# DA6400: Reinforcement Learning Programming Assignment 1 Report

Kush Shah CH21B053        Sai Teja MS ME21B171

March 30, 2025

**Abstract**

This report describes the implementation and experimental evaluation of SARSA (with $\epsilon$-greedy exploration) and Q-Learning (with Softmax exploration) on the Gymnasium environments: CartPole-v1 and MountainCar-v0. We compare each algorithm with its own variant over 5 random seeds and report the mean and variance of the episodic returns. Hyperparameter tuning results and inferences are also discussed.

## 1 Introduction

In this assignment, we implemented two tabular RL algorithms: SARSA and Q-Learning. The primary goal was to train these algorithms on various environments (CartPole-v1, MountainCar-v0, and optionally MiniGrid-Dynamic-Obstacles-5x5-v0) and evaluate their performance. We adhered to the following guidelines:

- Use $\gamma = 0.99$ for all experiments.

- Tune hyperparameters to minimize regret.

- Use 5 random seeds per experiment and plot the mean and variance of episodic returns.

## 2 Environment Descriptions

### 2.1 CartPole-v1

CartPole-v1 consists of a cart moving along a frictionless track with a pole attached by an un-actuated joint. The goal is to balance the pole by applying forces left or right.

### 2.2 MountainCar-v0

MountainCar-v0 involves a car positioned in a valley with a sinusoidal track. The objective is to drive the car to the top of the right hill by carefully accelerating left or right.

### 2.3 MiniGrid-Dynamic-Obstacles-5x5-v0 (Bonus)

This environment is an empty room with moving obstacles. The agent must reach a goal square without colliding with obstacles.

## 3 Algorithm Implementations

We implemented SARSA and Q-Learning in separate Python modules. Below are snippets of the key parts of our code.

## 3.1 SARSA Implementation

```
1  # Update rule for SARSA
2  Q[tuple(state)+(action,)] += alpha * (reward + gamma * Q[tuple(next_state)+(
       next_action,)] - Q[tuple(state)+(action,)])
```

Listing 1: Key SARSA Update Code

## 3.2 Our SARSA Code Implementation

```
1  def sarsa(env, Q, episodes, gamma=0.99, alpha=0.1, epsilon=1.0, epsilon_decay
       =0.995, epsilon_min=0.05, print_freq=100, plot_heat=False, choose_action=
       choose_action_epsilon):
2
3      episode_rewards = np.zeros(episodes)
4      steps_to_completion = np.zeros(episodes)
5
6
7      for ep in tqdm(range(episodes)):
8          tot_reward = 0
9          steps = 0
10
11          # Reset the environment and get the initial discretized state.
12          state = env.reset()
13
14          action = choose_action(Q, state, epsilon)
15
16          done = False
17          while not done:
18              # Take a step in the environment.
19              state_next, reward, done = env.step(action)
20
21              # Select next action using the same policy.
22              next_action = choose_action(Q, state_next, epsilon)
23
24              # Sarsa update
25              Q[tuple(state) + (action,)] += alpha * (reward + gamma * Q[tuple(
                   state_next) + (next_action,)] - Q[tuple(state) + (action,)])
26
27              tot_reward += reward
28              steps += 1
29
30              state, action = state_next, next_action
31
32          episode_rewards[ep] = tot_reward
33          steps_to_completion[ep] = steps
34
35          if (ep + 1) % print_freq == 0:
36              avg_reward = np.mean(episode_rewards[ep-print_freq+1:ep+1])
37              avg_steps = np.mean(steps_to_completion[ep-print_freq+1:ep+1])
38              print(f"Episode {ep+1}: Avg Reward: {avg_reward:.2f}, Avg Steps: {
                   avg_steps:.2f}, Qmax: {Q.max():.2f}, Qmin: {Q.min():.2f}")
39
40          # Decay epsilon after each episode
41          epsilon = max(epsilon * epsilon_decay, epsilon_min)
42
43      return episode_rewards, steps_to_completion
```

Listing 2: Our SARSA Code Implementation

## 3.3 Q-Learning Implementation

```
1 # Update rule for Q-Learning
2 best_next_value = np.max(Q[tuple(next_state)])
3 Q[tuple(state)+(action,)] += alpha * (reward + gamma * best_next_value - Q[
    tuple(state)+(action,)])
```

