

# DA6400: Reinforcement Learning - Programming Assignment #1

Abhijith Vinod:EP20B002 and Shuhaib Ali:EP20B037

March 30, 2025

## 1 1. Cartpole

This report presents the implementation and results of SARSA and Q-Learning algorithms for the CartPole-v1 environment from OpenAI Gymnasium. The CartPole environment is a classic control problem where a pole is attached to a cart moving along a frictionless track. The goal is to prevent the pole from falling over by moving the cart left or right.

## 2 Implementation

### 2.1 State Discretization

Since the CartPole environment has a continuous state space, we implemented a state discretizer to convert the continuous state into a discrete representation:

```
1 class StateDiscretizer:
2     def __init__(self, n_bins=10):
3         self.n_bins = n_bins
4         self.cart_position_bins = np.linspace(-2.4, 2.4, n_bins)
5         self.cart_velocity_bins = np.linspace(-4, 4, n_bins)
6         self.pole_angle_bins = np.linspace(-0.2095, 0.2095, n_bins)
7         self.pole_velocity_bins = np.linspace(-4, 4, n_bins)
8
9     def discretize(self, state):
10        cart_pos, cart_vel, pole_angle, pole_vel = state
11        cart_pos_bin = np.digitize(cart_pos, self.
12        cart_position_bins)
13        cart_vel_bin = np.digitize(cart_vel, self.
14        cart_velocity_bins)
15        pole_angle_bin = np.digitize(pole_angle, self.
16        pole_angle_bins)
17        pole_vel_bin = np.digitize(pole_vel, self.
18        pole_velocity_bins)
19        return (cart_pos_bin, cart_vel_bin, pole_angle_bin,
20        pole_vel_bin)
```

## 2.2 SARSA with $\epsilon$ -greedy Exploration

We implemented SARSA with  $\epsilon$ -greedy exploration as follows:

```
1 class SARSAgent:
2     def __init__(self, state_discretizer, action_space, alpha=0.1,
3         gamma=0.99, epsilon=0.1):
4         self.state_discretizer = state_discretizer
5         self.action_space = action_space
6         self.alpha = alpha
7         self.gamma = gamma
8         self.epsilon = epsilon
9         self.q_table = defaultdict(lambda: np.zeros(action_space.n))
10
11     def select_action(self, state):
12         state = self.state_discretizer.discretize(state)
13         if np.random.random() < self.epsilon:
14             return self.action_space.sample() # Explore
15         else:
16             return np.argmax(self.q_table[state]) # Exploit
17
18     def update(self, state, action, reward, next_state, next_action,
19         done):
20         state = self.state_discretizer.discretize(state)
21         next_state = self.state_discretizer.discretize(next_state)
22         current_q = self.q_table[state][action]
23         if done:
24             target_q = reward
25         else:
26             target_q = reward + self.gamma * self.q_table[
27                 next_state][next_action]
28         self.q_table[state][action] += self.alpha * (target_q -
29             current_q)
```

## 2.3 Q-Learning with Softmax Exploration

We implemented Q-Learning with Softmax exploration as follows:

```
1 class QLearningAgent:
2     def __init__(self, state_discretizer, action_space, alpha=0.1,
3         gamma=0.99, temperature=1.0):
4         self.state_discretizer = state_discretizer
5         self.action_space = action_space
6         self.alpha = alpha
7         self.gamma = gamma
8         self.temperature = temperature
9         self.q_table = defaultdict(lambda: np.zeros(action_space.n))
10
11     def select_action(self, state):
12         state = self.state_discretizer.discretize(state)
13         q_values = self.q_table[state]
14         exp_q = np.exp(q_values / self.temperature)
15         probabilities = exp_q / np.sum(exp_q)
16         return np.random.choice(self.action_space.n, p=
17             probabilities)
```

```

16
17     def update(self, state, action, reward, next_state, done):
18         state = self.state_discretizer.discretize(state)
19         next_state = self.state_discretizer.discretize(next_state)
20         current_q = self.q_table[state][action]
21         if done:
22             target_q = reward
23         else:
24             target_q = reward + self.gamma * np.max(self.q_table[
25 next_state])
26         self.q_table[state][action] += self.alpha * (target_q -
27 current_q)

