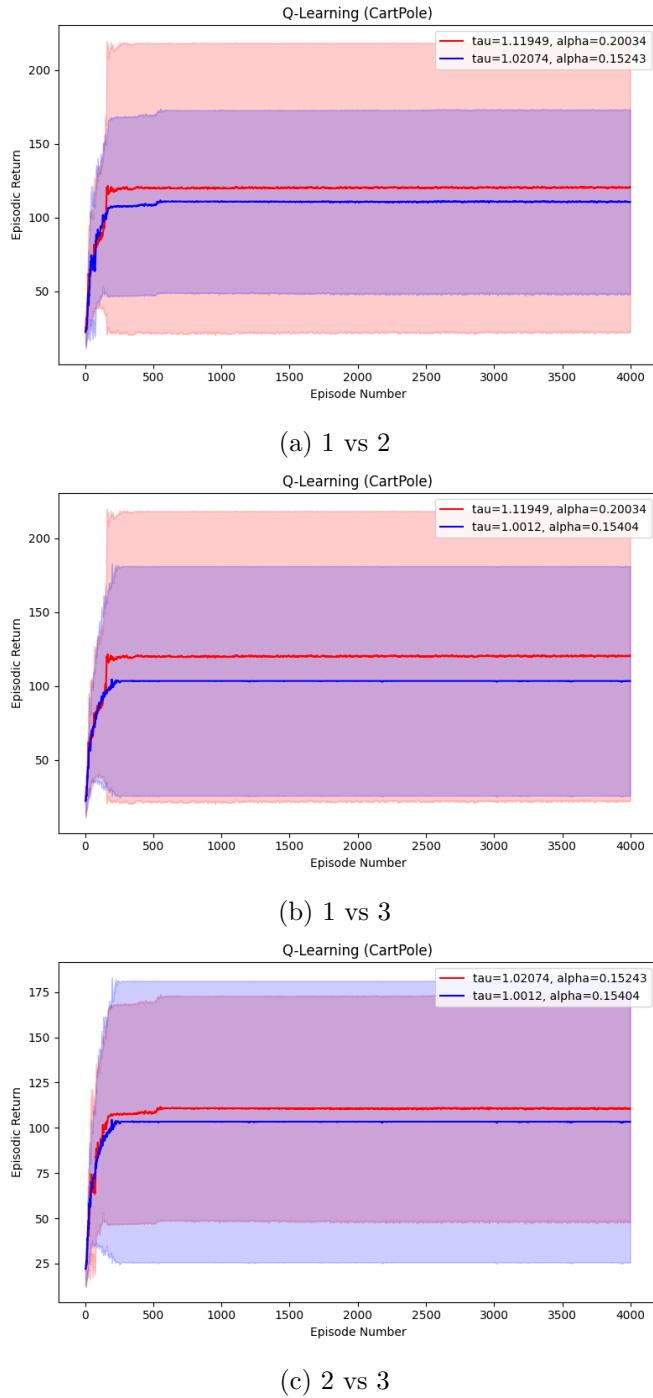


Figure 4: Performance comparison of the three best Q-Learning parameter configurations: (a) 1 vs 2, (b) 1 vs 3, and (c) 2 vs 3.

### 2.1.6 Result(SARSA vs Q-Learning)

The SARSA implementation achieved a best average reward of  $-255$  over 2000 episodes, whereas Q-learning achieved  $-259$  over the same number of episodes. Hence SARSA with  $\epsilon$ -greedy and Q-Learning both solves the environment and show a similar results.

## 2.2 MountainCar-v0

### 2.2.1 Code Snippets

```
1     current_reward = 0
```

Listing 7: Q-table

```
1 def getState(state, env_low=env_low, env_high=env_high, bins=bins):
2     """Returns the discretized position and velocity of an
3     observation"""
4     discretized_env = (env_high - env_low) / bins
5     discretized_pos = int((state[0] - env_low[0]) /
6                           discretized_env[0])
6     discretized_vel = int((state[1] - env_low[1]) /
7                           discretized_env[1])
8     discretized_pos = np.clip(discretized_pos, 0, bins - 1)
9     discretized_vel = np.clip(discretized_vel, 0, bins - 1)
10    return discretized_pos, discretized_vel
```

