

DA6400 : Reinforcement Learning
Programming Assignment #1
Report

Ritabrata Mandal
EE24E009

March 30, 2025

Contents

1	Introduction	3
1.1	Environments	3
1.2	Algorithms	3
2	Implementation	4
2.1	CartPole-v1	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	MountainCar-v0	12
2.2.1	Code Snippets	12
2.2.2	SARSA hyper-parameter tuning	15
2.2.3	SARSA best 3 results	16
2.2.4	Q-Learning hyper-parameter tuning	18
2.2.5	Q-Learning best 3 results	18
2.2.6	Result(SARSA vs Q-Learning)	20
2.3	MiniGrid-Dynamic-Obstacles-5x5-v0	20
2.3.1	Code Snippets	20
2.3.2	SARSA hyper-parameter tuning	23
2.3.3	SARSA best 3 results	23
2.3.4	Q-Learning hyper-parameter tuning	25
2.3.5	Q-Learning best 3 results	25
2.3.6	Result(SARSA vs Q-Learning)	27
3	Github link	27
4	Requirements and Hyperparameter Configurations	27
4.1	Requirements	27
4.2	Hyperparameter Configurations	27
5	Steps to Recreate	28

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

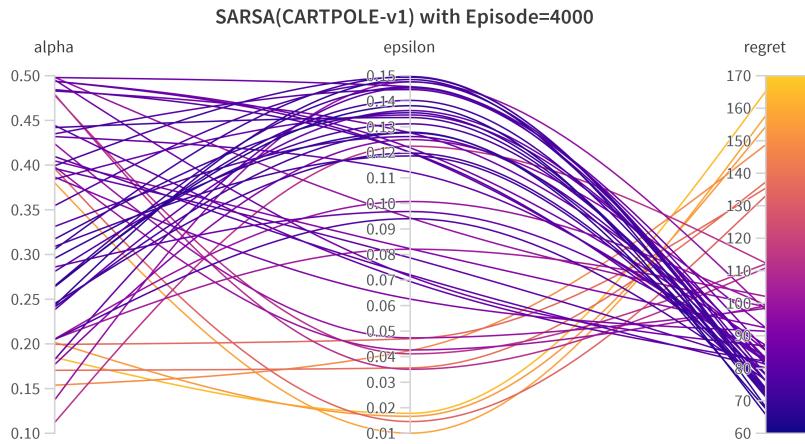
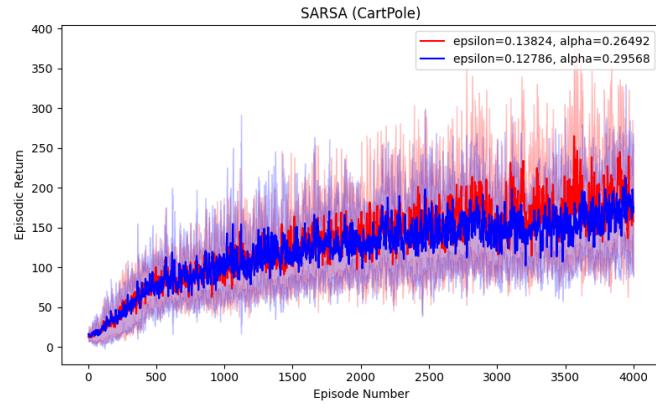


Figure 1: SARSA hyper-parameter sweeps for α and ϵ in CartPole-v1(4000 episodes), with lines color-coded by regret.

