

CS6700: Reinforcement Learning

Programming Assignment - II Report

HARSHAVARDHAN SRINIVAS PENTAKOTA [ME20B083]

PUTREVU SAI SATWIK [ME20B138]
[GITHUB LINK.](#)

April 22, 2024

Table of Contents

1 Settings of Environment	2
2 SMDP Q-Learning:	4
3 IOQL Q-Learning:	9
4 Inferences:	15
4.1 Reasons for the Working of Hierarchical Learning Algorithms:	15
5 Comparsion of IOQL vs SMDP:	16
6 New Set of Options:	17

1 Settings of Environment

The environment for this task is the taxi domain.

Goal: To pick up the passenger and Drop him at the Destination.

Primitive Actions:

- 0: move south
- 1: move north
- 2: move east
- 3: move west
- 4: pickup passenger
- 5: drop off passenger

No of Passenger locations: 5

No of Destinations: 4

Map: 5 x 5 Grid

So for each (Passanger location, Destination) state we have another 25 possible states

Options: Options to move the taxi to each of the four designated locations, executable when the taxi is not already there.

So, No of options = 4;

Options: The taxi goes to a specific goal location (R, G, Y, B). After reaching there, it has the option of picking up or dropping the passenger. If it's not at the goal location, it uses epsilon-greedy on option qvalues to choose any primitive action (up,down,right,left).

Q values of SMDP or IOQL is 20 sub-states x N options

Code snippets for Option Definition:

```

1  '''
2  We have 4 goals and an option to go to each goal.
3  '''
4
5  num_options = 4
6
7  # Position of goals
8  goal_pos = env.unwrapped.locs
9
10
11  '''
12  Goal is R, G, Y, B
13  '''
14
15  def Deliver_policy(Q_option, goal_pos, goal, state, epsilon):
16      # Decode the state into taxi position, passenger location, and drop
17      ↪ location
18      taxi_X, taxi_Y, Passenger, DropLoc = env.decode(state)
19
20      # Initialize the option termination flag
21      optdone = False
22
23      # Check if the taxi is at the goal location
24      if (taxi_X == goal_pos[goal][0] and taxi_Y == goal_pos[goal][1]):
25          optdone = True
26
27      # If the passenger is at the goal location, pick them up
28      if (Passenger == goal):
29          optact = 4 # Pick up the passenger at the goal
30      # If the drop location is at the goal location, drop the
31      ↪ passenger
32      elif (DropLoc == goal):
33          optact = 5 # Drop the passenger at the goal
34      else: # If it's just a (R, G, Y, B) location but not a pickup
35          ↪ or drop location
36          # Choose the action based on epsilon-greedy policy
37          optact = epsilon_policy(Q_option[goal], taxi_X * 5 + taxi_Y,
38          ↪ epsilon)
39
40      else:
41          # If not at the goal location, choose the action based on
42          ↪ epsilon-greedy policy
43          optact = epsilon_policy(Q_option[goal], taxi_X * 5 + taxi_Y,
44          ↪ epsilon)
45
46      return [optact, optdone]

```

2 SMDP Q-Learning:

For SMDP Q learning, we have 4 options as we mentioned before.

Each Option consists of Q values to decide the primitive action needed to move towards the Goal for each given substate (Grid Position of Taxi[25 states])

Q values option = 4 options * (25 taxi loc * 4 Primitive actions)

```

1 num_of_options = 4
2 N_options = 4
3 N_goals = 4
4 N_passenger_location = 5
5 N_row = 5
6 N_col = 5
7 Q_values_SMDP = np.ones((N_passenger_location*N_goals,N_options))
8 Q_values_options = np.ones((N_goals,5*5,no_of_actions-2))
9 alpha = 0.4
10 gamma = 0.99

```

SMDP ALGO:

```

1 def SMDP(seed, Deliver_policy=Deliver_policy):
2     # Load the environment and set the state and action spaces
3     env, state_space, action_space = LoadingEnv(seed=seed)
4
5     # Initialize rewards and other variables
6     rewards = []
7     episodes = 1000
8     eps = 0.1
9     eps_options = {i: 0.1 for i in range(N_options)} # Epsilon for
10    ↪ each option
11    count_success = 0
12
13    # Loop through episodes
14    for i in tqdm(range(episodes)):
15        state = env.reset()
16        done = False
17        total_reward = 0
18
19        # Loop until the episode is done
20        while not done:
21            # Decode the current state to get taxi location, passenger
22            ↪ location, and drop location
23            taxi_X, taxi_Y, Passenger, DropLoc = env.decode(state)
24
25            # Determine the sub-state for transforming 500 states to
26            ↪ fit in N_passenger_location * N_goals