Listing 8: Discretized sates

```
1     """Choose action based on an epsilon greedy strategy"""
2     if random.random() < epsilon:
3         action = env.action_space.sample()
4     else:
5         action = np.argmax(q_table[pos][vel])
6     return action
```

Listing 9: epsilon-greedy action selection

```
1     print(f"\n==== Training with Seed: {seed} ===")
2     np.random.seed(seed)
3     random.seed(seed)
4     env.reset(seed=seed)
5
6     q_table_sarsa = np.zeros((bins + 1, bins + 1, env.
7                               action_space.n))
8
```

```

9     rewards_sarsa = []
10
11    for ep in range(episode):
12        current_reward = 0
13        done = False
14        truncated = False
15        state, _ = env.reset(seed=seed)
16        pos, vel = getState(state)
17        action = chooseAction(pos, vel, q_table_sarsa, epsilon)
18
19        # while not (done or truncated):
20        while not done:
21            next_state, reward, done, truncated, _ = env.step(
22                action)
23            next_pos, next_vel = getState(next_state)
24            next_action = chooseAction(next_pos, next_vel,
25                q_table_sarsa, epsilon)
26
27            if done:
28                q_table_sarsa[pos][vel][action] += alpha * (
29                    reward - q_table_sarsa[pos][vel][action])
30            else:
31                q_table_sarsa[pos][vel][action] += alpha * (
32                    reward + gamma * q_table_sarsa[next_pos][
33                    next_vel][next_action] - q_table_sarsa[pos][vel][action]
34                )
35
36            pos, vel = next_pos, next_vel
37            action = next_action
38            current_reward += reward
39            rewards_sarsa.append(current_reward)
40            print(f'seed {seed} Episode {ep+1}/{episode}, Reward: {current_reward}')
41            all_rewards.append(rewards_sarsa)
42
43    env.close()

```

Listing 10: SARSA implementation

```

1 def softmax(q_values, temperature):
2     """Returns a probability distribution over actions using
3     softmax"""
4     exp_q = np.exp(q_values / temperature)
5     return exp_q / np.sum(exp_q)
6
7 def chooseAction(pos, vel, q_table, temperature):
8     """Choose action based on Softmax exploration"""
9     q_values = q_table[pos][vel]
10    action_probs = softmax(q_values, temperature)

```

```

10     action = np.random.choice(np.arange(len(q_values)), p=
11                               action_probs)
12
13     return action

```

Listing 11: Softmax action selection

```

1     np.random.seed(seed)
2     random.seed(seed)
3     env.reset(seed=seed)
4
5     q_table_qlearn = np.zeros((bins + 1, bins + 1, env.
6                                action_space.n))
6     rewards_qlearn = []
7
8     for ep in range(episode):
9         current_reward = 0
10        done = False
11        truncated = False
12        state, _ = env.reset(seed=seed)
13        pos, vel = getState(state)
14
15        # while not (done or truncated):
16        while not done:
17            action = chooseAction(pos, vel, q_table_qlearn,
18                                   temperature)
19            next_state, reward, done, truncated, _ = env.step(
20                action)
21            next_pos, next_vel = getState(next_state)
22            # if done or truncated:
23            if done:
24                q_table_qlearn[pos][vel][action] += alpha * (
25                    reward - q_table_qlearn[pos][vel][action])
26            else:
27                max_next_q = np.max(q_table_qlearn[next_pos][
28                    next_vel])
29                q_table_qlearn[pos][vel][action] += alpha * (
30                    reward + gamma * max_next_q -
31                    q_table_qlearn[pos][vel][action]
32                )
33
34            pos, vel = next_pos, next_vel
35            current_reward += reward
36
37            rewards_qlearn.append(current_reward)
38            print(f'Seed {seed} Episode {ep+1}/{episode}, Reward: {current_reward}')
39
40    all_rewards.append(rewards_qlearn)

```

Listing 12: Q-Learning implementation

### 2.2.2 SARSA hyper-parameter tuning

Using Weights & Biases (wandb) with a sweep method based on Bayesian optimization, we identified the best-performing hyper-parameters. Specifically, we set  $\alpha \in [0.1, 0.5]$  and  $\epsilon \in [0.01, 0.15]$ , and ran 2000 episodes while minimizing the regret, defined as  $-200 -$  (all-time average return). See the wandb report on this environment here. Additionally, Figure 5 displays the results from 50 sweeps, illustrating the relationship between  $\alpha$ ,  $\epsilon$ , and the regret.

