

# DA6400: Introduction to Reinforcement Learning

## Assignment 1

Jayagowtham J ME21B078, Lalit Jayanti ME21B096

## Algorithms and Training

### Agent

This is the common code block for both of our SARSA and Q Learning Agents. It encapsulates initializing discretizer, agent and agent environment interaction parameters. The functions *update\_* override the default parameters giving us more control over our environment.

```
1 # scripts/agents.py
2 class SARSAgent: / class QLearningAgent:
3     def __init__(self, state_space, action_space, seed):
4         ''' Initialises our agent with discretizer, agent and
5             agent environment parameters
6         '''
7
8     def reset(self):
9         '''Updates all the hyperparameters and related objects
10        '''
11
12    def update_hyperparameters(self, **kwargs):
13        '''Updates hyperparameters at the start of training
14        '''
15
16    def update_agent_parameters(self):
17        '''Updates agent training parameters based on decay type
18        '''
```

### SARSA

This code block deals with the state action value update step according to the SARSA algorithm. We execute the given action from our current state, observe the next state and reward and then choose the next action based on the  $\epsilon$ -greedy policy. Now, we update the state action value function based on the next state and chosen action value function.

```
1 # scripts/agents.py
2 class SARSAgent:
3     .
4     .
5
6     def step(self, state, action, reward, next_state, done):
7         '''Updates Q-values according to the SARSA
8             algorithm
9         '''
10    q_sa = self.q_table[state, action]
11
12    next_action_values = self.q_table[next_state]
13
14    # Choosing the next action
15    self.next_action = epsilon_greedy(
16        action_values=next_action_values,
17        action_size=self.action_size,
18        eps=self.eps
19    )
20
```

```

21     q_next_sa = next_action_values[self.next_action]
22
23     # SARSA update
24     self.q_table[state, action] = (
25         q_sa + self.LR*(reward + self.GAMMA*q_next_sa - q_sa)
26     )
27
28     def act(self, state):
29         '''Takes an action from current state based on
30         epsilon-greedy policy according to current Q values.
31         '''
32         action_values = self.q_table[state]
33         eps_action = self.next_action
34         return eps_action, action_values

```

## Q Learning

This code block deals with the state action value update step according to the Q Learning algorithm. We choose an action according to the softmax policy derived from our Q Table, observe the next state and reward. Now, we update the current state action value function based on the maximum possible state action value function.

```

1 # scripts/agents.py
2 class QLearningAgent:
3
4     .
5
6     def step(self, state, action, reward, next_state, done):
7         '''Updates Q-values according to the Q-learning
8         algorithm
9         '''
10        q_sa = self.q_table[state, action]
11
12        # Maximization of Q-value over actions
13        q_next_sa = np.max(self.q_table[next_state])
14
15        # Q-Learning update
16        self.q_table[state, action] = (
17            q_sa + self.LR*(reward + self.GAMMA*q_next_sa - q_sa)
18        )
19
20    def act(self, state):
21        '''Takes an action from current state based on
22        softmax policy according to current Q values.
23        '''
24        action_values = self.q_table[state]
25        softmax_action = softmax(
26            action_values=action_values,
27            action_size=self.action_size,
28            tau=self.tau
29        )
30        return softmax_action, action_values

```

## Tile Coding

We represent the states with their tile coordinates in multiple tilings hashed to give a single list of hashes. These hashes correspond to the indices in our Q Table which performs coarse coding (multiple updates for the same observation), resulting in a much better generalization. The following code blocks depict the algorithms used. We have based our implementation based on the description in section 9.5.4 of Reinforcement Learning: An Introduction. Further we have chosen to opt for asymmetrical tiling and used the offset vectors as the set of first n odd integers, where n is the number of dimensions of the state space, this choice is based on the work of Miller and Glanz (1996).

```

1 # scripts/tilecoding.py
2 class TileCoder:
3     '''Discretizes continuous states into discrete bins

```

```

4     with multiple tilings for state aggregation
5     Based on the description in section 9.5.4 of Reinforcement Learning: An Introduction
6     '',
7     def __call__(self, x):
8         eps = 1e-10
9         # Finding coordinates of states in different tilings
10        tile_coords = ((x - self.lower_lim) / (self.tile_width + eps)
11                      + self.offsets).astype(np.int32)
12
13        # Hashing the tile coords into a single number per tiling, size (num_tilings,1)
14        tile_coord_hash = (
15            np.dot(self.base_hash, tile_coords.T)
16        )
17        # Return a list of hashes of the size (num_tilings,1)
18        return list(self.base_tile_index + tile_coord_hash)
19
20 class QTable:
21     '''Maintains a 2 dimensional Q Table to store multiple representations
22     of a state from different tilings by hashing. Eases retrieval and
23     update of Q-values.
24     '''
25     def __getitem__(self, key):
26         """To get q_value = q_table[state] or q_value = q_table[state, action]
27         """
28         if isinstance(key, tuple) and len(key) == 2:
29             state, action = key
30             idx_s = self.tile_coder(state)
31             return self.table[idx_s, action].mean()
32         else:
33             state = key
34             idx_s = self.tile_coder(state)
35             return self.table[idx_s].mean(axis=0)
36
37     def __setitem__(self, key, value):
38         """To set q_table[state, action] = value
39         """
40         if not (isinstance(key, tuple) and len(key) == 2):
41             raise ValueError("Key must be a tuple (state, action)")
42
43         state, action = key
44         print(state, action)
45         idx_s = self.tile_coder(state)
46         self.table[idx_s, action] = value

```

## Training

This code block consists of the main training loop for the agent, it handles the agent's actions, stepping the environment forward and also updating the agent's parameters for learning. Further this code block also takes care of logging the results which are subsequently used to perform analysis and also hyperparameter tuning.

```

1 # scripts/training.py
2 class Trainer:
3
4     def training(self, env, agent, n_episodes=10000, process_training_info=lambda *args, **
5                 kwargs: (False, {})):
6         """
7             To train an agent in the given environment.
8
9             Args:
10                - env: The environment for training.
11                - agent: An agent with '.step()', '.act()', and '.update_agent_parameters()'.
12                - n_episodes (int, optional): Number of training episodes. Defaults to 10000.
13                - process_training_info (function, optional): Runs after each episode.
14                    - First return value must be a 'bool' for early stopping.
15                    - Second return value must be a 'dict' to update the progress bar's postfix.
16
17             Returns:

```

```

17         - dict: Summary of the training process.
18         """
19
20     # Initialize variables for logging here
21
22     progress_bar = tqdm(range(1, n_episodes+1), desc="Training")
23
24     # Training loop
25     for i_episode in progress_bar:
26         state, _ = env.reset()
27         score = 0
28         total_reward = 0
29         terminated, truncated = False, False
30         episode_history = []
31
32         # Running an episode
33         while not (terminated or truncated):
34             action, action_vals = agent.act(state)
35             next_state, reward, terminated, truncated, _ = env.step(action)
36             episode_history.append((state, action, reward, next_state))
37             agent.step(state, action, reward, next_state, terminated)
38             state = next_state
39             score += self.compute_score(reward)
40             total_reward += reward
41
42             agent.update_agent_parameters()
43
44         # Code for logging here
45
46     return {
47         "computation_time": end_time - begin_time,
48         "scores": np.array(history_scores),
49         "total_rewards": np.array(history_total_rewards)
50     }
51
52     def compute_score(self, reward):
53         return reward

```

## Metrics and Hyperparameters

To quantitatively compare the performance of the algorithms in an environment we have used the following two metrics:

- **Cumulative regret**: The sum of difference between return assuming optimal policy and the 5 run average return at each step. For this assignment we have tuned the hyperparameters to minimize the regret in all experiments.
- **Maximum rolling average score**: The maximum of the rolling average return over the past 100 episodes averaged over 5 runs

Further, based on the description in section 1, the agents are characterized by a few hyperparameters:

- **num\_tilings\_per\_feature**: Determines the total discrete bins dividing the continuous state-space.
- **num\_tilings**: Determines the total number of overlapping grids used in while using Tile Coding.
- **learning\_rate**: The step size with which the Q-values are updated.
- **eps\_start**: The starting value of  $\epsilon$  in the  $\epsilon$ -greedy policy.
- **eps\_end**: The minimum value  $\epsilon$  can take while training.
- **tau\_start**: The starting value of  $\tau$  in the softmax policy.
- **tau\_end**: The minimum value  $\tau$  can take while training.

- **decay\_type** : Determines the rate at which the parameter  $\epsilon$  or  $\tau$  is decayed over episodes, this chosen either to be linear or exponential for this case.
- **frac\_episodes\_to\_type** : Determines the fraction of total number of episodes where  $\epsilon$  or  $\tau$  is decayed, after which these parameters are held at their minimum values.

We have tuned the hyperparameters using wandb to minimize regret in all experiments, by performing a bayesian search over the search-space of hyperparameters. The exact configurations of the search-space for the hyperparameters can be found in the config/ folder.

## CartPole-v1

This environment's observation space consists of four continuous state variables, i.e. cart position, cart velocity, pole angle, and pole angular velocity. These variables range within  $[-4.8, -\infty, -0.41887903, -\infty]$  to  $[4.8, \infty, 0.41887903, \infty]$ .

To address this environment, we process the state as follows. We first clip each variable to the range  $[-5, 5]$  and then apply tile coding for discretization before passing it to the agent, as described in subsection (iv). We have observed empirically that the loss in information due to clipping does not lead to significant performance loss and yields satisfactory results in practice. We train the agent over a course of 10000 episodes.

### Part (i): SARSA

To tune the hyperparameters we ran a parameter sweep consisting of 20 samples using wandb and performed a bayesian search over the search-space of hyperparameters to minimize average cumulative regret, from these we selected the top 3 configurations, presented in Table 1.

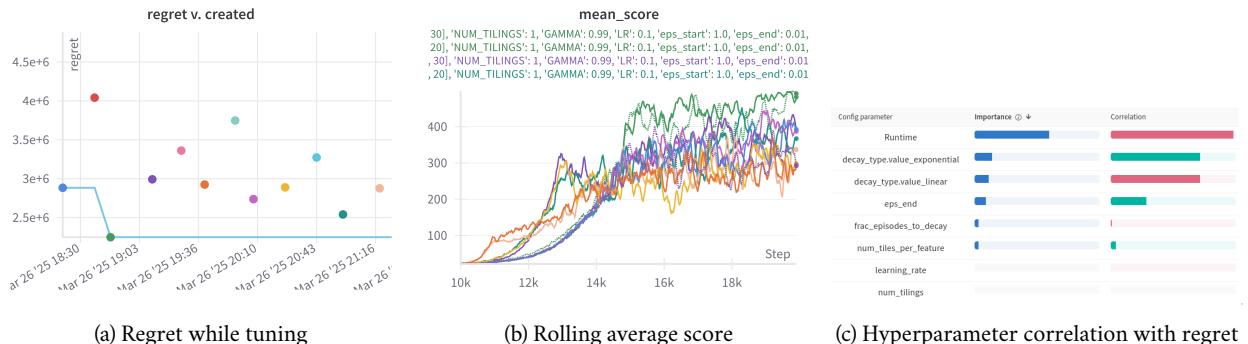


Figure 1: Tuning for CartPole-SARSA

Hyperparameters	Top	2nd Best	3rd Best
num_tiles_per_feature	30	20	30
num_tilings	1	1	1
learning_rate	0.1	0.1	0.1
eps_start	1	1	1
eps_end	0.01	0.01	0.01
decay_type	linear	linear	linear
frac_episodes_to_decay	0.3	0.5	0.5
Cumulative regret	2470134.6	2671750.8	2710446.4
Max rolling average score	443.366	444.438	449.752