Listing 3: Key Q-Learning Update Code

## 3.4 Our Q-Learning Code Implementation

```
1 def q_learning(env, Q, episodes, gamma=0.99, alpha=0.1,
2                tau=1.0, tau_decay=0.995, tau_min=0.1,
3                print_freq=100, plot_heat=False, choose_action=
                    choose_action_softmax):
4
5     episode_rewards = np.zeros(episodes)
6     steps_to_completion = np.zeros(episodes)
7
8     for ep in tqdm(range(episodes), desc="Q-learning"):
9         tot_reward = 0
10        steps = 0
11
12        # Reset the environment and get the initial discretized state.
13        state = env.reset()
14        action = choose_action(Q, state, tau)
15        done = False
16
17        while not done:
18            # Take a step in the environment.
19            state_next, reward, done = env.step(action)
20
21            # Q-learning update
22            best_next_value = np.max(Q[tuple(state_next)])
23            Q[tuple(state) + (action,)] += alpha * (reward + gamma *
                best_next_value - Q[tuple(state) + (action,)])
24
25            tot_reward += reward
26            steps += 1
27
28            # Update state and select next action.
29            state = state_next
30            action = choose_action(Q, state, tau)
31
32        episode_rewards[ep] = tot_reward
33        steps_to_completion[ep] = steps
34
35        if (ep + 1) % print_freq == 0:
36            avg_reward = np.mean(episode_rewards[ep - print_freq + 1 : ep + 1])
37            avg_steps = np.mean(steps_to_completion[ep - print_freq + 1 : ep +
                1])
38            print(f"Episode {ep+1}: Avg Reward: {avg_reward:.2f}, Avg Steps: {
                avg_steps:.2f}, Qmax: {Q.max():.2f}, Qmin: {Q.min():.2f}")
39
40        tau = max(tau * tau_decay, tau_min)
41
42    return episode_rewards, steps_to_completion
```

Listing 4: Our Q-Learning Implementation

## 3.5 Policy Functions

We used two policy functions:

- `choose_action_epsilon`: for $\epsilon$-greedy exploration.

- `choose_action_softmax`: for softmax (Boltzmann) exploration.

```python
def choose_action_epsilon(Q, state, epsilon, rg=rg):
    q_values = Q[tuple(state)]
    # If all Q-values are zero or with probability epsilon, choose a random
        action
    if not q_values.any() or rg.rand()<epsilon:
        return rg.choice(Q.shape[-1])
    else:
        return np.argmax(q_values)

def choose_action_softmax(Q, state, tau, rg=rg):
    q_values = Q[tuple(state)]
    probs = softmax(q_values - np.max(q_values)) ** (1/tau)
    probs = probs / np.sum(probs)
    return rg.choice(Q.shape[-1], p=probs)
```

Listing 5: Our policy Function Code

# 4 Experimental Setup

We discretized the continuous state space using a custom function `create_bins()` and `discretize_state()` (see `utils.py`). The Q-table is built based on the number of discretization bins per state dimension.

## 4.1 State Discretization

To apply tabular reinforcement learning methods, we discretized the continuous state space using the following helper functions:

```python
def create_bins(env, n_bins):
    """
    Create bins for each state variable. This is required for discretizing the
        state for sarsa and q-learning.
    """
    if env.spec.id == "CartPole-v1":
        bins = []
        bins.append(np.linspace(-4.8, 4.8, n_bins))     # Cart position
        bins.append(np.linspace(-3.0, 3.0, n_bins))     # Cart velocity
        bins.append(np.linspace(-0.418, 0.418, n_bins))# Pole angle
        bins.append(np.linspace(-3.5, 3.5, n_bins))     # Pole angular velocity
        return bins
    else:
        obs_low = env.observation_space.low
        obs_high = env.observation_space.high
        bins = []
        for i in range(len(obs_low)):
            bins.append(np.linspace(obs_low[i], obs_high[i], n_bins))
        return bins

def discretize_state(state, bins):
    """
    Discretize the state.
    """
```

```
24      indices=[]
25      for i, val in enumerate(state):
26          # Subtract 1 as digitize is 1-indexed
27          index = np.digitize(val, bins[i]) - 1
28          # Ensure the index is within the bounds
29          index = min(max(index, 0), len(bins[i]) - 2)
30          indices.append(index)
31      return tuple(indices)
```

<div align="center">Listing 6: State Discretization Functions</div>

The function `create_bins()` creates bins for each state variable, allowing us to convert continuous state values into discrete indices. The function `discretize_state()` maps an observed state to its corresponding discretized index using the bins. This discretization enables us to use a Q-table for SARSA and Q-learning.