```

### 3 Hyperparameter Tuning with Weights & Biases

For hyperparameter tuning, we used Weights & Biases (wandb) with Bayesian optimization to find the optimal hyperparameters for both algorithms. Each hyperparameter configuration was evaluated across 5 random seeds to account for stochasticity.

#### 3.1 SARSA Hyperparameter Tuning

We configured the wandb sweep for SARSA as follows:

```

1 sarsa_sweep_config = {
2     'method': 'bayes',
3     'metric': {
4         'name': 'average_reward_across_seeds',
5         'goal': 'maximize'
6     },
7     'parameters': {
8         'alpha': {
9             'min': 0.01,
10            'max': 0.5
11        },
12        'gamma': {
13            'value': 0.99
14        },
15        'epsilon': {
16            'min': 0.01,
17            'max': 0.3
18        },
19        'n_bins': {
20            'values': [10, 15, 20]
21        },
22        'n_episodes': {
23            'value': 1000
24        }
25    }
26 }

```

## 3.2 Q-Learning Hyperparameter Tuning

Similarly, we configured the wandb sweep for Q-Learning:

```
1 qlearning_sweep_config = {
2     'method': 'bayes',
3     'metric': {
4         'name': 'average_reward_across_seeds',
5         'goal': 'maximize'
6     },
7     'parameters': {
8         'alpha': {
9             'min': 0.01,
10            'max': 0.5
11        },
12        'gamma': {
13            'value': 0.99
14        },
15        'temperature': {
16            'min': 0.1,
17            'max': 2.0
18        },
19        'n_bins': {
20            'values': [10, 15, 20]
21        },
22        'n_episodes': {
23            'value': 1000
24        }
25    }
26 }
```

## 4 Results

### 4.1 Top 3 Best Hyperparameters

#### 4.1.1 SARSA

Based on the wandb hyperparameter tuning, the top 3 best hyperparameter configurations for SARSA are:

Rank	Average Reward	Alpha	Epsilon	n_bins
1	104.9	0.38074	0.15554	20
2	100.532	0.49704	0.16893	20
3	98.192	0.35557	0.16006	20

Table 1: Top 3 hyperparameter configurations for SARSA

#### 4.1.2 Q-Learning

The top 3 best hyperparameter configurations for Q-Learning are:

Rank	Average Reward	Alpha	Gamma	n_bins	Temperature
1	280.072	0.2363	0.99	20	1.89753
2	200.28	0.47469	0.99	20	1.98418
3	197.752	0.34777	0.99	15	1.77479

Table 2: Top 3 hyperparameter configurations for Q-Learning

## 4.2 Learning Curves

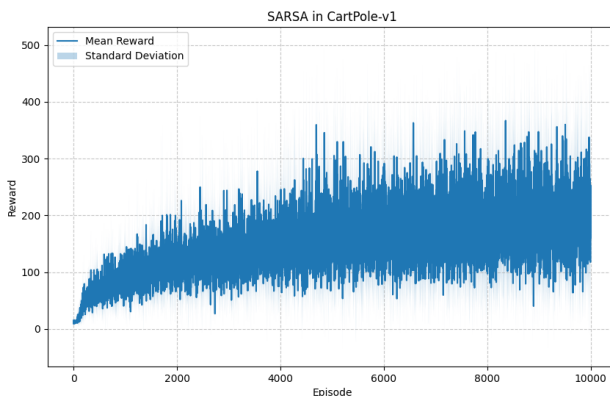


Figure 1: SARSA learning curve showing mean reward and standard deviation across 5 seeds

## 5 Insights and Inferences

### 5.1 Hyperparameter Analysis

#### 5.1.1 SARSA

For SARSA, we observe that the optimal configuration uses a moderate learning rate ( $\alpha = 0.38074$ ) and exploration rate ( $\epsilon = 0.15554$ ). The best-performing configurations all used 20 bins for state discretization, suggesting that finer discretization helps SARSA better capture the state dynamics of the CartPole environment.