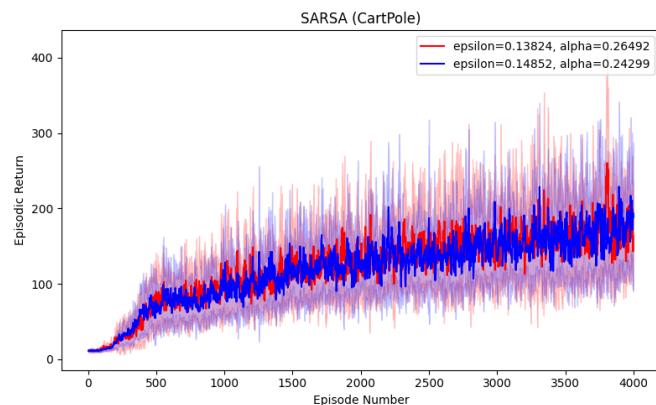
2.1.3 SARSA best 3 results

	α	ϵ	γ	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

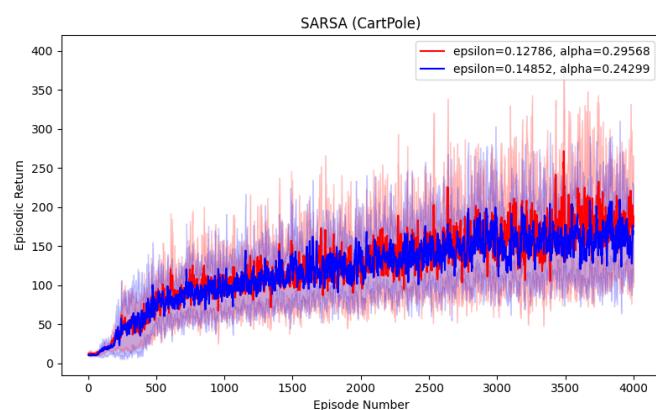
Figure 2 include the episode vs reward plot for best three hyper-parameter combination.



(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 α , τ , and the regret.

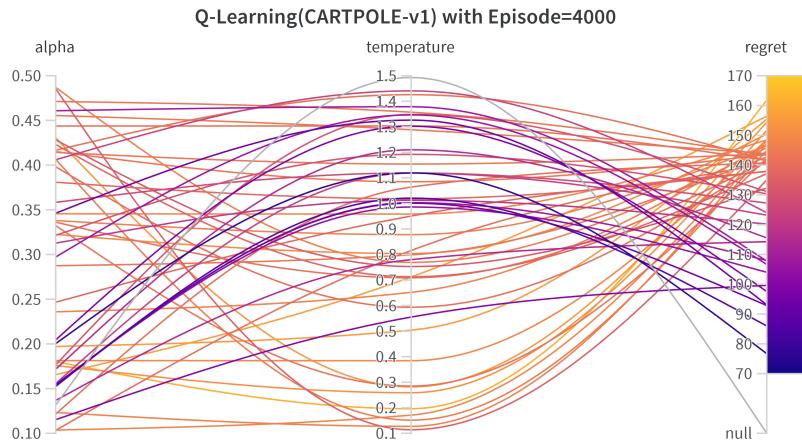


Figure 3: Q-Learning hyper-parameter sweeps for α and τ in Cartpole-v1(4000 episodes), with lines color-coded by regret.

2.1.5 Q-Learning best 3 results

	α	τ	γ	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.