Experiments were run over 5 random seeds. For each experiment, we recorded the episodic return and computed the mean and standard deviation across seeds.

# 5 Results and Plots

The main() function code is `run_experiment.py` is:

```
1  def main():
2      parser = argparse.ArgumentParser(description="Run RL algorithms on a Gym
           environment.")
3      parser.add_argument("--env", type=str, default="MountainCar-v0",
4                          help="Gym environment to run (e.g., MountainCar-v0,
                              CartPole-v1)")
5      parser.add_argument("--algorithm", type=str, default="sarsa",
6                          choices=["sarsa", "qlearning", "both"],
7                          help="Algorithm to run: sarsa, qlearning, or both")
8      parser.add_argument("--episodes", type=int, default=12000, help="Number of
           episodes to run")
9      parser.add_argument("--n_bins", type=int, default=20, help="Number of bins
           per state dimension")
10     parser.add_argument("--tau_decay", type=float, default=0.999, help="
           Tau_decay")
11     parser.add_argument("--epsilon_decay", type=float, default=0.995, help="
           Epsilon_decay")
12     parser.add_argument("--alpha", type=float, default=0.1, help="Learning rate
           ")
13     args = parser.parse_args()
14
15     seeds = [42, 43, 44, 45, 46]
16     gamma = 0.99
17
18     epsilon = 1.0
19     epsilon_min = 0.05
20
21     # For softmax exploration (Q-learning), use tau (temperature)
22     tau = 2.0
23     tau_min = 0.1
24
25     if args.algorithm in ["sarsa", "both"]:
26         print("Running SARSA with epsilon-greedy exploration on", args.env)
27         sarsa_rewards = run_experiment(args.env, sarsa, choose_action_epsilon,
               args.episodes, seeds,
28                                  gamma, args.alpha, epsilon, epsilon_decay=args.
                                      epsilon_decay,
29                                  epsilon_min=epsilon_min, n_bins=args.n_bins)
```

```
30          plot_results(sarsa_rewards, f"SARSA on {args.env}")
31
32      if args.algorithm in ["qlearning", "both"]:
33          print("Running Q-Learning with softmax exploration on", args.env)
34          q_rewards = run_experiment(args.env, q_learning, choose_action_softmax,
                args.episodes, seeds,
35                              gamma, args.alpha, epsilon, tau=tau, tau_decay=args.
                                  tau_decay, tau_min=tau_min,
36                              n_bins=args.n_bins)
37          plot_results(q_rewards, f"Q-Learning on {args.env}")
```

The **run_experiment** function is implemented as follows:

```
1  def run_experiment(env_name, algorithm, policy_fn, episodes, seeds, gamma,
       alpha,
2                      epsilon, epsilon_decay=None, epsilon_min=None,
3                      tau=None, tau_decay=None, tau_min=None, n_bins=20):
4      all_rewards = np.zeros((len(seeds), episodes))
5      print_freq = 100
6
7      for i, seed in enumerate(seeds):
8          env = gym.make(env_name)
9          env.reset(seed=seed)
10
11         bins = create_bins(env, n_bins)
12         # Discretize the environment
13         d_env = DiscretizedEnv(env, bins)
14
15         n_actions = env.action_space.n
16         # Initialize Q-table
17         state_shape = (n_bins-1,)*len(env.observation_space.low)
18         Q = np.zeros(state_shape + (n_actions,))
19
20         if policy_fn == choose_action_epsilon:
21             rewards, _ = algorithm(d_env, Q, episodes, gamma=gamma, alpha=alpha
                   , epsilon=epsilon,
22                                     epsilon_decay=epsilon_decay, epsilon_min=
                                         epsilon_min,
23                                     print_freq=print_freq, choose_action=
                                         policy_fn)
24         else:
25             rewards, _ = algorithm(d_env, Q, episodes, gamma=gamma, alpha=alpha
                   , tau=tau,
26                                     tau_decay=tau_decay, tau_min=tau_min,
27                                     print_freq=print_freq, choose_action=
                                         policy_fn)
28         all_rewards[i, :] = rewards
29
30     return all_rewards
```