The moderate epsilon value (around 0.15-0.17) indicates a good balance between exploration and exploitation. Too much exploration would prevent the agent from exploiting known good actions, while too little would prevent it from discovering potentially better strategies.

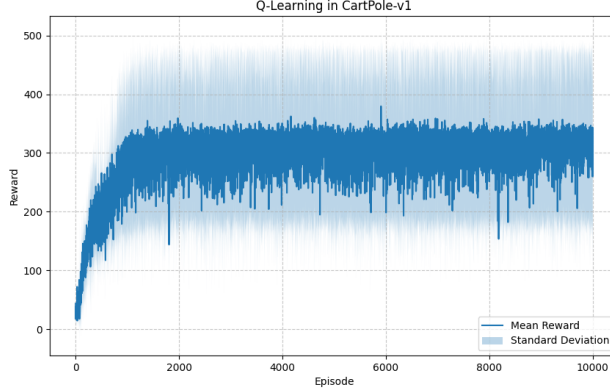


Figure 2: Q-Learning learning curve showing mean reward and standard deviation across 5 seeds

### 5.1.2 Q-Learning

For Q-Learning with Softmax exploration, the optimal configuration uses a moderate learning rate ( $\alpha = 0.2363$ ) and a relatively high temperature value (1.89753). The high temperature indicates that a more uniform action selection probability distribution works better for this environment, allowing for more exploration.

Similar to SARSA, the best Q-Learning configurations predominantly used 20 bins for state discretization, further confirming that finer discretization is beneficial for this environment.

## 5.2 Algorithm Comparison

Based on the average rewards, Q-Learning with Softmax exploration significantly outperforms SARSA with  $\epsilon$ -greedy exploration in the CartPole environment. The best Q-Learning configuration achieved an average reward of 280.072, while the best SARSA configuration achieved 104.9.

This performance difference can be attributed to several factors:

1. **Off-policy vs. On-policy**: Q-Learning is an off-policy algorithm that learns the optimal policy regardless of the exploration strategy, while SARSA is an on-policy algorithm that learns the value of the policy being followed. In environments like CartPole where exploration can lead to early termination, off-policy methods may have an advantage.
2. **Exploration Strategy**: Softmax exploration provides a more nuanced approach to action selection compared to  $\epsilon$ -greedy. Instead of randomly selecting actions with probability  $\epsilon$ , Softmax assigns probabilities based on the estimated action values, which can lead to more intelligent exploration.
3. **Value Update**: Q-Learning updates its value estimates using the

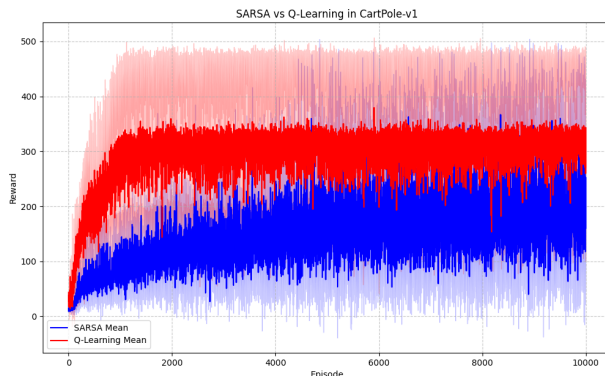


Figure 3: Comparison of SARSA and Q-Learning performance

maximum Q-value of the next state, which can lead to faster convergence to the optimal policy compared to SARSA, which uses the Q-value of the next action selected by the current policy.