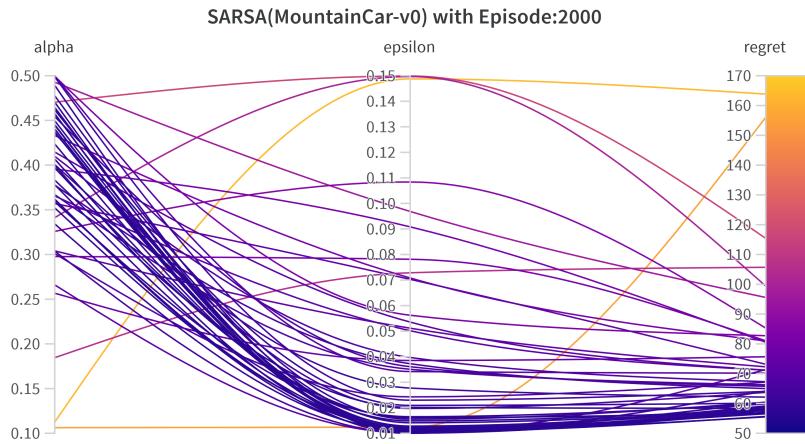
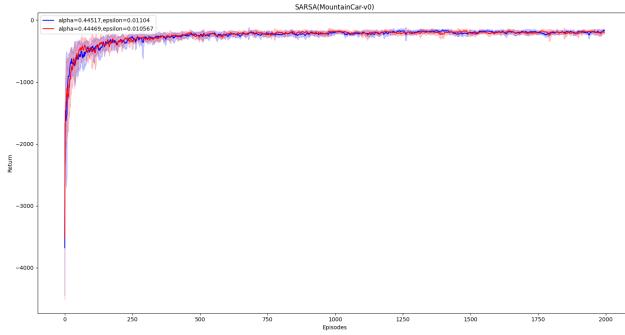


Figure 5: SARSA hyper-parameter sweeps for  $\alpha$  and  $\epsilon$  in MountainCar-v0(2000 episodes), with lines color-coded by regret.

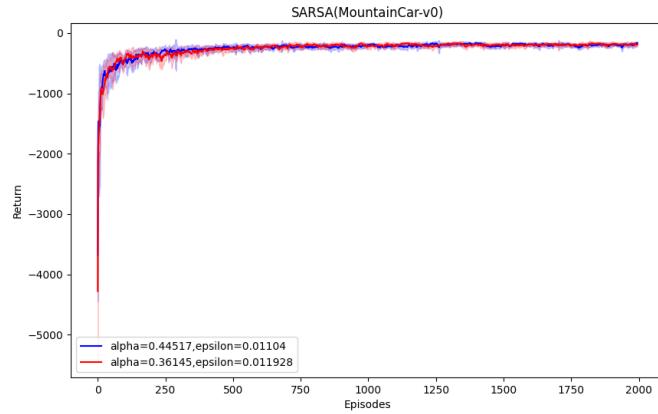
### 2.2.3 SARSA best 3 results

	$\alpha$	$\epsilon$	$\gamma$	regret
1	0.44517	0.01104	0.99	55.5099
2	0.44469	0.010567	0.99	55.6604
3	0.36145	0.011928	0.99	56.547

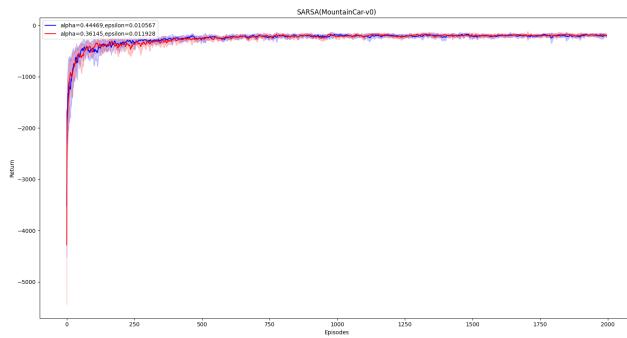
Figure 6 include the episode vs reward plot for best three hyper-parameter combination.