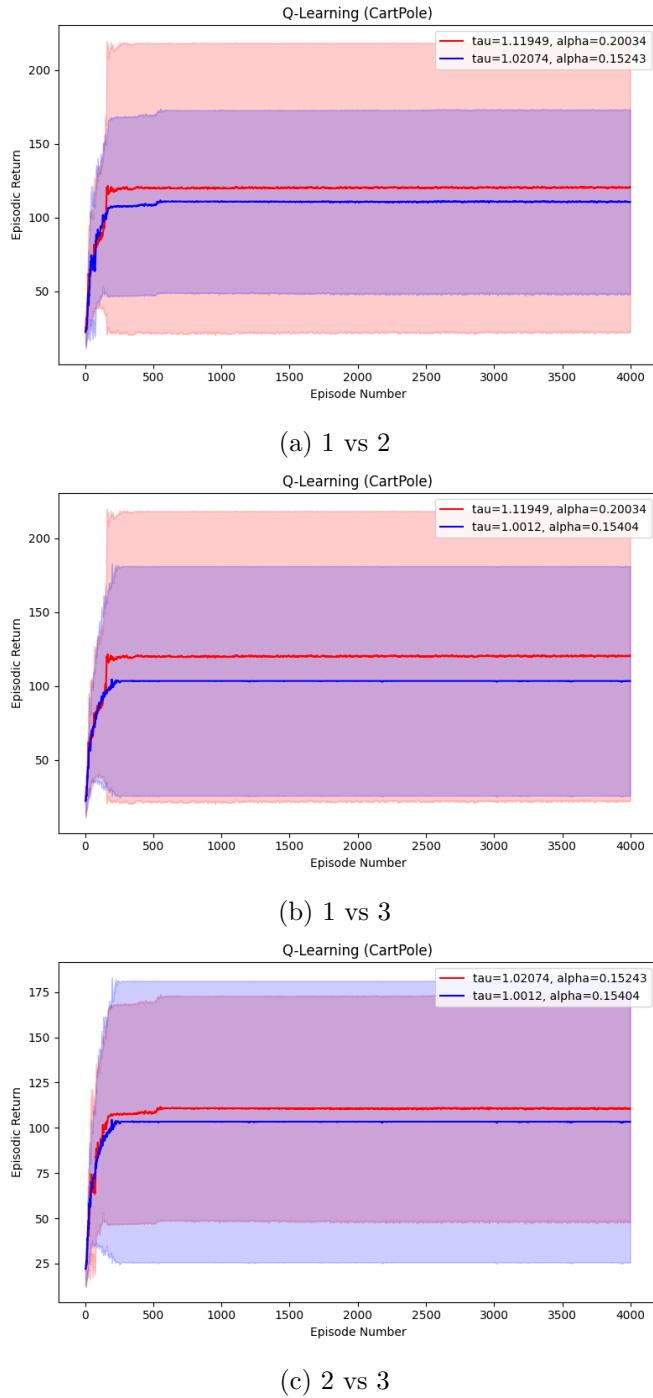


Figure 4: Performance comparison of the three best Q-Learning parameter configurations: (a) 1 vs 2, (b) 1 vs 3, and (c) 2 vs 3.

2.1.6 Result(SARSA vs Q-Learning)

The SARSA implementation achieved the highest average reward of 130 over 4000 episodes, while the Q-Learning implementation achieved an average reward of 119 over the same number of episodes.

From the graphs, we observe the following trends:

- SARSA: It takes longer to converge but ultimately achieves a higher average reward. The slower convergence indicates that SARSA explores more thoroughly, leading to better long-term performance.
- Q-Learning: Although Q-Learning converges faster, it tends to saturate at a lower reward, especially when different hyperparameters are used. This faster convergence comes at the cost of reduced average performance compared to SARSA.

In summary, SARSA demonstrates better stability and higher rewards in the long run, while Q-Learning converges more quickly but results in a slightly lower average reward.

2.2 MountainCar-v0

2.2.1 Code Snippets

```
1 q_table = np.zeros((bins + 1, bins + 1, env.  
action_space.n))
```

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  
observation"""  
3     discretized_env = (env_high - env_low) / bins  
4     discretized_pos = int((state[0] - env_low[0]) /  
discretized_env[0])  
5     discretized_vel = int((state[1] - env_low[1]) /  
discretized_env[1])  
6     discretized_pos = np.clip(discretized_pos, 0, bins - 1)  
7     discretized_vel = np.clip(discretized_vel, 0, bins - 1)  
8     return discretized_pos, discretized_vel
```

Listing 8: Discretized states

```

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

```

Listing 9: epsilon-greedy action selection

```

1 for seed in seeds:
2     print(f"\n==== Training with Seed: {seed} ===")
3     np.random.seed(seed)
4     random.seed(seed)
5     env.reset(seed=seed)
6
7
8     q_table_sarsa = np.zeros((bins + 1, bins + 1, env.
action_space.n))
9
10    rewards_sarsa = []
11
12    for ep in range(episode):
13        current_reward = 0
14        done = False
15        truncated = False
16        state, _ = env.reset(seed=seed)
17        pos, vel = getState(state)
18        action = chooseAction(pos, vel, q_table_sarsa, epsilon)
19
20        # while not (done or truncated):
21        while not done:
22            next_state, reward, done, truncated, _ = env.step(
action)
23            next_pos, next_vel = getState(next_state)
24            next_action = chooseAction(next_pos, next_vel,
q_table_sarsa, epsilon)
25
26            if done:
27                q_table_sarsa[pos][vel][action] += alpha * (
reward - q_table_sarsa[pos][vel][action])
28            else:
29                q_table_sarsa[pos][vel][action] += alpha * (
reward + gamma * q_table_sarsa[next_pos][
next_vel][next_action] - q_table_sarsa[pos][vel][action]
)
30
31
32            pos, vel = next_pos, next_vel
33            action = next_action
34

```

```

35         current_reward += reward
36     rewards_sarsa.append(current_reward)
37     print(f'seed {seed} Episode {ep+1}/{episode}, Reward: {current_reward}')
38     all_rewards.append(rewards_sarsa)

```

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

```

Listing 11: Softmax action selection

```

1 for seed in seeds:
2     np.random.seed(seed)
3     random.seed(seed)
4     env.reset(seed=seed)
5
6     q_table_qlearn = np.zeros((bins + 1, bins + 1, env.
7         action_space.n))
8     rewards_qlearn = []
9
10    for ep in range(episode):
11        current_reward = 0
12        done = False
13        truncated = False
14        state, _ = env.reset(seed=seed)
15        pos, vel = getState(state)
16
17        # while not (done or truncated):
18        while not done:
19            action = chooseAction(pos, vel, q_table_qlearn,
20                temperature)
21            next_state, reward, done, truncated, _ = env.step(
22                action)
23            next_pos, next_vel = getState(next_state)
24            # if done or truncated:
25            if done:

```

```

23             q_table_qlearn[pos][vel][action] += alpha * (
24                 reward - q_table_qlearn[pos][vel][action])
25             else:
26                 max_next_q = np.max(q_table_qlearn[next_pos][
27                     next_vel])
28             q_table_qlearn[pos][vel][action] += alpha * (
29                 reward + gamma * max_next_q -
30                 q_table_qlearn[pos][vel][action]
31             )
32
33             pos, vel = next_pos, next_vel
34             current_reward += reward
35
36             rewards_qlearn.append(current_reward)
37             print(f'Seed {seed} Episode {ep+1}/{episode}, Reward: {current_reward}')
38
39         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 α , ϵ , and the regret.

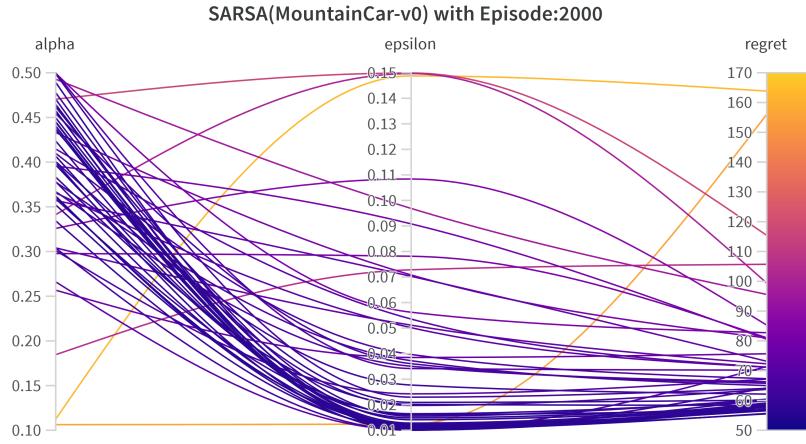
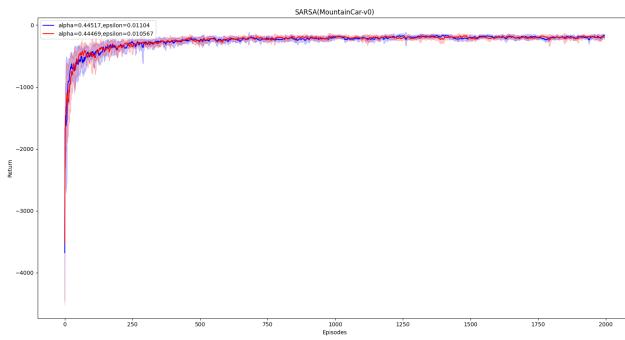


Figure 5: SARSA hyper-parameter sweeps for α and ϵ in MountainCar-v0(2000 episodes), with lines color-coded by regret.

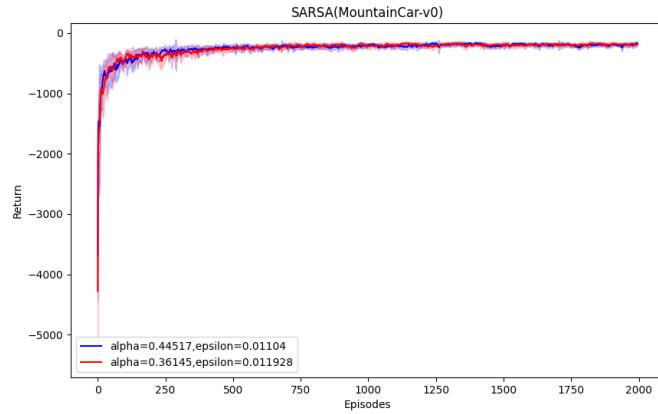
2.2.3 SARSA best 3 results

	α	ϵ	γ	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

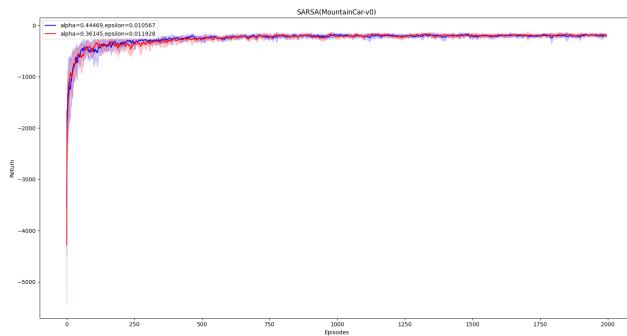
Figure 6 include the episode vs reward plot for best three hyper-parameter combination.



(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 α , τ , and the regret.

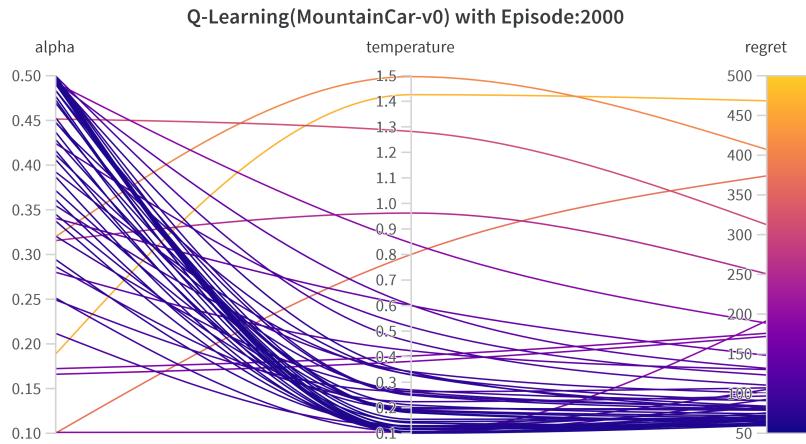
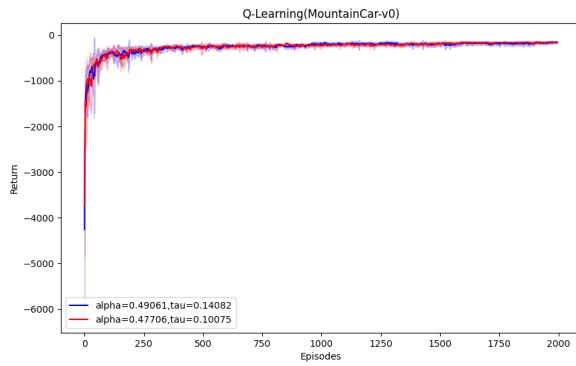


Figure 7: Q-Learning hyper-parameter sweeps for α and τ in MountainCar-v0(2000 episodes), with lines color-coded by regret.

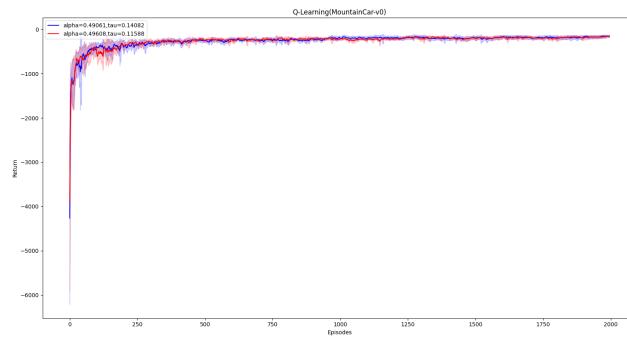
2.2.5 Q-Learning best 3 results

	α	τ	γ	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

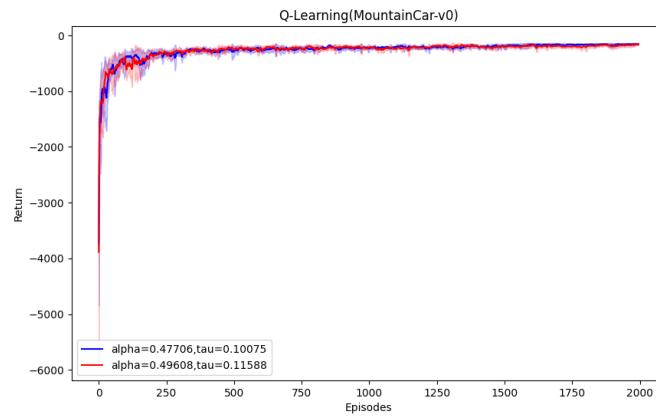
Figure 8 include the episode vs reward plot for best three hyper-parameter combination.



(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 the best average reward of -255 over 2000 episodes, while the Q-Learning implementation achieved an average reward of -259 over the same period.

From these results, we can conclude the following:

- SARSA with ϵ -greedy: It performs slightly better, achieving a marginally higher average reward.