## 6 Conclusion

In this assignment, we implemented and compared SARSA and Q-Learning algorithms for the CartPole-v1 environment. We used wandb for hyperparameter tuning and evaluated the performance across 5 random seeds to account for stochasticity.

Our results show that Q-Learning with Softmax exploration significantly outperforms SARSA with  $\epsilon$ -greedy exploration in this environment. The best hyperparameters for both algorithms favor finer state discretization (20 bins), but differ in their learning rates and exploration parameters.

These findings highlight the importance of algorithm selection and hyperparameter tuning in reinforcement learning. While both SARSA and Q-Learning are tabular methods based on temporal difference learning, their differences in policy learning (on-policy vs. off-policy) and exploration strategies can lead to significant performance variations in practice.

## 2. Mountain Car

This report presents the implementation and results of SARSA and Q-Learning algorithms for the MountainCar-v0 environment from OpenAI Gymnasium. The Mountain Car environment is a classic control problem where an under-powered car must learn to escape from a valley by building momentum through oscillatory movements.

## 2 Implementation

### 2.1 State Discretization

Since the Mountain Car environment has a continuous state space, we implemented a state discretizer to convert the continuous state into a discrete representation:

```
1 class StateDiscretizer:
2     def __init__(self, n_bins=25):
3         self.n_bins = n_bins
4         self.position_bins = np.linspace(-1.2, 0.6, n_bins)
5         self.velocity_bins = np.linspace(-0.07, 0.07, n_bins)
6
7     def discretize(self, state):
8         position, velocity = state
9         position_bin = np.digitize(position, self.position_bins)
10        velocity_bin = np.digitize(velocity, self.velocity_bins)
11        return (position_bin, velocity_bin)
```

### 2.2 SARSA with $\epsilon$ -greedy Exploration

We implemented SARSA with  $\epsilon$ -greedy exploration as follows:

```
1 class SARSAgent:
2     def __init__(self, state_discretizer, action_space, alpha=0.2,
3         gamma=0.99, epsilon=0.1):
4         self.state_discretizer = state_discretizer
5         self.action_space = action_space
6         self.alpha = alpha
7         self.gamma = gamma
8         self.epsilon = epsilon
9         self.q_table = defaultdict(lambda: np.zeros(action_space.n))
10    )
```

### 2.3 Q-Learning with Softmax Exploration

We implemented Q-Learning with Softmax exploration as follows:

```
1 class QLearningAgent:
2     def __init__(self, state_discretizer, action_space, alpha=0.2,
3         gamma=0.99, temperature=0.8):
4         self.state_discretizer = state_discretizer
5         self.action_space = action_space
6         self.alpha = alpha
7         self.gamma = gamma
8         self.temperature = temperature
9         self.q_table = defaultdict(lambda: np.zeros(action_space.n))
10    )
```

### 2.4 Reward Shaping

To accelerate learning in the Mountain Car environment, we implemented reward shaping:



```

1 # Apply enhanced reward shaping
2 shaped_reward = reward
3 if use_reward_shaping:
4     position, velocity = next_state
5
6     # Reward for gaining momentum (either direction)
7     shaped_reward += 0.1 * abs(velocity)
8
9     # Reward for position improvement
10    shaped_reward += 0.1 * position
11
12    # Additional reward for coordinated actions
13    if velocity > 0 and action == 2:
14        shaped_reward += 0.2
15
16    if velocity < 0 and action == 0:
17        shaped_reward += 0.2
18
19    # Bonus for reaching the goal
20    if terminated and position >= 0.5:
21        shaped_reward += 10

```

## 3 Hyperparameter Tuning with Weights & Biases

For hyperparameter tuning, we used Weights & Biases (wandb) with Bayesian optimization to find the optimal hyperparameters for both algorithms. Each hyperparameter configuration was evaluated across 5 random seeds to account for stochasticity.