## 5.1  Episodic Returns

**The results are plotted for the best hyperparameters. The best hyperparameters are found using hyperparameter tuning, which is in the next section. The method to reproduce these results is in the readme.md**

Figure 1 shows the episodic return (mean ± variance) for SARSA on CartPole-v1.
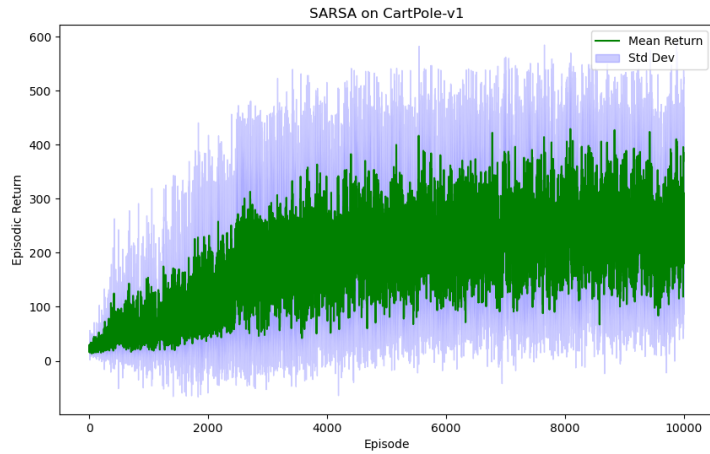
Figure 1: SARSA on CartPole-v1: Mean and standard deviation of episodic return over 5 seeds.

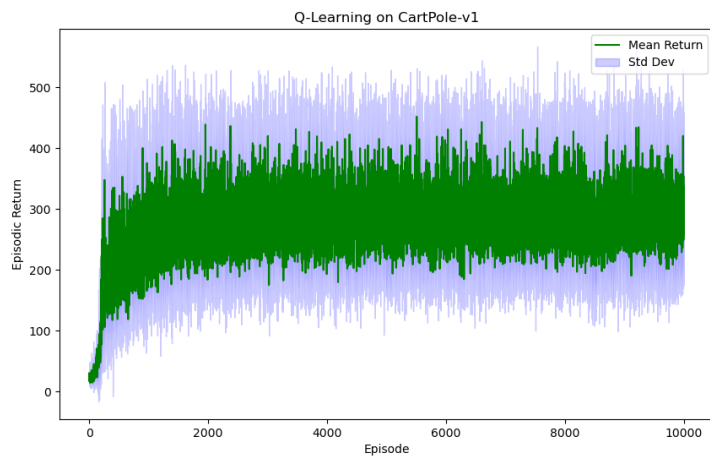Figure 2 shows the corresponding plot for Q-Learning on Cartpole-v1.



Figure 2: Q-Learning on CartPole-v1: Mean and standard deviation of episodic return over 5 seeds.

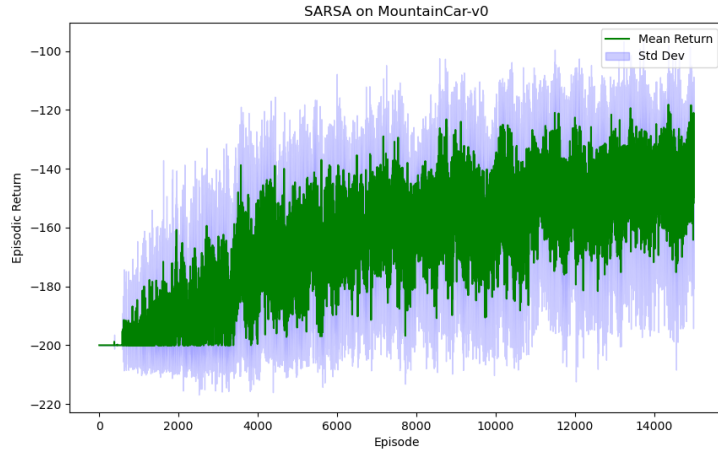Figure 3 shows the episodic return (mean $\pm$ variance) for SARSA on MountainCar-v0.