Table 1: Top 3 Hyperparameter configurations for CartPole - SARSA

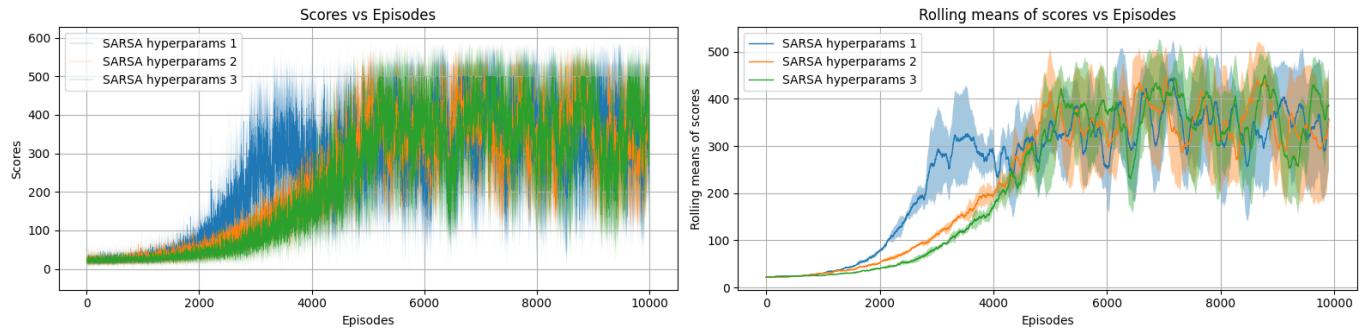


Figure 2: CartPole SARSA experiments for Top 3 Hyperparameter configurations

## Part (ii): Q-Learning

Similar to SARSA case, we ran another sweep for Q-Learning and selected the top 3 configurations from 20 samples, see Table 2. For each sample we averaged the scores and cumulative regret over 5 runs.

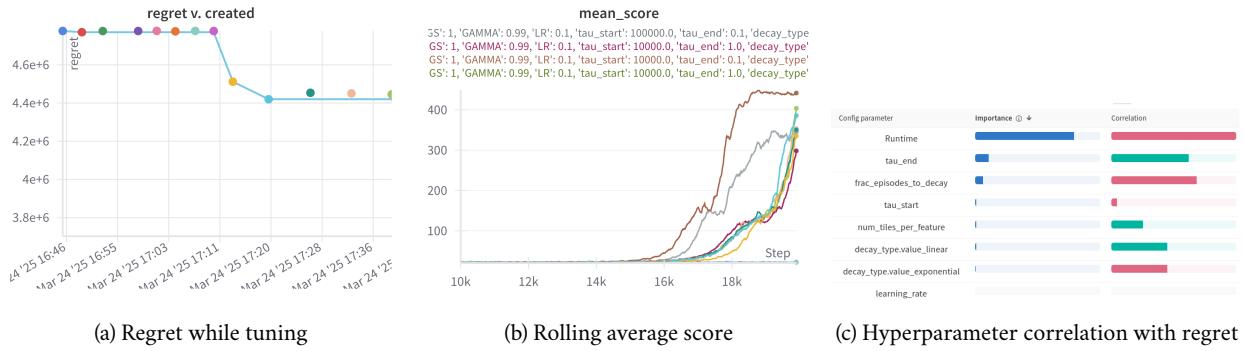


Figure 3: Tuning for CartPole-Q-Learning

Hyperparameters	Top	2nd Best	3rd Best
num_tiles_per_feature	30	20	30
num_tilings	1	1	1
learning_rate	0.1	0.1	0.1
tau_start	10000	10000	1000
tau_end	0.1	0.1	0.1
decay_type	exponential	exponential	exponential
frac_episodes_to_decay	0.7	0.7	1.0
Cumulative regret	2847478.6	3285443.4	3403292.0
Max rolling average score	432.358	340.840	453.070

Table 2: Top 3 Hyperparameter configurations for CartPole - Q-Learning

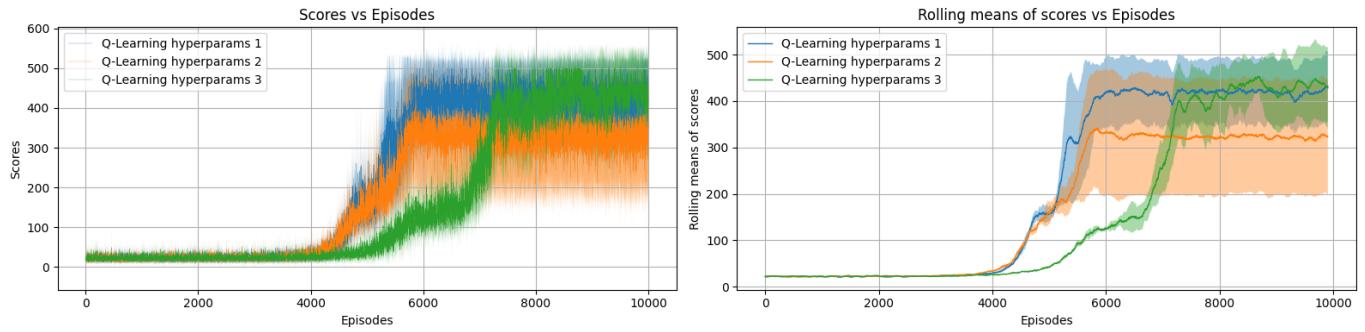


Figure 4: CartPole Q-Learning experiments for Top 3 Hyperparameter configurations

### Part (iii): Comparison and Inferences

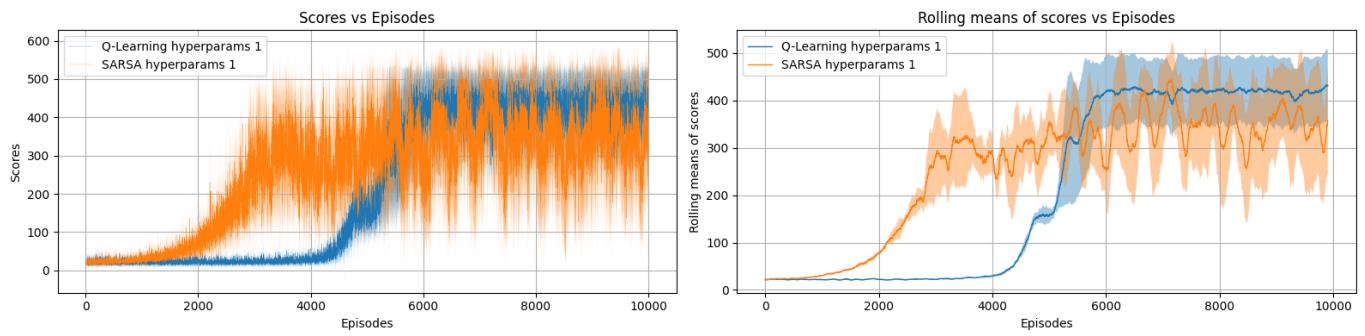


Figure 5: Comparing SARSA and Q-Learning for Top Hyperparameter configurations

In this experiment, we observe that SARSA initially outperforms Q-Learning. However, after approximately 7000 episodes, Q-Learning surpasses SARSA by a significant margin. This can be attributed to the fact that Q-Learning is an off-policy algorithm, where the exploratory nature of the softmax policy with a high  $\tau$  leads to greater regret.