(a) 1 vs 2



(b) 1 vs 3



(c) 2 vs 3

Figure 6: Performance comparison of the three best SARSA parameter configurations: (a) 1 vs 2, (b) 1 vs 3, and (c) 2 vs 3.

#### 2.2.4 Q-Learning hyper-parameter tuning

Using Weights & Biases (wandb) with a sweep method based on Bayesian optimization, we identified the best-performing hyper-parameters. Specifically, we set  $\alpha \in [0.1, 0.5]$  and  $\tau \in [0.1, 1.5]$ , and ran 2000 episodes while minimizing the regret, defined as  $-200 -$  (all-time average return). See the wandb report on this environment here. Additionally, Figure 7 displays the results from 50 sweeps, illustrating the relationship between  $\alpha$ ,  $\tau$ , and the regret.

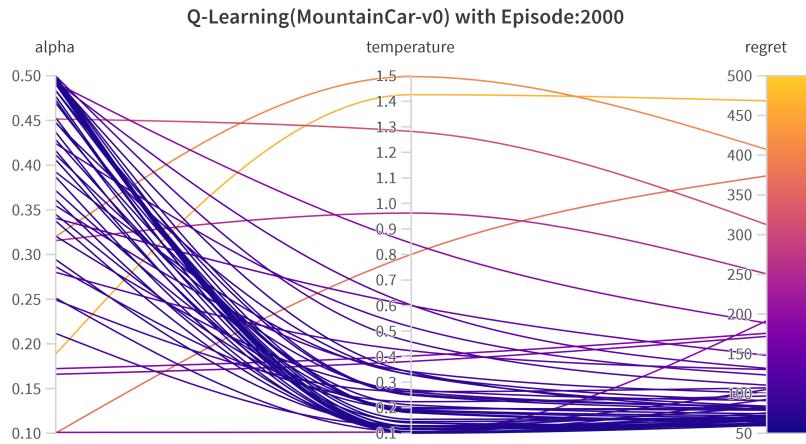
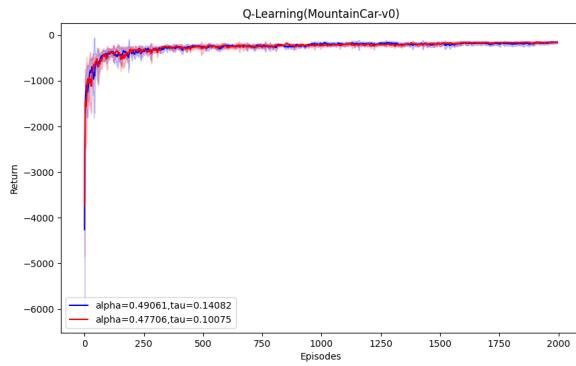


Figure 7: Q-Learning hyper-parameter sweeps for  $\alpha$  and  $\tau$  in MountainCar-v0(2000 episodes), with lines color-coded by regret.

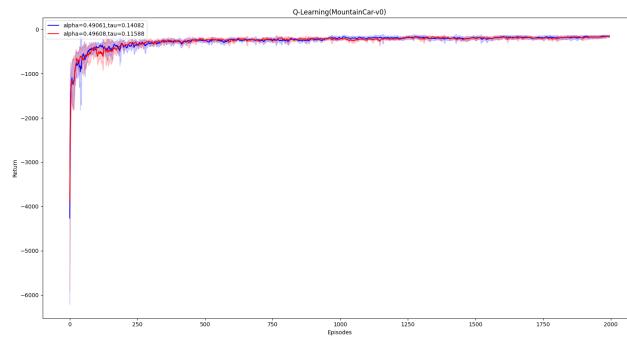
#### 2.2.5 Q-Learning best 3 results

	$\alpha$	$\tau$	$\gamma$	regret
1	0.49061	0.14082	0.99	59.6092
2	0.47706	0.10075	0.99	60.0469
3	0.49608	0.11588	0.99	60.4235