- Q-Learning: Although it converges to a similar solution, its average reward remains slightly lower than that of SARSA.

In summary, both SARSA with ϵ -greedy and Q-Learning successfully solve the environment, producing comparable results with only a minor difference in performance.

2.3 MiniGrid-Dynamic-Obstacles-5x5-v0

2.3.1 Code Snippets

Action selection using ϵ -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     for seed in seeds:
2         env = gym.make('MiniGrid-Dynamic-Obstacles-Random-5x5-
v0')
3         env.reset(seed=seed)
4         # Q-table dimensions: action x (5x5 grid) x agent
direction (4) x front cell flag (2)
5         q_value = np.zeros((3, 25, 4, 2))
6         total_reward = np.zeros(episodes)
7
8         for ep in range(episodes):
9             env.reset()
10            terminated, truncated = False, False
11            x1 = env.agent_pos[0] * 5 + env.agent_pos[1]
12            x2 = env.agent_dir
13            front_cell = env.grid.get(*env.front_pos)
14            x3 = 1 if (front_cell and front_cell.type != "goal"
) else 0
```

```

15         action = take_action(q_value[:, x1, x2, x3],
16                           epsilon)
17
18         while not (terminated or truncated):
19             observation, reward, terminated, truncated,
20             info = env.step(action)
21             new_x1 = env.agent_pos[0] * 5 + env.agent_pos
22             [1]
23             new_x2 = env.agent_dir
24             front_cell = env.grid.get(*env.front_pos)
25             new_x3 = 1 if (front_cell and front_cell.type
26             != "goal") else 0
27             new_action = take_action(q_value[:, new_x1,
28             new_x2, new_x3], epsilon)
29
30             q_value[action, x1, x2, x3] += alpha * (
31                 reward + gamma * q_value[new_action, new_x1
32                 , new_x2, new_x3] - q_value[action, x1, x2, x3]
33                 )
34
35             x1, x2, x3, action = new_x1, new_x2, new_x3,
36             new_action
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 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     for seed in seeds:
2         env = gym.make('MiniGrid-Dynamic-Obstacles-Random-5x5-
3         v0')
4         env.reset(seed=seed)