As  $\tau$  decays to its final value (here 0.1) around 7000 episodes (as  $\text{frac\_episodes\_to\_decay} = 0.7$ ), the Q-Learning algorithm has effectively learned the optimal policy. At this point, it shifts toward exploitation, using the optimal Q-values to outperform SARSA.

## MountainCar-v0

The environment's observation space consists of two continuous state variables: the car's position along the x-axis and its velocity. These values range from  $[-1.2, -0.07]$  to  $[0.6, 0.07]$ . As before, we apply tile coding for discretization before passing the state to the agent, as described in subsection (iv).

To encourage exploration, we initialize all Q-values to 0, an optimistic value. This indirectly incentivizes the agent to explore previously un-explored actions since, in this experiment, all Q-values for states action pairs are negative, as all rewards are negative in this environment. We train the agent over a course of 10000 episodes.

### Part (i): SARSA

Similar to the previous environment, to tune the hyperparameters we ran a parameter sweep consisting of 20 samples using wandb and performed a bayesian search over the search-space of hyperparameters to minimize cumulative regret, from these we selected the top 3 configurations, presented in Table 3.

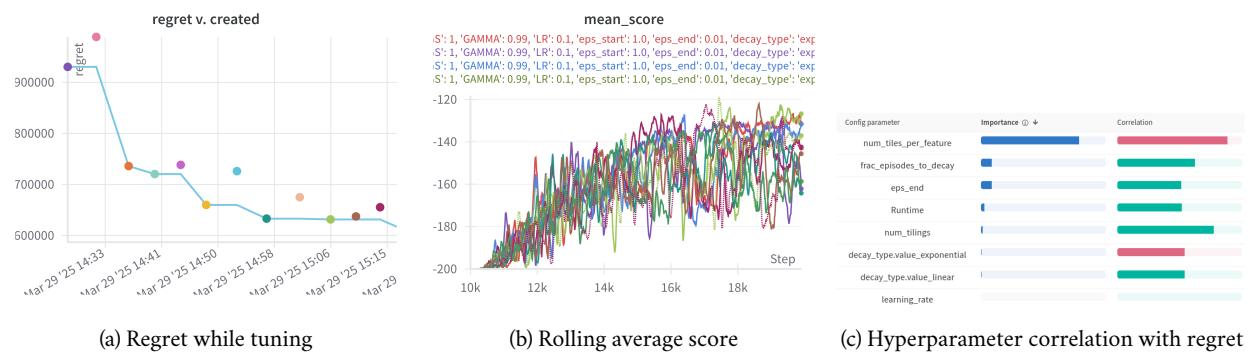


Figure 6: Tuning for MountainCar-SARSA

Hyperparameters	Top	2nd Best	3rd Best
num_tiles_per_feature	20	20	20
num_tilings	1	1	4
learning_rate	0.1	0.1	0.1
eps_start	1	1	1
eps_end	0.01	0.01	0.01
decay_type	exponential	exponential	exponential
frac_episodes_to_decay	0.1	0.3	0.5
Cumulative regret	608150.6	618654.0	685039.0
Max rolling average score	-135.624	-135.514	-138.364

Table 3: Top 3 Hyperparameter configurations for MountainCar - SARSA

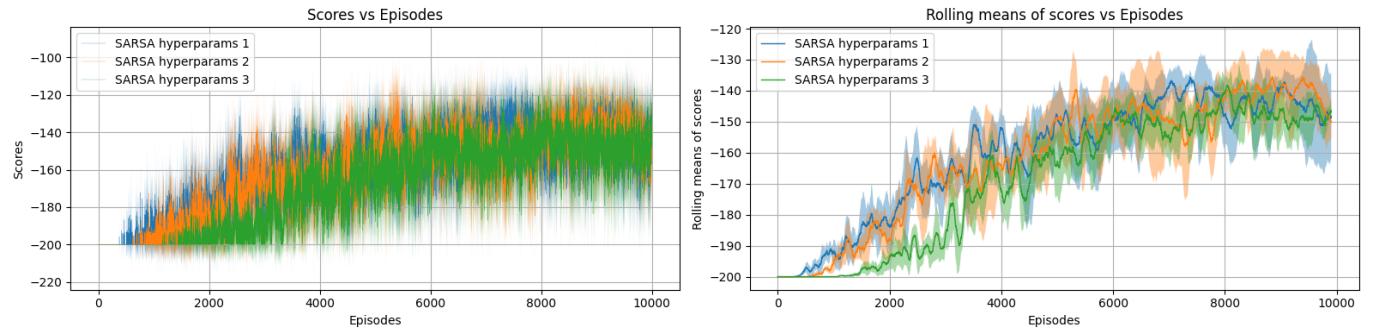


Figure 7: MountainCar SARSA experiments for Top 3 Hyperparameter configurations

## Part (ii): Q-Learning

Similar to the SARSA case, we ran another sweep for Q-Learning and selected the top 3 configurations from 20 samples, see Table 4.