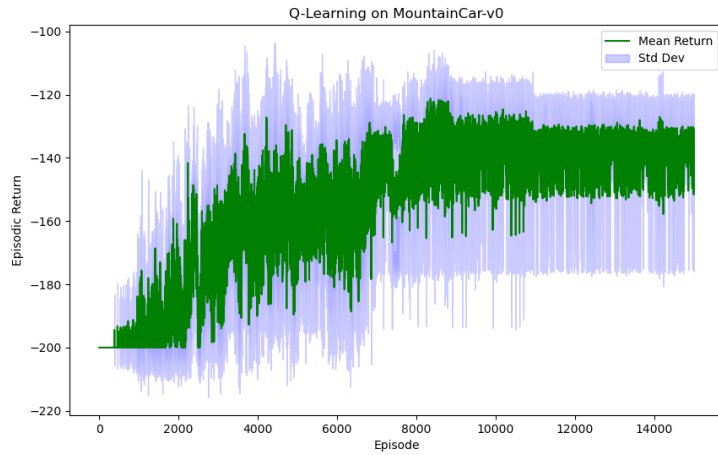Figure 3: SARSA on MountainCar-v0: Mean and standard deviation of episodic return over 5 seeds.

Figure 4 shows the episodic return (mean ± variance) for Q-Learning on MountainCar-v0.



Figure 4: Q-Learning on MountainCar-v0: Mean and standard deviation of episodic return over 5 seeds.

# 6 Hyperparameter Tuning

We tested various hyperparameter settings. The top 3 settings for each algorithm were:

## 6.1 Parallel Coordinates Plots and Reward Images

The hyperparameters taken for SARSA on CartPole-v1 are:

- **Learning rate**
- **Epsilon Decay**
- **Number of bins**

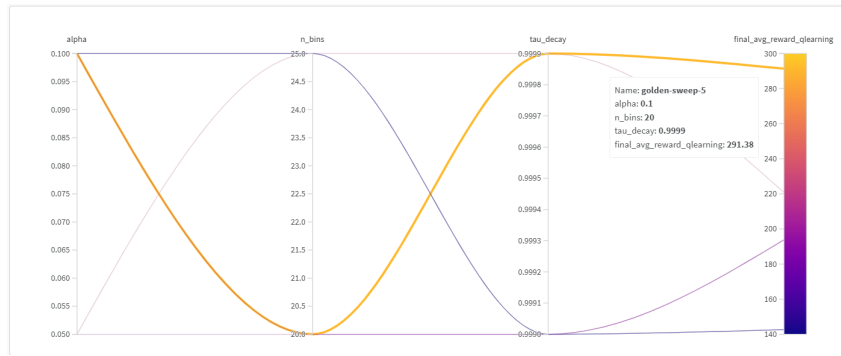Figure 5 shows the final parallel coordinates plot for SARSA on CartPole-v1.

Figure 5: SARSA on CartPole-v1:Parallel coordinates plot.

Figure 6 shows the rewards for different runs for SARSA on CartPole-v1.
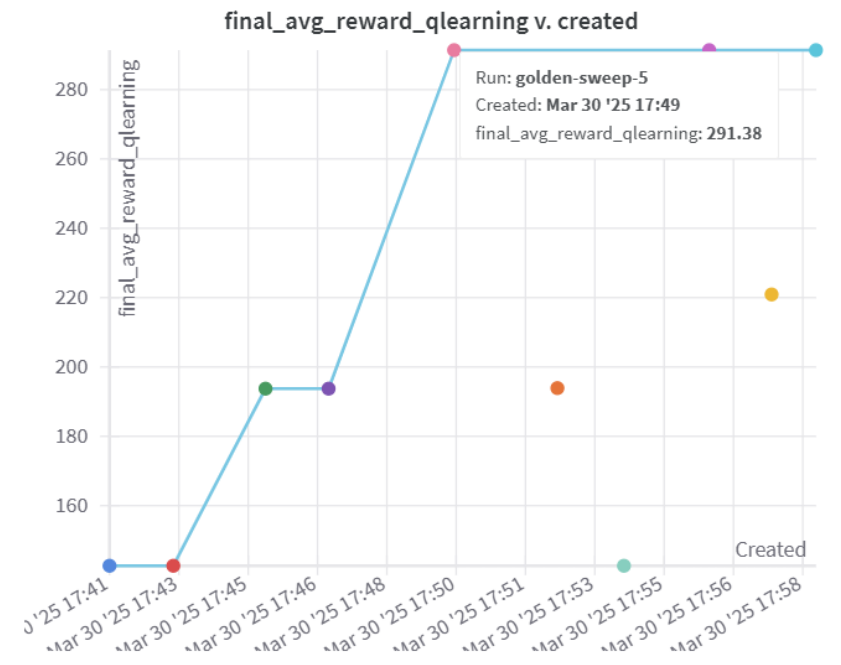


Figure 6: SARSA on CartPole-v1: Rewards for different runs.