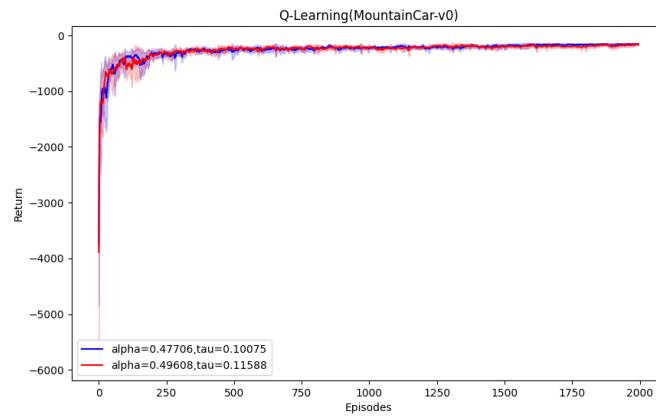
Figure 8 include the episode vs reward plot for best three hyper-parameter combination.



(a) 1 vs 2



(b) 1 vs 3



(c) 2 vs 3

Figure 8: Performance comparison of the three best Q-Learning parameter configurations: (a) 1 vs 2, (b) 1 vs 3, and (c) 2 vs 3.

## 2.2.6 Result(SARSA vs Q-Learning)

The SARSA implementation achieved a best average reward of  $-255$  over 2000 episodes, whereas Q-learning achieved  $-259$  over the same number of episodes. Hence SARSA with  $\epsilon$ -greedy and Q-Learning both solves the environment and show a similar results.

## 2.3 MiniGrid-Dynamic-Obstacles-5x5-v0

### 2.3.1 Code Snippets

Action selection using  $\epsilon$ -greedy

```
1 def take_action(q_value, epsilon):
2     if np.random.random() < epsilon:
3         return np.random.randint(0, 3)
4     return np.argmax(q_value)
```

Listing 13: epsilon-greedy action selection

```
1     q_value = np.zeros((3, 25, 4, 2))
2     total_reward = np.zeros(episodes)
3
4     for ep in range(episodes):
5         env.reset()
6         terminated, truncated = False, False
7         x1 = env.agent_pos[0] * 5 + env.agent_pos[1]
8         x2 = env.agent_dir
9         front_cell = env.grid.get(*env.front_pos)
10        x3 = 1 if (front_cell and front_cell.type != "goal")
11        ) else 0
12
13        action = take_action(q_value[:, x1, x2, x3],
14        epsilon)
15
16        while not (terminated or truncated):
17            observation, reward, terminated, truncated,
18            info = env.step(action)
19            new_x1 = env.agent_pos[0] * 5 + env.agent_pos
20            [1]
21            new_x2 = env.agent_dir
22            front_cell = env.grid.get(*env.front_pos)
23            new_x3 = 1 if (front_cell and front_cell.type
24            != "goal") else 0
25            new_action = take_action(q_value[:, new_x1,
26            new_x2, new_x3], epsilon)
27
28            q_value[action, x1, x2, x3] += alpha * (
```

```

23             reward + gamma * q_value[new_action, new_x1
24     , new_x2, new_x3] - q_value[action, x1, x2, x3]
25         )
26
27     x1, x2, x3, action = new_x1, new_x2, new_x3,
new_action
28     total_reward[ep] += reward
29
30     print(f"Seed: {seed} Episode: {ep+1} Reward: {total_reward[ep]}")
31
32     all_rewards.append(total_reward)
33     env.close()
34
35     all_rewards = np.array(all_rewards)
36     mean_rewards = np.mean(all_rewards, axis=0)
37     all_time_avg_reward = np.mean(mean_rewards)
38
39     return all_time_avg_reward, mean_rewards

```

Listing 14: SARSA implementation

```

1 def softmax_action(q_value, temperature):
2     exp_values = np.exp(q_value / temperature)
3     probabilities = exp_values / np.sum(exp_values)
4     return np.random.choice(len(q_value), p=probabilities)

```

Listing 15: Softmax action selection

```

1     env = gym.make('MiniGrid-Dynamic-Obstacles-Random-5x5-
v0')
2     env.reset(seed=seed)
3     # Q-table dimensions: actions x (5x5 grid) x agent
direction (4) x front cell flag (2)
4     q_value = np.zeros((3, 25, 4, 2))
5     total_reward = np.zeros(episodes)
6
7     for ep in range(episodes):
8         env.reset()
9         terminated, truncated = False, False
10        x1 = env.agent_pos[0] * 5 + env.agent_pos[1]
11        x2 = env.agent_dir
12        front_cell = env.grid.get(*env.front_pos)
13        x3 = 1 if (front_cell and front_cell.type != "goal"
) else 0
14
15        action = softmax_action(q_value[:, x1, x2, x3],
temperature)
16
17        while not (terminated or truncated):