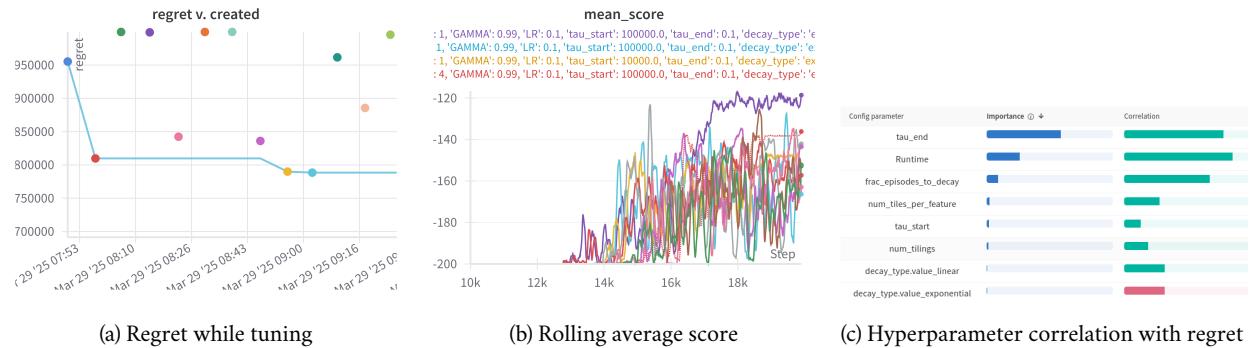


Figure 8: Tuning for MountainCar-Q-Learning

Hyperparameters	Top	2nd Best	3rd Best
num_tiles_per_feature	20	10	20
num_tilings	1	1	4
learning_rate	0.1	0.1	0.1
tau_start	100000	100000	100000
tau_end	0.1	0.1	0.1
decay_type	exponential	exponential	exponential
frac_episodes_to_decay	0.3	0.5	0.3
Cumulative regret	758822.0	809055.8	784852.2
Max rolling average score	-130.074	-146.634	-144.930

Table 4: Top 3 Hyperparameter configurations for MountainCar - Q-Learning

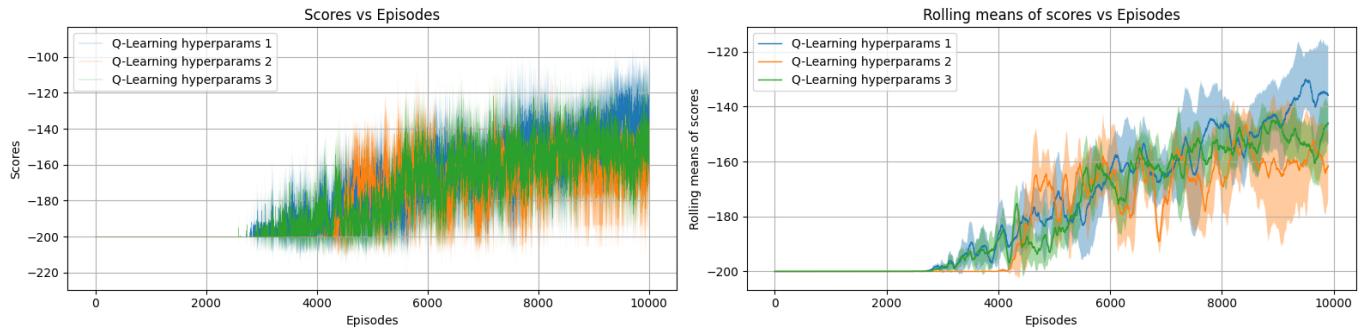


Figure 9: MountainCar Q-Learning experiments for Top 3 Hyperparameter configurations

### Part (iii): Comparison and Inferences

While solving the MountainCar environment, we observed that the agent takes approximately 1500 episodes to explore and reach the goal for the first time. We infer from this although initializing Q-values with optimistic values encourages exploration, it may take a long time before the agent performs expected actions to reach the goal. After reaching the goal once, the agent's regret decreases as it learns to achieve the goal more efficiently.

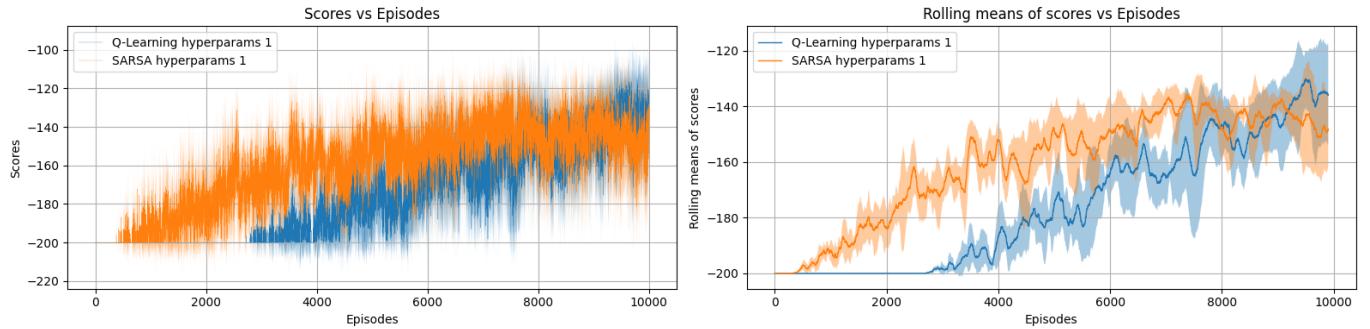


Figure 10: Comparing SARSA and Q-Learning for Top Hyperparameter configurations

### Part (iv): Solving MountainCar with Reward Shaping

As discussed in the previous section, the agent takes a significant amount of time to reach the goal for the first time. To address this, we consider an alternative approach using reward shaping. Based on our experiments, we observed that the agent must be incentivized to accumulate higher velocities (or, informally, gain more momentum). This helps the agent escape the valley and reach the goal. Therefore, we introduce a new reward scheme:

$$R_{t+1} = R_v - 1$$

where  $R_v$  represents the reward for taking actions that accelerate the agent in its direction of motion, and  $R_{t+1}$  is the reward for taking an action at time step  $t$ . Here  $R_v$  is defined as:

$$R_v = \begin{cases} -k \cdot v_t, & \text{if action} = 0 \text{ (accelerate left)} \\ 0, & \text{if action} = 1 \text{ (no acceleration)} \\ k \cdot v_t, & \text{if action} = 2 \text{ (accelerate right)} \end{cases}$$

where  $k$  is a positive constant that determines the weight for incentivizing acceleration in the direction of motion, and  $v_t$  is the velocity of the agent at time step  $t$ . For this experiment we set  $k$  to be 100.

In a more compact form:

$$R_v = k \cdot (\text{action} - 1) \cdot v_t$$

For demonstration purposes, we use the same hyperparameters as described in Table 3 and Table 4.