```

```

4      # Q-table dimensions: actions x (5x5 grid) x agent
5      direction (4) x front cell flag (2)
6      q_value = np.zeros((3, 25, 4, 2))
7      total_reward = np.zeros(episodes)
8
9      for ep in range(episodes):
10         env.reset()
11         terminated, truncated = False, False
12         x1 = env.agent_pos[0] * 5 + env.agent_pos[1]
13         x2 = env.agent_dir
14         front_cell = env.grid.get(*env.front_pos)
15         x3 = 1 if (front_cell and front_cell.type != "goal"
16 ) else 0
17
18         action = softmax_action(q_value[:, x1, x2, x3],
19         temperature)
20
21         while not (terminated or truncated):
22             observation, reward, terminated, truncated,
23             info = env.step(action)
24             new_x1 = env.agent_pos[0] * 5 + env.agent_pos
25             [1]
26             new_x2 = env.agent_dir
27             front_cell = env.grid.get(*env.front_pos)
28             new_x3 = 1 if (front_cell and front_cell.type
29             != "goal") else 0
30
31             best_next_action = np.argmax(q_value[:, new_x1,
32             new_x2, new_x3])
33             q_value[action, x1, x2, x3] += alpha * (
34                 reward + gamma * q_value[best_next_action,
35             new_x1, new_x2, new_x3] - q_value[action, x1, x2, x3]
36             )
37
38             x1, x2, x3 = new_x1, new_x2, new_x3
39             action = softmax_action(q_value[:, x1, x2, x3],
40             temperature)
41             total_reward[ep] += reward
42
43             print(f"Seed: {seed} Episode: {ep+1} Reward: {total_reward[ep]}")
44
45             all_rewards.append(total_reward)
46             env.close()
47
48             all_rewards = np.array(all_rewards)
49             mean_rewards = np.mean(all_rewards, axis=0)
50             all_time_avg_reward = np.mean(mean_rewards)

