

DA6400 : Reinforcement Learning

Programming Assignment #1

Ritabrata Mandal
EE24E009

March 25, 2025

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 above 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 SARSA Code Snippets:

```
1 import numpy as np
2 from q_table import discretize_state
3 from policy import EpsilonGreedyPolicy
4
5 class SarasLearner:
6     """
7         Description: This class implements the SARSA algorithm.
8         Args:
9             alpha      : The learning rate.
10            gamma       : The discount factor.
11            epsilon     : The exploration rate.
12            q_table     : The q-table for the cartpole-v1
13            environment.
14            bins        : The bins for discretizing the state
15            space.
16            env         : The cartpole-v1 environment.
17        """
18    def __init__(self, alpha, gamma, epsilon, q_table, bins,
19                env, seed):
20        self.alpha = alpha
21        self.gamma = gamma
22        self.epsilon = epsilon
23        self.q_table = q_table
24        self.env = env
25        self.bins = bins
26        self.seed = seed
27        self.policy = EpsilonGreedyPolicy(self.epsilon, self.
28            q_table, self.env)
29
30    def compute_td_error(self, state, action, next_state,
31                        next_action, reward):
32        """
33            Description: This function computes the TD error.
34            Args:
35                state      : The current state.
36                action     : The current action.
37                next_state  : The next state.
38                reward     : The reward.
39            Returns:
40                td_error   : The TD error.
41        """
42        return reward + self.gamma * self.q_table[next_state
43            [0], next_state[1], next_state[2], next_state[3],
```

```

38         next_action] - \
           self.q_table[state[0], state[1], state[2], state[3],
           action]
39
40     def update_q_table(self, state, action, td_error):
41         """
42             Description: This function updates the q-table.
43             Args:
44                 state      : The current state.
45                 action     : The current action.
46                 td_error   : The TD error.
47             Equations:
48                  $Q(s, a) \leftarrow Q(s, a) + \alpha * (td\_error)$ 
49         """
50         self.q_table[state[0], state[1], state[2], state[3],
51         action] += self.alpha * td_error
52
53     def learn(self, num_episodes, num_steps):
54         reward_list = []
55         for episode in range(num_episodes):
56             state, _ = self.env.reset(seed=self.seed)
57             state_discrete = discretize_state(state, self.bins)
58             action = self.policy.get_action(state_discrete) #
59             # SARSA selects initial action
60             total_reward = 0
61
62             for step in range(num_steps):
63                 next_state, reward, done = self.env.step(action
64                 )[:3]
65                 next_state_discrete = discretize_state(
66                 next_state, self.bins)
67                 next_action = self.policy.get_action(
68                 next_state_discrete) # Select next action
69
70                 td_error = self.compute_td_error(state_discrete
71                 , action, next_state_discrete, next_action, reward)
72                 self.update_q_table(state_discrete, action,
73                 td_error)
74
75                 state_discrete, action = next_state_discrete,
76                 next_action # SARSA moves (s, a) -> (s', a')
77                 total_reward += reward
78
79                 if done:
80                     break
81
82             print(f"Episode: {episode + 1}/{num_episodes},
83             Total Reward: {total_reward}")
84             reward_list.append(total_reward)

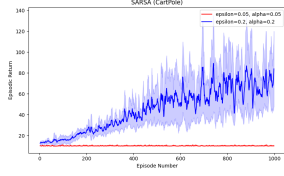
```

```
76  
77     return reward_list
```

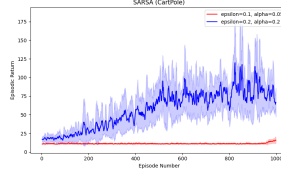
Listing 1: SARSA - Agent

2.1.2 SARSA Runs

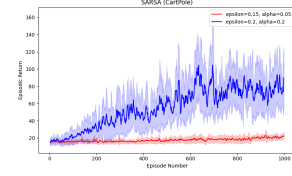
below are the graphs of runs of SARSA with different hyper-parameter



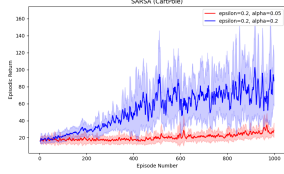
(a) $\alpha = 0.05$



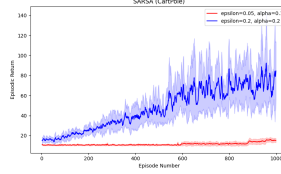
(b) $\epsilon = 0.1$



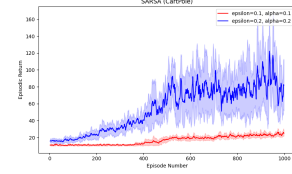
(c) $\epsilon = 0.15$



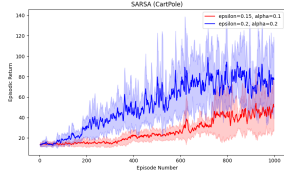
(d) $\epsilon = 0.2$



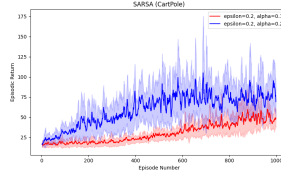
(e) $\epsilon = 0.25$



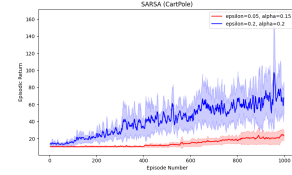
(f) $\epsilon = 0.3$



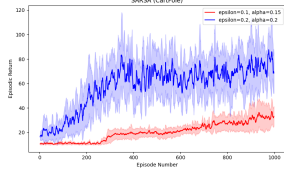
(g) $\epsilon = 0.35$



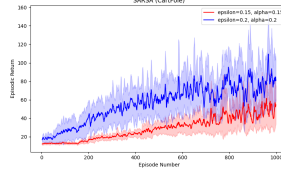
(h) $\epsilon = 0.4$



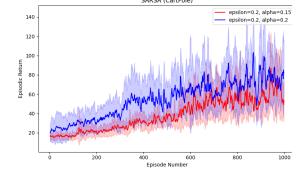
(i) $\epsilon = 0.45$



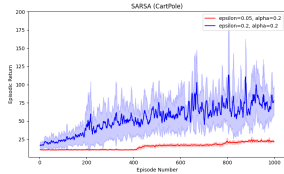
(j) $\epsilon = 0.5$



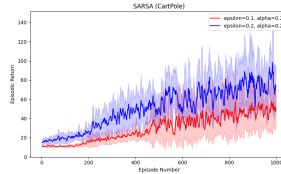
(k) $\epsilon = 0.55$



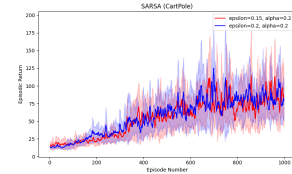
(l) $\epsilon = 0.6$



(m) $\epsilon = 0.65$



(n) $\epsilon = 0.7$



(o) $\epsilon = 0.75$

Figure 1: Phase portraits for different values of ϵ