### 3.1 SARSA Hyperparameter Tuning

We configured the wandb sweep for SARSA as follows:

```

1 sarsa_sweep_config = {
2     'method': 'bayes',
3     'metric': {
4         'name': 'average_reward_across_seeds',
5         'goal': 'maximize'
6     },
7     'parameters': {
8         'alpha': {
9             'min': 0.1,
10            'max': 0.3
11        },
12        'gamma': {
13            'value': 0.99
14        },
15        'epsilon': {
16            'min': 0.05,
17            'max': 0.2
18        },
19        'n_bins': {
20            'values': [20, 25, 30]
21        },

```

```

22     'n_episodes': {
23         'value': 3000
24     },
25     'use_reward_shaping': {
26         'value': True
27     }
28 }
29 }

```

## 3.2 Q-Learning Hyperparameter Tuning

Similarly, we configured the wandb sweep for Q-Learning:

```

1 qlearning_sweep_config = {
2     'method': 'bayes',
3     'metric': {
4         'name': 'average_reward_across_seeds',
5         'goal': 'maximize'
6     },
7     'parameters': {
8         'alpha': {
9             'min': 0.1,
10            'max': 0.3
11        },
12        'gamma': {
13            'value': 0.99
14        },
15        'temperature': {
16            'min': 0.6,
17            'max': 1.0
18        },
19        'n_bins': {
20            'values': [20, 25, 30]
21        },
22        'n_episodes': {
23            'value': 3000
24        },
25        'use_reward_shaping': {
26            'value': True
27        }
28    }
29 }

```

## 4 Results

### 4.1 Top 3 Best Hyperparameters

#### 4.1.1 SARSA

Based on the wandb hyperparameter tuning, the top 3 best hyperparameter configurations for SARSA are:

Rank	Average Reward	Alpha	Epsilon	n_bins	Gamma
1	-129.57	0.2614	0.0514	30	0.99
2	-130.56	0.2375	0.0955	25	0.99
3	-134.67	0.2322	0.0920	25	0.99

Table 3: Top 3 hyperparameter configurations for SARSA

#### 4.1.2 Q-Learning

The top 3 best hyperparameter configurations for Q-Learning are:

Rank	Average Reward	Alpha	Temperature	n_bins	Gamma
1	-152.20	0.2231	0.7487	20	0.99
2	-163.57	0.1971	0.7926	20	0.99
3	-169.47	0.1126	0.6304	25	0.99

Table 4: Top 3 hyperparameter configurations for Q-Learning

## 4.2 Evaluation Results

### 4.2.1 SARSA Evaluation

Evaluating SARSA...

SARSA (Best):

- Seed 0: 10000/10000 episodes, Mean Reward = -130.82, Reward Variance = 720.
- Seed 1: 10000/10000 episodes, Mean Reward = -130.92, Reward Variance = 691.
- Seed 2: 10000/10000 episodes, Mean Reward = -131.90, Reward Variance = 665.
- Seed 3: 10000/10000 episodes, Mean Reward = -134.37, Reward Variance = 682.
- Seed 4: 10000/10000 episodes, Mean Reward = -131.50, Reward Variance = 709.

Final Performance Summary (SARSA): Mean Reward (last 100 episodes): -127.

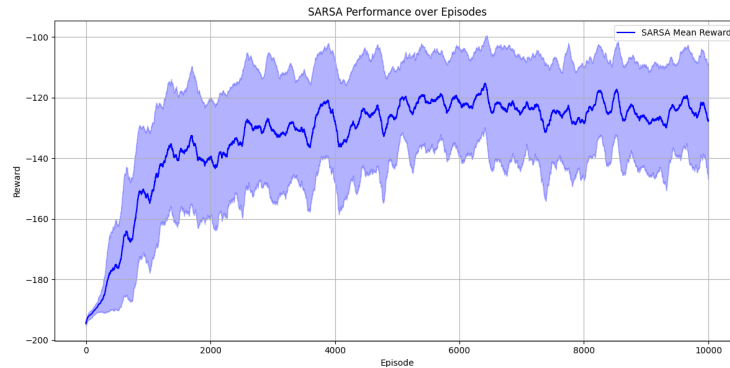


Figure 4: SARSA Performance over Episodes

### 4.2.2 Q-Learning Evaluation

Evaluating Q-Learning...