The best hyperparameters for SARSA on CartPole-v1 from the parallel coordinates plot were found to be:

- **Learning rate:** 0.1

- **Epsilon Decay:** 0.995

- **Number of bins:** 25

The hyperparameters taken for Q-Learning on CartPole-v1 are:

- **Learning rate**

- **Tau Decay**

- **Number of bins**

Figure 7 shows the final parallel coordinates plot for Q-Learning on CartPole-v1.



Figure 7: Q-Learning on CartPole-v1: Parallel coordinates plot.

Figure 8 shows the rewards for different runs for Q-Learning on CartPole-v1.



Figure 8: Q-Learning on CartPole-v1: Rewards for different runs.

The best hyperparameters for Q-Learning on CartPole-v1 from the parallel coordinates plot were found to be:

- **Learning rate:** 0.1

- **Tau Decay:** 20

- **Number of bins:** 0.9999

The hyperparameters taken for SARSA on MountainCar-v0 are:

- **Learning rate**

- **Epsilon Decay**

- **Number of bins**

Figure 9 shows the final parallel coordinates plot for SARSA on MountainCar-v0.
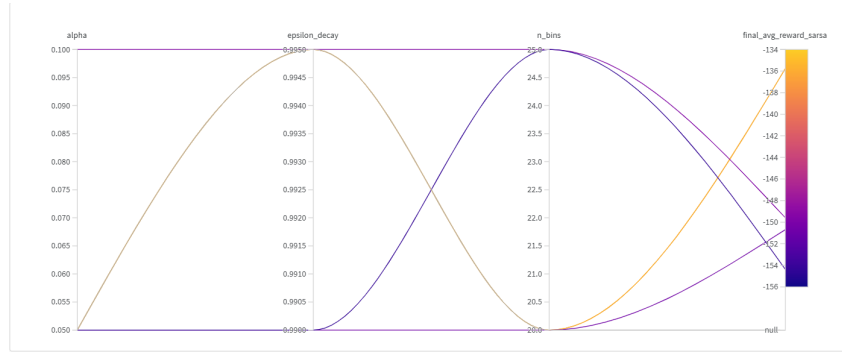


Figure 9: SARSA on MountainCar-v0: Parallel coordinates plot

Figure 10 shows the rewards for different runs for SARSA on MountainCar-v0 with.
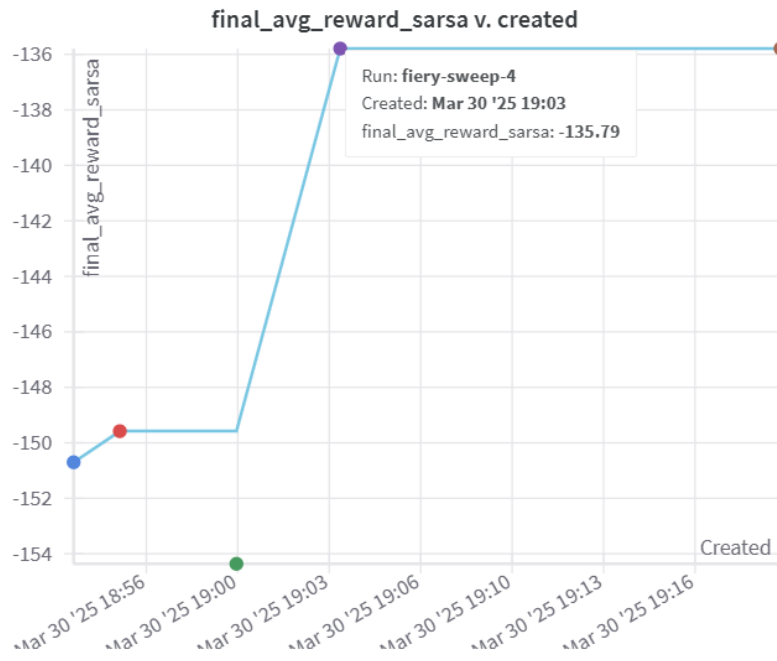


Figure 10: SARSA on MountainCar-v0: Rewards for different runs.