```

```

18         observation, reward, terminated, truncated,
19     info = env.step(action)
20         new_x1 = env.agent_pos[0] * 5 + env.agent_pos
21             [1]
22             new_x2 = env.agent_dir
23             front_cell = env.grid.get(*env.front_pos)
24             new_x3 = 1 if (front_cell and front_cell.type
25 != "goal") else 0
26
27             best_next_action = np.argmax(q_value[:, new_x1,
28             new_x2, new_x3])
29             q_value[action, x1, x2, x3] += alpha * (
30                 reward + gamma * q_value[best_next_action,
31             new_x1, new_x2, new_x3] - q_value[action, x1, x2, x3]
32             )
33
34             x1, x2, x3 = new_x1, new_x2, new_x3
35             action = softmax_action(q_value[:, x1, x2, x3],
36             temperature)
37             total_reward[ep] += reward
38
39             print(f"Seed: {seed} Episode: {ep+1} Reward: {total_reward[ep]}")
40
41             all_rewards.append(total_reward)
42             env.close()
43
44             all_rewards = np.array(all_rewards)
45             mean_rewards = np.mean(all_rewards, axis=0)
46             all_time_avg_reward = np.mean(mean_rewards)
47             return all_time_avg_reward, mean_rewards

```

Listing 16: Q-Learning implementation

### 2.3.2 SARSA hyper-parameter tuning

Using Weights & Biases (wandb) with a sweep method based on Bayesian optimization, we identified the best-performing hyper-parameters. Specifically, we set  $\alpha \in [0.1, 0.5]$  and  $\epsilon \in [0.01, 0.15]$ , and ran 2000 episodes while minimizing the regret, defined as  $1 - (\text{all-time average return})$ . See the wandb report on this environment here. Additionally, Figure 9 displays the results from 50 sweeps, illustrating the relationship between  $\alpha$ ,  $\epsilon$ , and the regret.

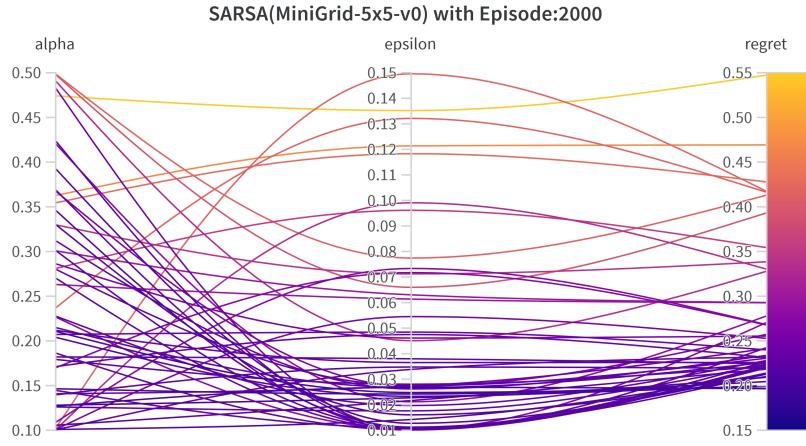
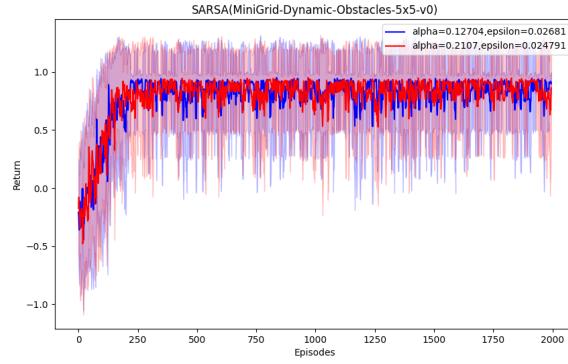


Figure 9: SARSA hyper-parameter sweeps for  $\alpha$  and  $\epsilon$  in MiniGrid-Dynamic-Obstacles-5x5-v0(2000 episodes), with lines color-coded by regret.