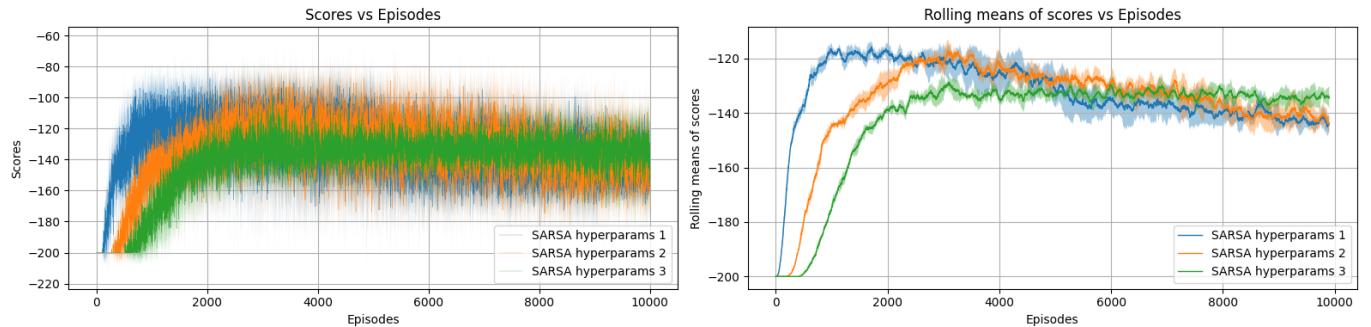


Figure 11: MountainCar SARSA experiments (with Reward shaping) for Top 3 Hyperparameter configurations

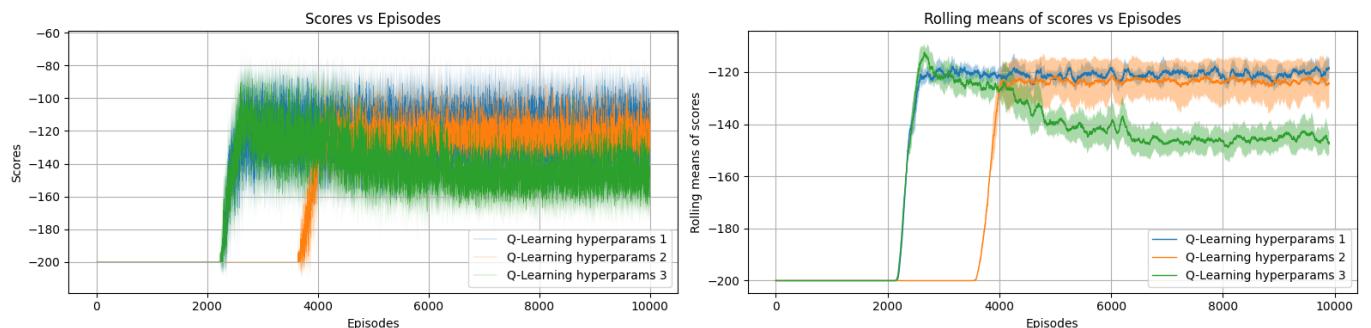


Figure 12: MountainCar Q-Learning experiments (with Reward shaping) for Top 3 Hyperparameter configurations

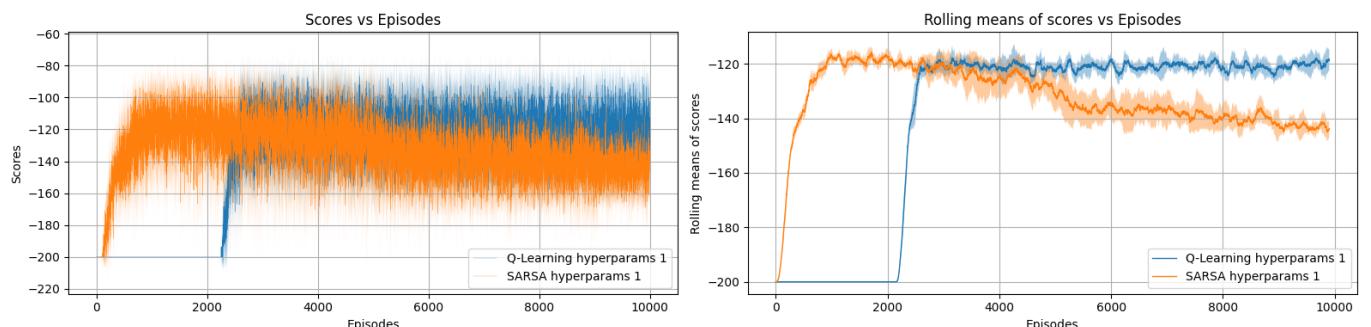


Figure 13: Comparing SARSA and Q-Learning (with Reward shaping) for Top Hyperparameter configurations

We observe the fact that the variance across five runs drastically decreases when the rewards are shaped. We can infer that the agent is consistently learning similar policies. We also observe a drastic improvement in rolling average performance, leading us to believe that the agent has learned a better policy than without shaping. SARSA seems to have a hold on the top policy for some time before it loses its grip. Q-Learning joins the race late, but clings to the top policy once it learns it. We believe that the drop in performance and stagnation at a sub-optimal score can be attributed to the agents optimizing for a different objective determined by  $R_v$ . As a result, an agent may spend more time oscillating in the valley to accumulate positive rewards before gaining enough velocity to escape, leading to sub-optimal behavior. We believe that this problem can be solved to some degree by tuning the weight  $k$  described earlier.