The best hyperparameters for SARSA on MountainCar-v0 from the parallel coordinates plot were found to be:

- **Learning rate:** 0.05

- **Epsilon Decay:** 0.995

- **Number of bins:** 20

The hyperparameters taken for Q-Learning on MountainCar-v0 are:

- **Learning rate**

- **Tau Decay**

- **Number of bins**

Figure 11 shows the final parallel coordinates plot for Q-Learning on MountainCar-v0.
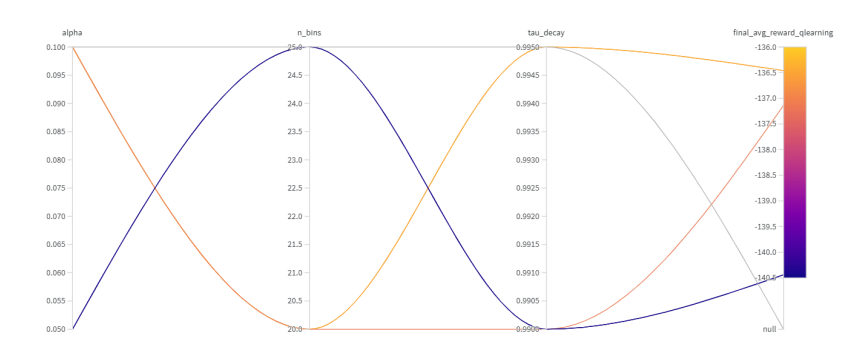


Figure 11: Q-Learning on MountainCar-v0: Parallel coordinates plot

Figure 12 shows the rewards for different runs for Q-Learning on MountainCar-v0.
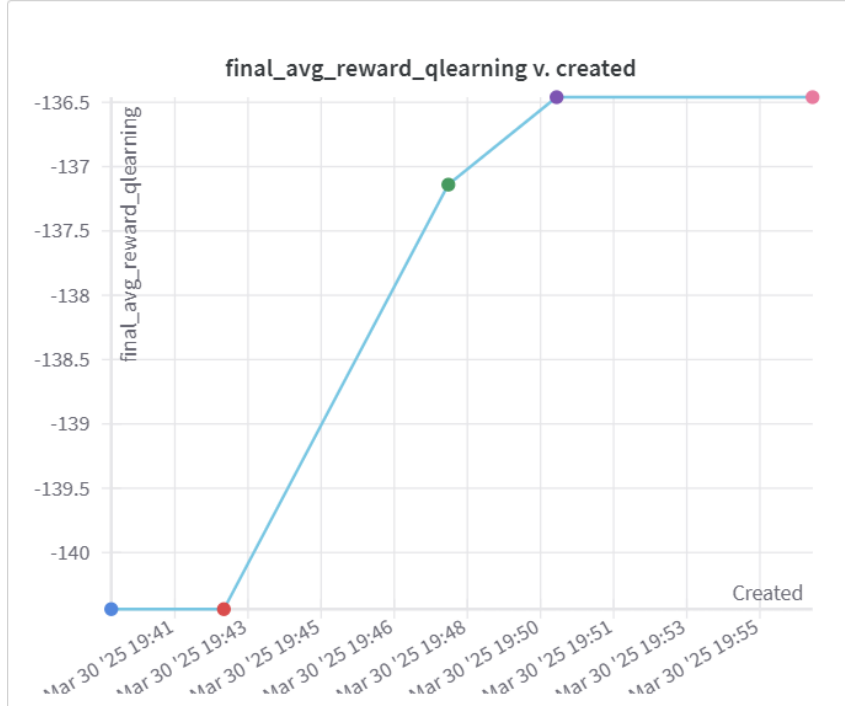


Figure 12: Q-Learning on MountainCar-v0:Rewards for different runs.

The best hyperparameters for Q-Learning on MountainCar-v0 from the parallel coordinates plot were found to be:

- **Learning rate:** 0.1

- **Tau Decay:** 20

- **Number of bins:** 0.995

# 7  Inferences and Conjectures

Based on our experiments, we observed the following:

- The number of bins can affect environmets like CartPole-v0 due to the definition of its states. A higher number of bins could help in better state representation as it can be seen in the hyperparameter tuning for CartPole-v0 using Sarsa.

- For both CartPole-v1 and MountainCar-v0, if the best hyperparameters are used, Q-Learning seems to give higher rewards than Sarsa.