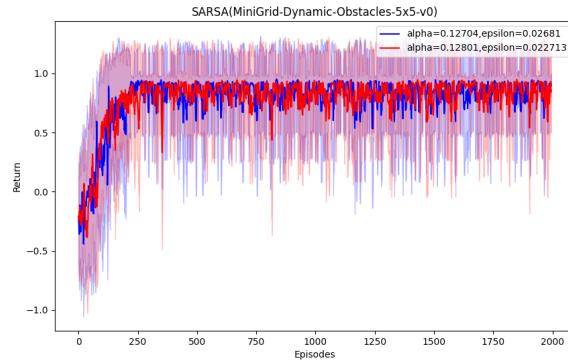
### 2.3.3 SARSA best 3 results

	$\alpha$	$\epsilon$	$\gamma$	regret
1	0.12704	0.02681	0.99	0.19653
2	0.2107	0.024791	0.99	0.19896
3	0.12801	0.022713	0.99	0.199997

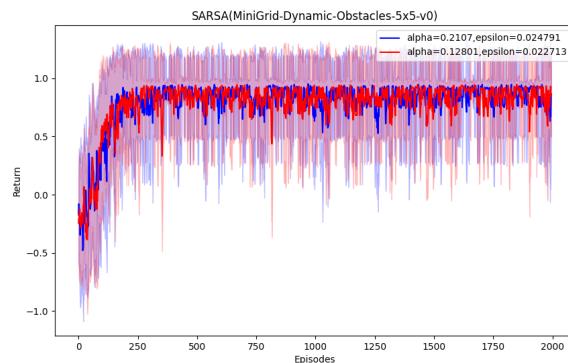
Figure 10 include the episode vs reward plot for best three hyper-parameter combination.



(a) 1 vs 2



(b) 1 vs 3



(c) 2 vs 3

Figure 10: Performance comparison of the three best SARSA parameter configurations: (a) 1 vs 2, (b) 1 vs 3, and (c) 2 vs 3.

### 2.3.4 Q-Learning hyper-parameter tuning

Using Weights & Biases (wandb) with a sweep method based on Bayesian optimization, we identified the best-performing hyper-parameters. Specifically, we set  $\alpha \in [0.1, 0.5]$  and  $\tau \in [0.1, 1.5]$ , and ran 2000 episodes while minimizing the regret, defined as  $1 - (\text{all-time average return})$ . See the wandb report on this environment here. Additionally, Figure 11 displays the results from 50 sweeps, illustrating the relationship between  $\alpha$ ,  $\tau$ , and the regret.

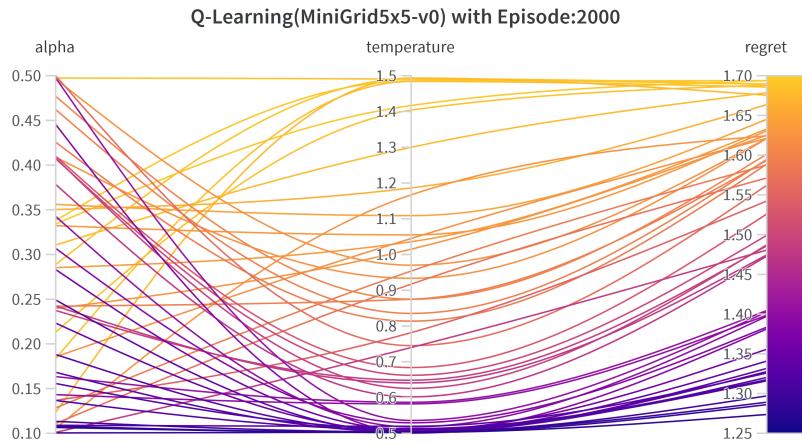
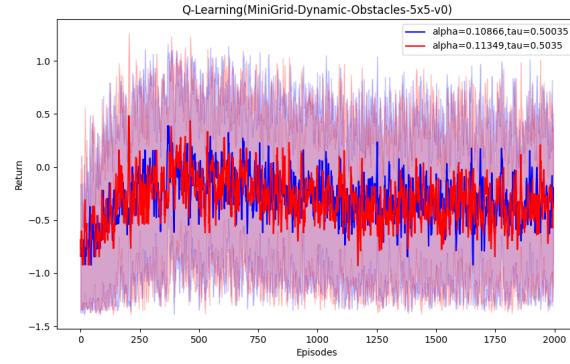


Figure 11: Q-Learning hyper-parameter sweeps for  $\alpha$  and  $\tau$  in MiniGrid-Dynamic-Obstacles-5x5-v0(2000 episodes), with lines color-coded by regret.