```

```
42     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 α , ϵ , and the regret.

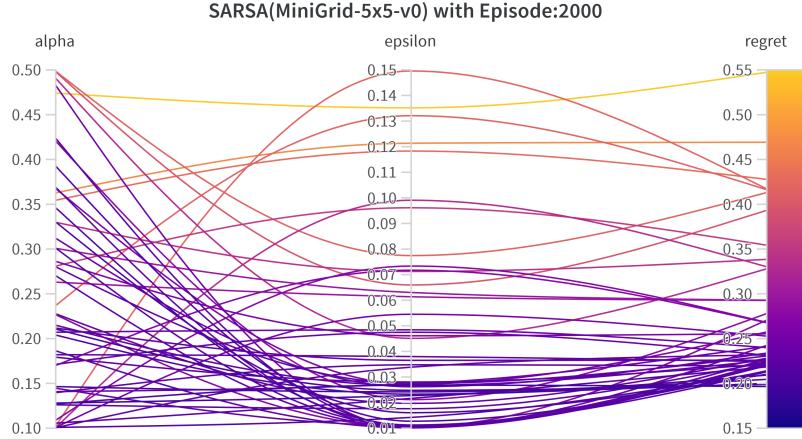
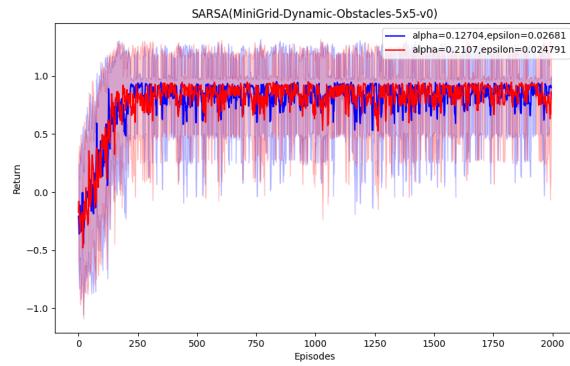


Figure 9: SARSA hyper-parameter sweeps for α and ϵ in MiniGrid-Dynamic-Obstacles-5x5-v0(2000 episodes), with lines color-coded by regret.

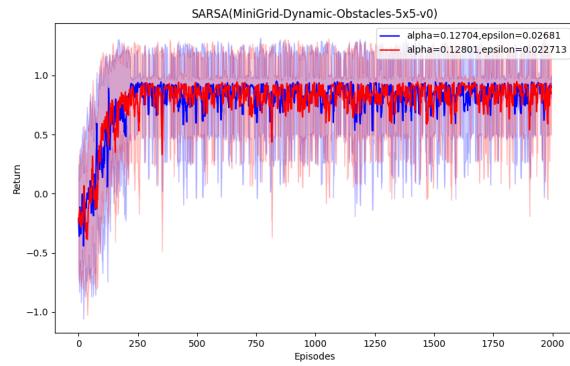
2.3.3 SARSA best 3 results

	α	ϵ	γ	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

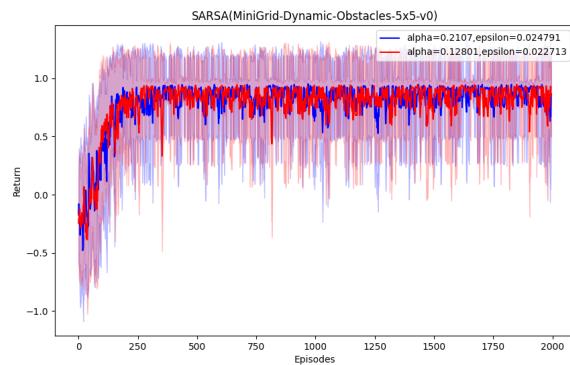
Figure 10 include the episode vs reward plot for best three hyper-parameter combination.



(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 α , τ , and the regret.

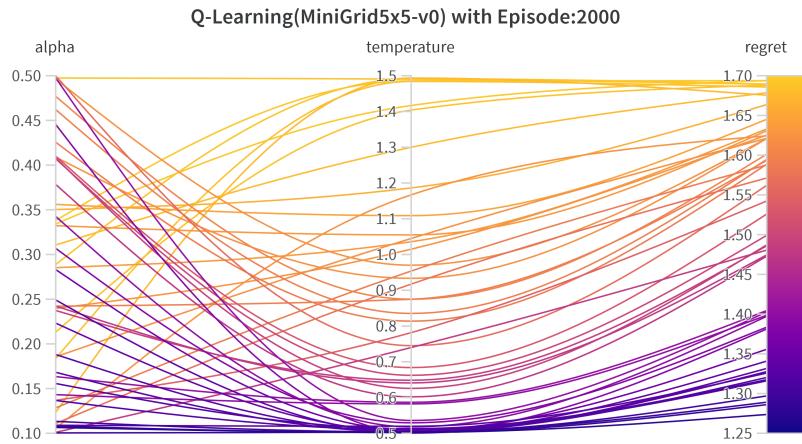
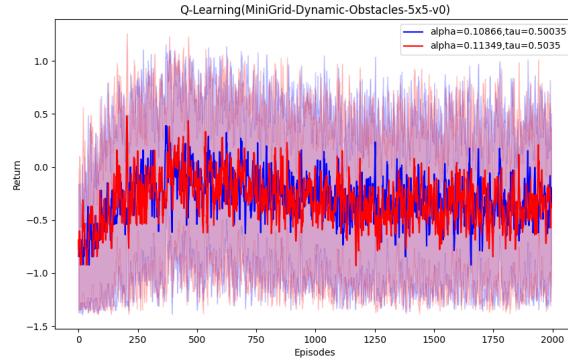


Figure 11: Q-Learning hyper-parameter sweeps for α and τ in MiniGrid-Dynamic-Obstacles-5x5-v0(2000 episodes), with lines color-coded by regret.

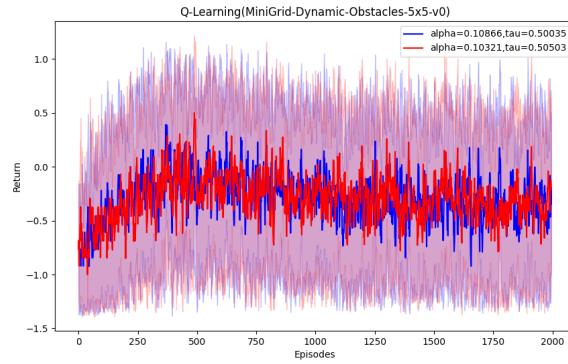
2.3.5 Q-Learning best 3 results

	α	τ	γ	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

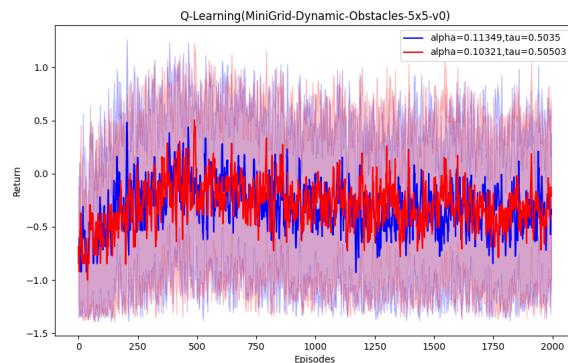
Figure 12 include the episode vs reward plot for best three hyper-parameter combination.



(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 ϵ -greedy solves the environment better than Q-Learning with softmax exploration strategy.

3 Github link

https://github.com/RitabrataMandal/RL-DA6400-assignment_1

4 Requirements and Hyperparameter Configurations

4.1 Requirements

The required packages and dependencies are listed in the `requirements.txt` file:

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
```