- Q-Learning also lconverges quicker than Sarsa, as it can be seen from the episodic return plots that the Q-Learning plots flatten much earlier than the Sarsa Plots

- Sarsa tends to take lesser time per episode.

- The Parallel Coordinate Plots are used to find the best hyperparameters.

# 8   Code Repository

Our complete code is available on GitHub: `https://github.com/DA6400-RL-JanMay2025/programming-assignment-01-ch21b053_me21b171`

# 9   Bonus Task

We have also done the bonus task of evaluating the performance of SARSA and Q-Learning in the MiniGrid-Dynamic-Obstacles-5x5-v0.

The set of best hyperparameters of SARSA is: alpha0.134, epsilon 0.3117, decay 0.986, min 0.165 and seed 46.

And the set of best hyperparameters of Q-Learning is: alpha 0.1, decay 0.995, tau 1.0, min 0.01, seed 42.

Figure 13 shows the final reward plot distribution for Sarsa on MiniGrid-Dynamic-Obstacles5x5-v0.
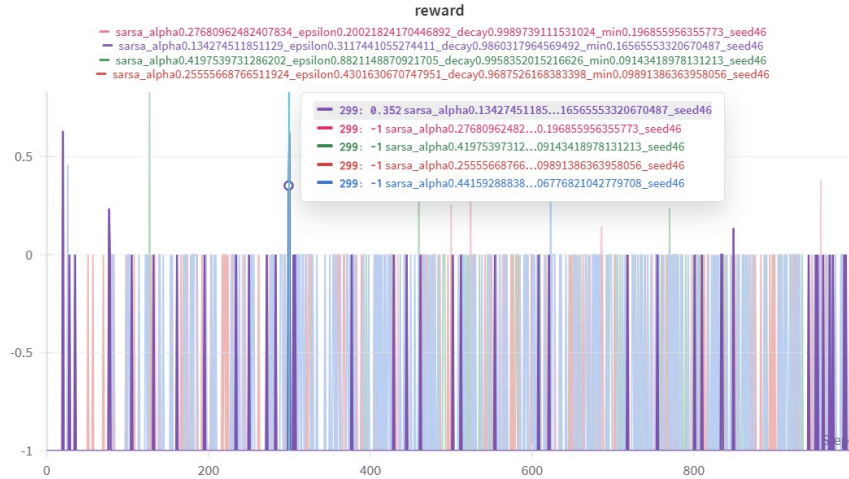


Figure 13: Sarsa on MiniGrid-Dynamic-Obstacles5x5-v0: Reward Plot Distribution.

Figure 14 shows the final reward plot distribution for Q-Learning on MiniGrid-Dynamic-Obstacles5x5-v0.
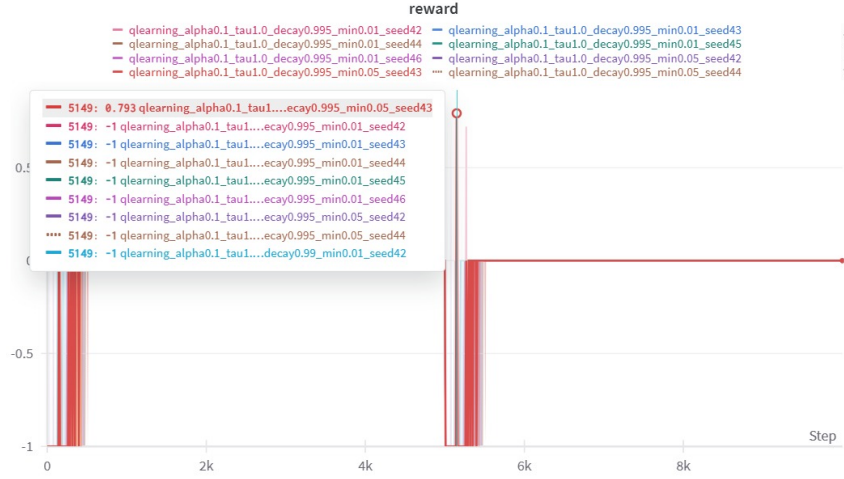
Figure 14: Q-Learning on MiniGrid-Dynamic-Obstacles5x5-v0: Reward Plot Distribution.

## 10 Conclusion

In this assignment, we implemented and compared SARSA and Q-Learning on multiple Gymnasium environments. Our experiments and hyperparameter tuning demonstrate the importance of appropriate discretization and exploration strategies in tabular reinforcement learning.