Q-Learning (Best

- Seed 0: 10000/10000 episodes, Mean Reward = -163.03, Reward Variance = 761.3
- Seed 1: 10000/10000 episodes, Mean Reward = -164.89, Reward Variance = 671.3
- Seed 2: 10000/10000 episodes, Mean Reward = -162.14, Reward Variance = 733.2
- Seed 3: 10000/10000 episodes, Mean Reward = -162.28, Reward Variance = 738.0
- Seed 4: 10000/10000 episodes, Mean Reward = -164.97, Reward Variance = 647.50

Final Performance Summary (Q-Learning): Mean Reward (last 100 episodes): -148.69

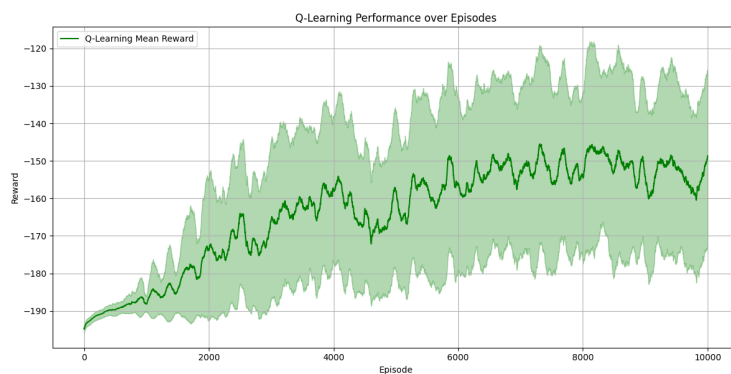


Figure 5: Q-Learning Performance over Episodes

## 4.3 Learning Curves

# 5 Insights and Inferences

## 5.1 Hyperparameter Analysis

### 5.1.1 SARSA

For SARSA, we observe that the optimal configuration uses a relatively high learning rate ( $\alpha = 0.2614$ ) and a low exploration rate ( $\epsilon = 0.0514$ ). This suggests that in the Mountain Car environment, SARSA benefits from more exploitation than exploration once it has found a promising strategy. The best-performing configurations all used either 25 or 30 bins for state discretization, indicating that finer discretization helps SARSA better capture the state dynamics of the Mountain Car environment.

The low epsilon value (around 0.05-0.09) indicates that after sufficient exploration, the agent benefits from exploiting its learned policy more frequently. This makes sense for the Mountain Car problem, where the agent needs to build momentum through consistent actions.