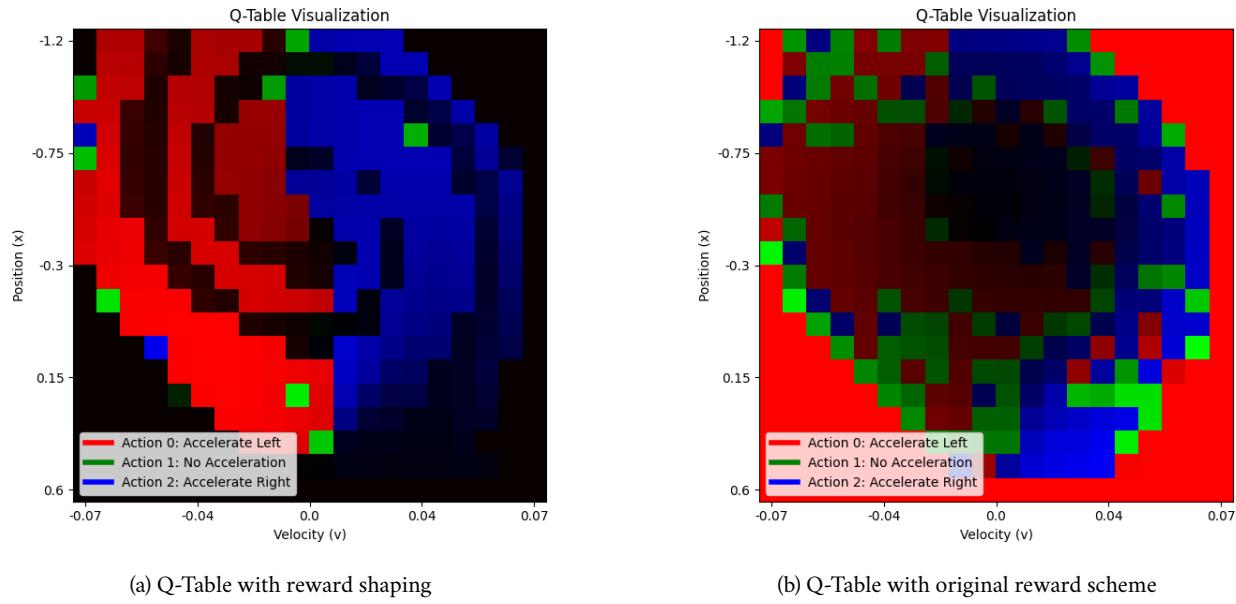


Figure 14: Visualizing the Q-Tables

To further analyze the policies learned by the agents with and without reward shaping, we plotted the optimal action for each state (position, velocity of the car), as shown in Figure 14.

We observed a spiral-shaped structure in the Q-table, centered approximately at the starting state. Additionally, we inferred that the Q-table using the original reward scheme has a higher level of noise, which might explain the suboptimal performance. On the other hand, the Q-table with reward shaping displays less noise and a more structured pattern, where the agent follows the spiral outwards, while taking less suboptimal steps toward the goal.

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

The environment consists of the following observation space:

- **direction**: Represents the agent's current direction.
- **image**: A partially observed  $7 \times 7$  grid, where each tile represents an object, which may be "wall", "empty", "ball", or "goal".
- **mission**: Specifies the environment's mission structure (as a string).

In this environment, the agent always starts from the same tile, and the goal is fixed at a specific location. To solve this environment, we apply state aggregation and extract three key features from the observation that we believe are sufficient for reaching the goal:

- **agent\_direction**: The direction the agent is currently facing, represented as an integer in the range  $[0, 3]$ .
- **agent\_pos\_encoding**: The agent's position is inferred from the partially observed grid by counting the number of tiles in front of and to the right of the agent. This is then encoded as a single number in the range  $[0, 8]$ .
- **agent\_path\_clear**: Determines whether the path directly ahead is clear for movement (i.e., no "wall" or "ball"). A value of 0 indicates the path is blocked, while 1 indicates it is clear.

By aggregating states, we reduce the state space to:

$$4 \times 9 \times 2 = 72 \text{ states}$$

This allows for effective generalization over a much larger state space, consisting of all possible locations of obstacles and the agent. Since the goal is always in a fixed position, we do not explicitly compute its location, instead, the agent is left to learn it over episodes. We train the agent over a course of 1000 episodes.

## Part (i): SARSA

Similar to the previous environment, to tune the hyperparameters we ran a parameter sweep consisting of 50 samples using wandb and performed a bayesian search over the search-space of hyperparameters to minimize cumulative regret, from these we selected the top 3 configurations, presented in Table 5.

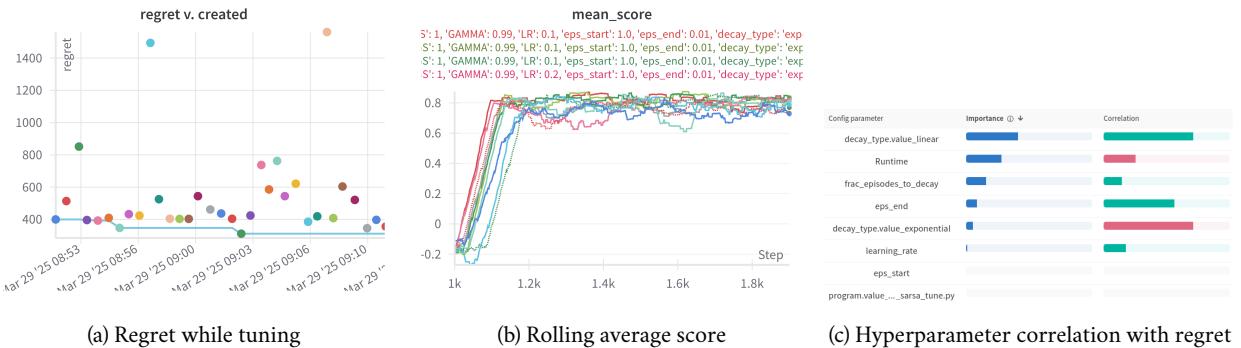


Figure 15: Tuning for Minigrid-SARSA

Hyperparameters	Top	2nd Best	3rd Best
learning_rate	0.1	0.1	0.1
eps_start	1	1	1
eps_end	0.01	0.01	0.01
decay_type	exponential	exponential	exponential
frac_episodes_to_decay	0.1	0.3	1
Cumulative regret	379.268	433.254	854.887
Max rolling average score	0.770	0.836	0.791

Table 5: Top 3 Hyperparameter configurations for Minigrid - SARSA

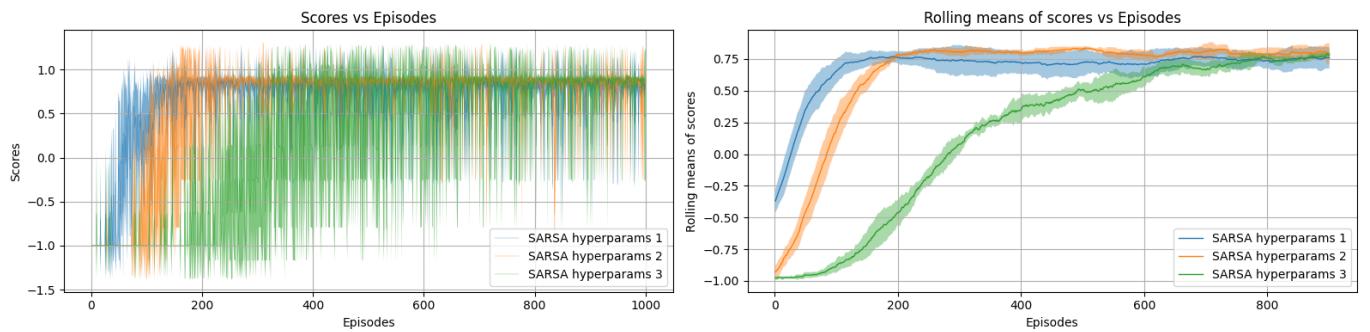


Figure 16: Minigrid SARSA experiments for Top 3 Hyperparameter configurations

## Part (ii): Q-Learning

Similar to the SARSA case, we ran another sweep for Q-Learning and selected the top 3 configurations from 50 samples, see Table 6.