4.2 Hyperparameter Configurations

The hyperparameter configurations for SARSA and Q-Learning are provided in the following YAML files:

- **SARSA Configuration:** The configuration for SARSA is defined in the file `sweep_sarsa.yaml`.

```
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: `sweep_sarsa.yaml`

- **Q-Learning Configuration:** The configuration for Q-Learning is defined in the file `sweep_qlearning.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: `sweep_qlearning.yaml`

5 Steps to Recreate

To set up and run the project, follow these steps:

1. Install Required Packages:

- Run the following command to install all necessary packages:

```
pip install -r requirements.txt
```

2. Run Experiments for Different Environments:

(a) CartPole-v1:

- Navigate to the `cartpole-v1` directory:

```
cd cartpole-v1
```

- Run the appropriate script:

- For SARSA implementation:

```
python sarsa.py
```

- For Q-Learning implementation:

```
python q_learning.py
```

(b) **Mountain_Car-v0:**

- Navigate to the `mountain_car-v0` directory:

```
cd mountain_car-v0
```

- Run the appropriate script:

- For SARSA implementation:

```
python sarsa.py
```

- For Q-Learning implementation:

```
python q_learning.py
```

(c) **MiniGrid-Dynamic-Obstacles-5x5-v0:**

- Navigate to the `minigrid_world` directory:

```
cd minigrid_world
```

- Run the appropriate script:

- For SARSA implementation:

```
python sarsa_epsilon_greedy.py
```

- For Q-Learning implementation:

```
python q_learning_softmax.py
```