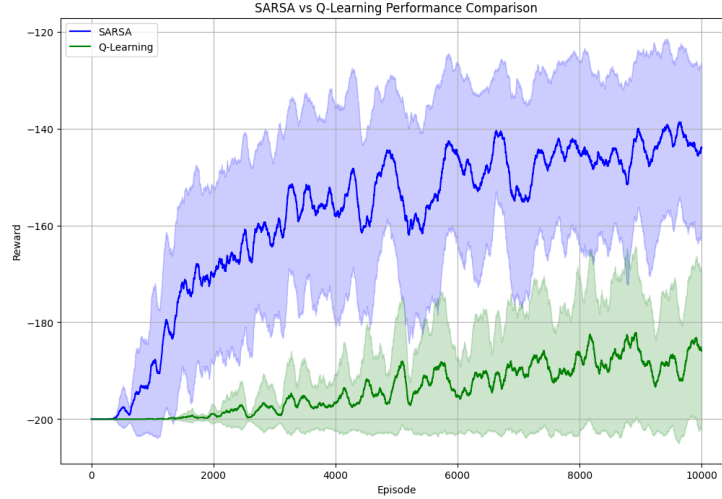


Figure 6: SARSA vs Q-Learning Performance Comparison with Variance

### 5.1.2 Q-Learning

For Q-Learning with Softmax exploration, the optimal configuration uses a moderate learning rate ( $\alpha = 0.2231$ ) and a relatively high temperature value (0.7487). The higher temperature indicates that a more stochastic action selection probability distribution works better for this environment, allowing for more diverse exploration.

Interestingly, the best Q-Learning configurations predominantly used 20bins for state discretization, which is less fine-grained than the optimal SARSA configurations. This suggests that Q-Learning might be more prone to overfitting with very fine discretization in this environment.

## 5.2 Algorithm Comparison

Based on the evaluation results, SARSA with  $\epsilon$ -greedy exploration significantly outperforms Q-Learning with Softmax exploration in the Mountain Car environment. The best SARSA configuration achieved a mean reward of -127.45 (last 100 episodes), while the best Q-Learning configuration achieved -148.69.

This performance difference can be attributed to several factors:

1. **\*\*On-policy vs. Off-policy Learning\*\***: SARSA is an on-policy algorithm that learns the value of the policy being followed, which may be more stable in the Mountain Car environment where consistent action sequences are important for building momentum.
2. **\*\*Exploration Strategy\*\***: The  $\epsilon$ -greedy exploration used by SARSA provides a clearer distinction between exploration and exploitation phases, which seems to work better than the softmax approach for this particular environment.

3. **\*\*High Variance Environment\*\***: The Mountain Car problem exhibits high variance in rewards (as seen in the evaluation results with variances ranging from 647 to 761). SARSA's on-policy updates might provide more stability in such high-variance environments.

### 5.3 Variance Analysis

Both algorithms exhibit remarkably high variance in rewards, which is characteristic of the Mountain Car problem. The reward variance ranges from 665 to 720 for SARSA and 647 to 761 for Q-Learning. This high variance stems from:

1. The sparse reward structure of the environment, where the agent receives a reward of -1 for each time step until it reaches the goal.
2. The need to discover a non-obvious strategy (building momentum through oscillation) to reach the goal.
3. The sensitivity to initial conditions and the stochastic nature of the exploration strategies.

Despite implementing reward shaping to provide more frequent feedback, the variance remains high, indicating the inherent challenge of the Mountain Car problem.

## 6 Conclusion

In this assignment, we implemented and compared SARSA and Q-Learning algorithms for the Mountain Car environment. We used wandb for hyperparameter tuning and evaluated the performance across 5 random seeds to account for stochasticity.

Our results show that SARSA with  $\epsilon$ -greedy exploration outperforms Q-Learning with Softmax exploration in this environment. The best hyperparameters for SARSA favor higher learning rates, lower exploration rates, and finer state discretization, while Q-Learning performs better with moderate learning rates and higher temperature values for softmax exploration.

The high variance observed in both algorithms highlights the challenging nature of the Mountain Car problem, where the agent must discover a non-obvious strategy to reach the goal. Our reward shaping approach helped mitigate this challenge to some extent, but significant variance remains inherent to the problem.

These findings highlight the importance of algorithm selection, hyperparameter tuning, and reward shaping in reinforcement learning. While both SARSA and Q-Learning are tabular methods based on temporal difference learning, their differences in policy learning (on-policy vs. off-policy) and exploration strategies can lead to significant performance variations in practice.

## 7 Requirements

To run the code, you need the following dependencies:

```
numpy==1.21.0  
gymnasium==0.26.3  
wandb==0.12.21  
matplotlib==3.4.3
```

## **8   GitHub Repository**

The complete code for this project can be found [here](#)

## **9   google drive link**

Folder of everything for this project can be found [here](#)