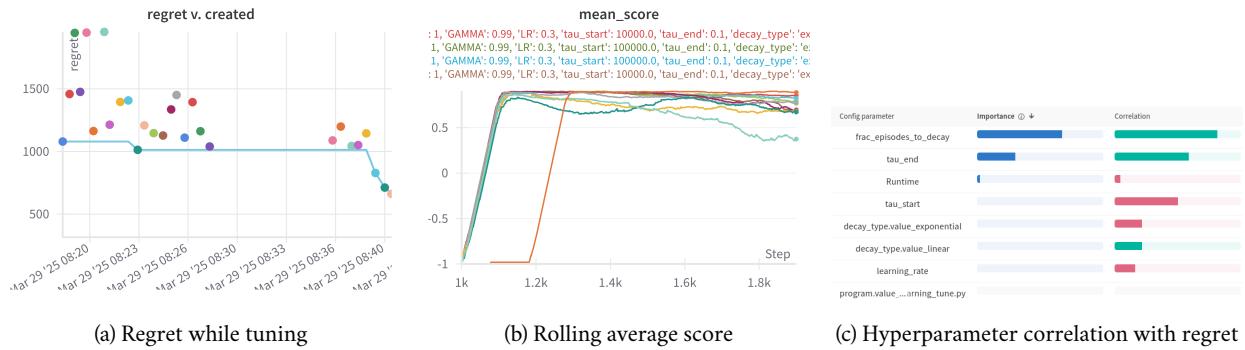


Figure 17: Tuning for Minigrid-Q-Learning

Hyperparameters	Top	2nd Best	3rd Best
learning_rate	0.3	0.3	0.3
tau_start	100000	10000	100000
tau_end	0.1	0.1	0.1
decay_type	exponential	exponential	exponential
frac_episodes_to_decay	0.1	0.1	0.3
Cumulative regret	382.606	412.651	735.958
Max rolling average score	0.872	0.875	0.844

Table 6: Top 3 Hyperparameter configurations for Minigrid - Q-Learning

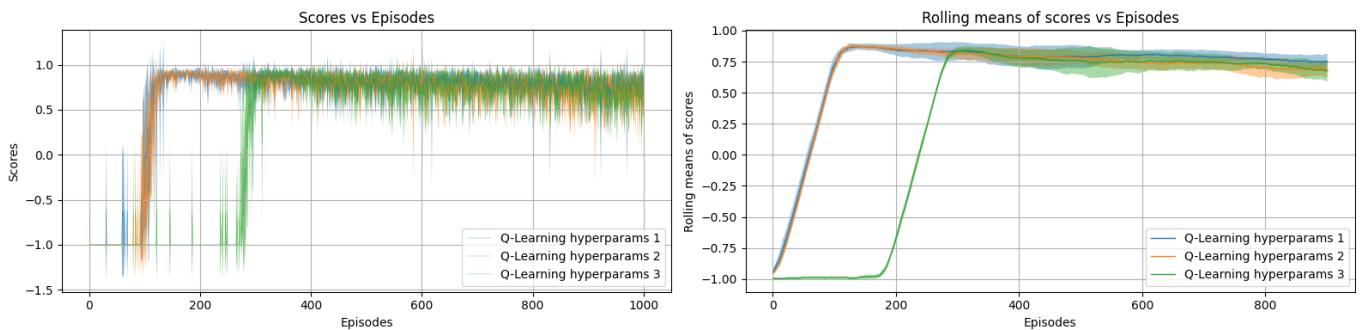


Figure 18: Minigrid Q-Learning experiments for Top 3 Hyperparameter configurations

### Part (iii): Comparison and Inferences

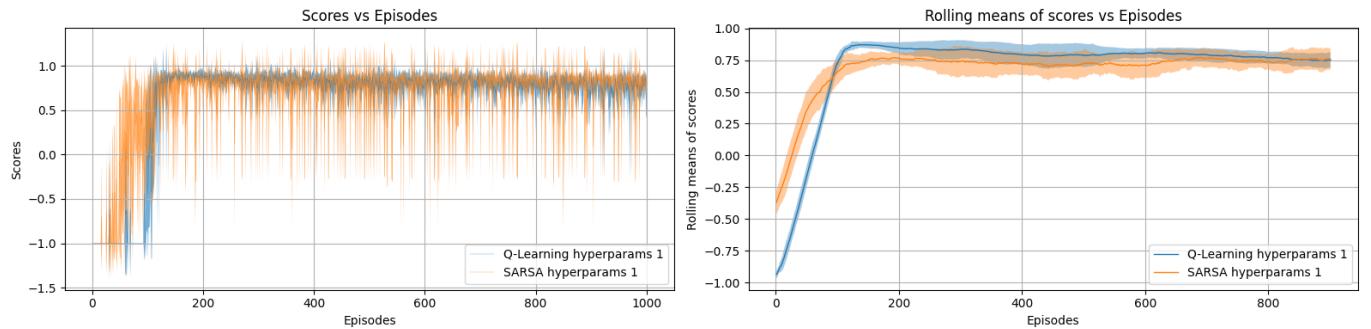


Figure 19: Comparing SARSA and Q-Learning for Top Hyperparameter configurations

It seems clear that the score we can expect is around 0.8. The speed which we can reach it is represented in the plots. The policy doesn't seem too complex to learn, as is evident from the plots and the corresponding `frac_episodes_to_decay`, we need to start exploitation as soon as possible. It is evident from the plot that in the initial stages, Q-Learning takes risky steps (as it evaluates a different policy) and hence converges a bit faster.