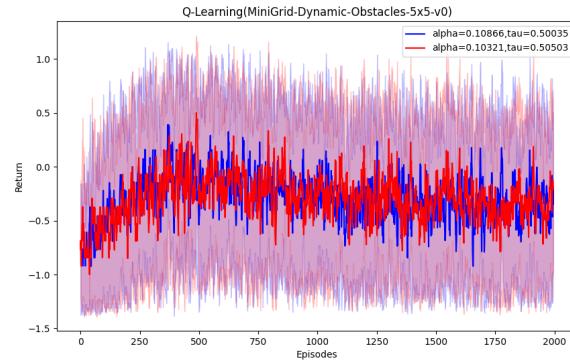
### 2.3.5 Q-Learning best 3 results

	$\alpha$	$\tau$	$\gamma$	regret
1	0.10866	0.50035	0.99	1.22505
2	0.11349	0.5035	0.99	1.25063
3	0.10321	0.50503	0.99	1.25763

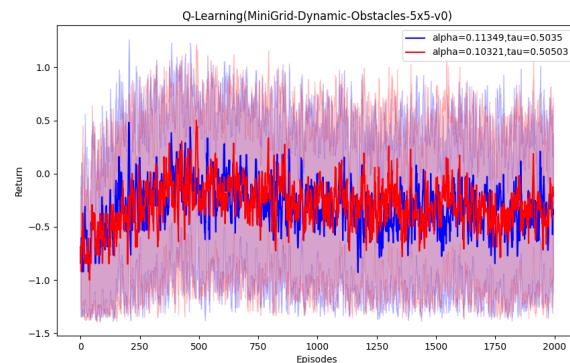
Figure 12 include the episode vs reward plot for best three hyper-parameter combination.



(a) 1 vs 2



(b) 1 vs 3



(c) 2 vs 3

Figure 12: Performance comparison of the three best Q-Learning parameter configurations: (a) 1 vs 2, (b) 1 vs 3, and (c) 2 vs 3.

### **2.3.6 Result(SARSA vs Q-Learning)**

The SARSA implementation achieved a best average reward of 0.8 over 2000 episodes, whereas Q-learning achieved  $-0.2$  over the same number of episodes. Hence SARSA with  $\epsilon$ -greedy solves the environment better than Q-Learning with softmax exploration strategy.

## **3 Github link**

[https://github.com/RitabrataMandal/RL-DA6400-assignment\\_1](https://github.com/RitabrataMandal/RL-DA6400-assignment_1)

## **4 Steps to recreate**

1. for CartPole-v1
  - cd to cartpole-v1
  - for sarsa implementation run sarsa.py
  - for Q-Learning run qlearning.py
2. for Mountain\_car-v0
  - cd mountain\_car-v0
  - for sarsa implementation run sarsa.py
  - for Q-Learning implementation run q\_learning.py
3. for MiniGrid-Dynamic-Obstacles-5x5-v0
  - cd minigrid\_world
  - for sarsa implementation run sarsa\_epsilon\_greedy.py
  - for Q-Learning implementation run q\_learning\_softmax.py

---

requirements.txt

---

```
gymnasium==0.29.1
numpy==1.26.4
pygame==2.5.2
matplotlib==3.8.4
argparse==1.1
wandb==0.19.8
```

---

```
1 method: bayes
2 metric:
3   name: regret
4   goal: minimize
5 parameters:
6   alpha:
7     min: 0.1
8     max: 0.5
9   epsilon:
10    min: 0.01
11    max: 0.15
```

Listing 17: SARSA.yaml

```
1 method: bayes
2 metric:
3   name: regret
4   goal: minimize
5 parameters:
6   temperature:
7     min: 0.5
8     max: 1.5
9   alpha:
10    min: 0.1
11    max: 0.5
```

Listing 18: Q-Learning.yaml