```

```

24     sub_state = Passenger * N_goals + DropLoc
25
26     # Choose an option using epsilon-greedy policy
27     option = epsilon_policy(Q_values_SMDP, sub_state,
        ↪     epsilon=eps)
28
29     # Initialize variables for option execution
30     reward_bar = 0
31     opt_done = False
32     steps = 0
33     prev_state = state
34
35     # Execute the option until it is done or the episode ends
36     while not opt_done and not done:
37         # Choose an action for the option
38         Taxi_curr_X, Taxi_curr_Y, _, _ = env.decode(state)
39         opt_act, opt_done = Deliver_policy(Q_values_options,
        ↪         goal_pos, option, state, eps_options[option])
40
41         # Execute the action and observe the next state and
        ↪ reward
42         next_state, reward, done, _ = env.step(opt_act)
43         Taxi_next_X, Taxi_next_Y, _, _ = env.decode(next_state)
44         # Update the surrogate reward based on whether the
        ↪ option is completed
45         if opt_done:
46             reward_surr = 20 # Surrogate reward for completing
        ↪ the option
47         else:
48             reward_surr = reward
49
50         # Update Q-values if the action is primitive
51         if opt_act < 4:
52             Q_values_options[option][5 * Taxi_curr_X +
        ↪             Taxi_curr_Y][opt_act] += alpha * (
53                 reward_surr + gamma *
        ↪             np.max(Q_values_options[option][5 *
        ↪             Taxi_next_X + Taxi_next_Y, :]) -
54             Q_values_options[option][5 * Taxi_curr_X +
        ↪             Taxi_curr_Y, opt_act])
55
56         # Update variables
57         steps += 1
58         reward_bar = gamma * reward_bar + reward
59         total_reward += reward
60         state = next_state
61
62         # Update SMDP Q-value
63         _, _, passenger, DropLoc = env.decode(state)

```

```

64     sub_state = N_goals * passenger + DropLoc
65
66     _, _, passenger, DropLoc = env.decode(prev_state)
67     prev_sub_state = N_goals * passenger + DropLoc
68
69     Q_values_SMDP[prev_sub_state, option] += alpha * (
70         reward_bar + (gamma ** steps) *
71         ↪ np.max(Q_values_SMDP[sub_state, :]) -
72         Q_values_SMDP[prev_sub_state, option])
73
74     # Append total reward for the episode
75     rewards.append(total_reward)
76
77     return rewards

```

SMDP Reward Plot:

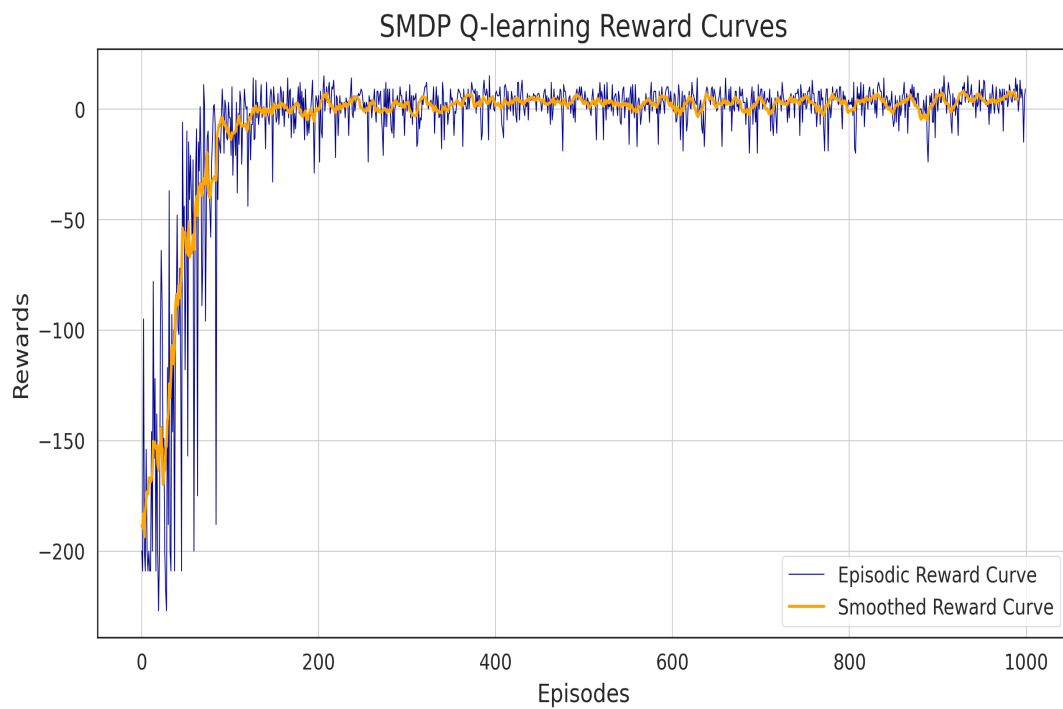


Figure 1: SMDP Q learning Reward Curve

Q values option[Move Red (option 0) is 25 x 4 Qmap denoting the Qmap to choose primitive action to move the taxi towards that option goal.

Visulaization of Polciy Qmap:

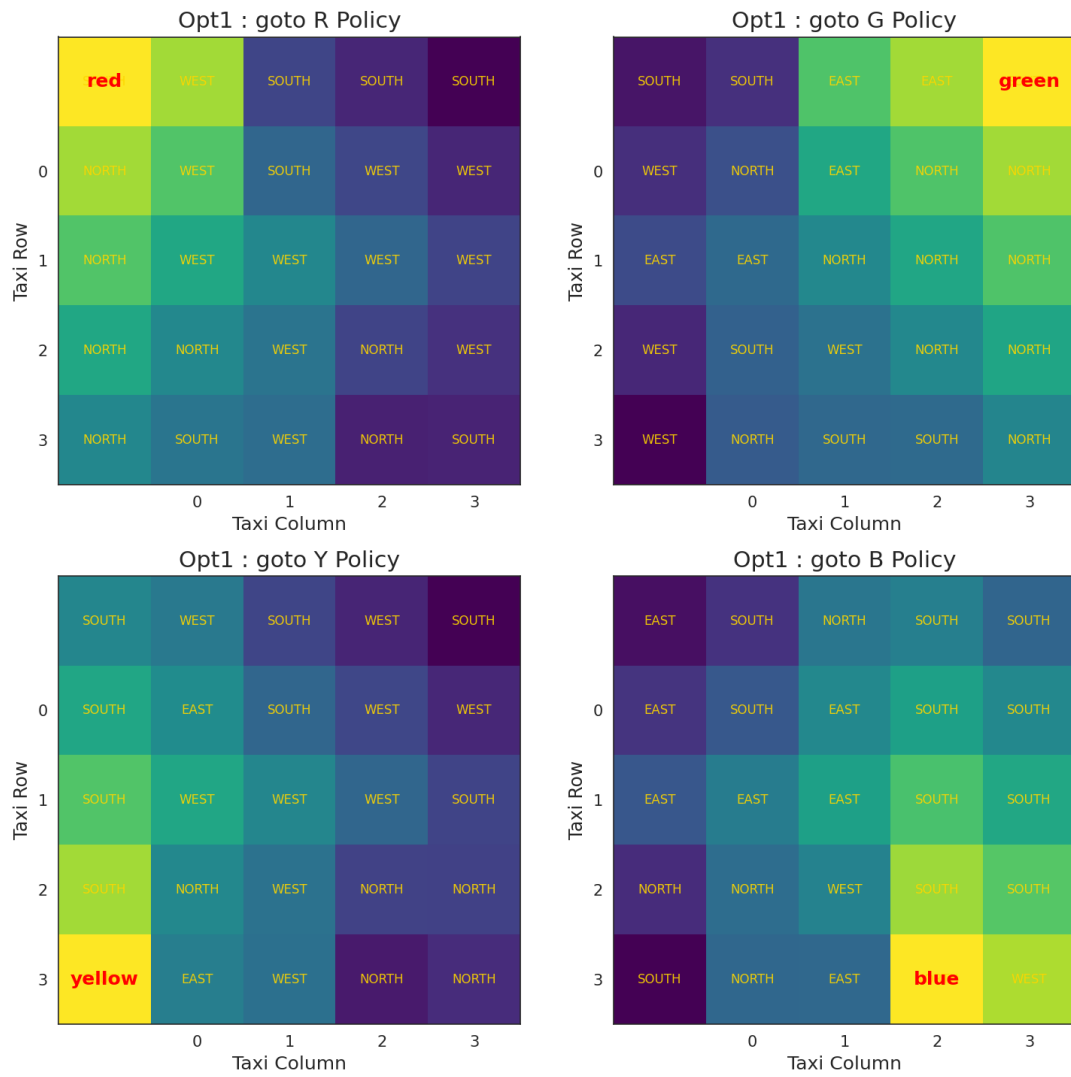


Figure 2: Qmap of each policy

Q values of SMDP or IOQL is 20 sub-states x N options. For each sub-state we choose which one of the 4 options is best.

Visulaization of SMDP Polciy Qmap:

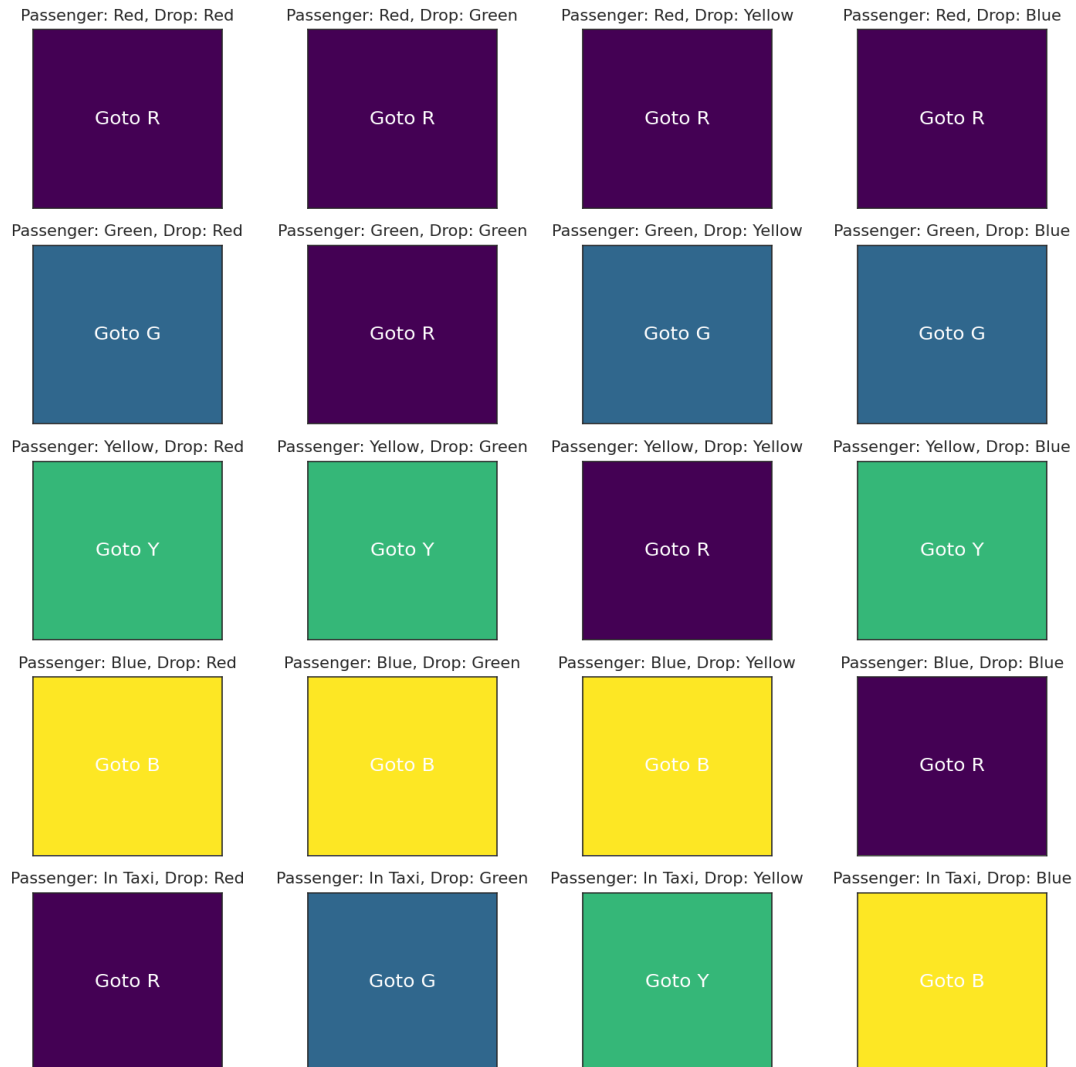


Figure 3: SMDP Qmap

3 IOQL Q-Learning:

For IOQL Q learning, we have 4 options as we mentioned before.

Each Option consists of Q values to decide the primitive action needed to move towards the Goal for each given substrate (Grid Position of Taxi[25 states])

Q values option = 4 options * (25 taxi loc * 4 Primitive actions)

```

1 num_of_options = 4
2 N_options = 4
3 N_goals = 4
4 N_passenger_location = 5
5 N_row = 5
6 N_col = 5
7 Q_values_IOQL = np.ones((N_passenger_location*N_goals,N_options))
8 Q_values_options = np.ones((N_goals,5*5,no_of_actions-2))
9 alpha = 0.1
10 gamma = 0.99

```

IOQL ALGO:

```

1 def IOQL(seed, Deliver_policy = Deliver_policy):
2     # Load the environment and set the state and action spaces
3     env, state_space, action_space = LoadingEnv(seed=seed)
4
5     # Initialize rewards and other variables
6     rewards = []
7     episodes = 1000
8     eps = 0.1
9     eps_options = {i: 0.1 for i in range(N_options)} # Epsilon for
10    ↪ each option
11    count_success = 0
12
13    # Loop through episodes
14    for i in tqdm(range(episodes)):
15        state = env.reset()
16        done = False
17        total_reward = 0
18
19        # Loop until the episode is done
20        while not done:
21            # Decode the current state to get taxi location, passenger
22            ↪ location, and drop location
23            taxi_X, taxi_Y, Passenger, DropLoc = env.decode(state)
24
25            # Calculate the sub-state to fit in N_passenger_location *
26            ↪ N_goals by considering only passenger_loc, DropLoc

```

```

24     sub_state = Passenger * N_goals + DropLoc
25
26     # Choose an option using epsilon-greedy policy
27     option = epsilon_policy(Q_values_IOQL, sub_state,
        ↪     epsilon=eps)
28
29     # Initialize variables for option execution
30     opt_done = False
31     steps = 0
32     prev_state = state
33
34     # Execute the option until it is done or the episode ends
35     while not opt_done and not done:
36         # Decode the current state to get the taxi location,
        ↪ passenger location, and drop location
37         Taxi_curr_X, Taxi_curr_Y, passenger, DropLoc =
        ↪     env.decode(state)
38
39         # Calculate the sub-state
40         sub_state = N_goals * passenger + DropLoc
41
42         # Choose an action for the option
43         opt_act, opt_done = Deliver_policy(Q_values_options,
        ↪     goal_pos, option, state, eps_options[option])
44
45         # Execute the action and observe the next state and
        ↪ reward
46         next_state, reward, done, _ = env.step(opt_act)
47
48         # Decode the next state to get the taxi location,
        ↪ passenger location, and drop location
49         Taxi_next_X, Taxi_next_Y, passenger1, DropLoc1 =
        ↪     env.decode(next_state)
50
51         # Calculate the next sub-state
52         next_sub_state = N_goals * passenger1 + DropLoc1
53
54         # Update the surrogate reward based on whether the
        ↪ option is completed
55         if opt_done:
56             reward_surr = 100 # Surrogate reward for
        ↪ completing the option
57         else:
58             reward_surr = reward
59
60         # Update Q-values if the action is primitive
61         if opt_act < 4:
62             Q_values_options[option][5 * Taxi_curr_X +
        ↪     Taxi_curr_Y][opt_act] += alpha * (reward_surr +
        ↪     gamma * np.max(Q_values_options[option][5 *
        ↪     Taxi_next_X + Taxi_next_Y])
        ↪     -Q_values_options[option][5 * Taxi_curr_X +
        ↪     Taxi_curr_Y, opt_act])

```

```

63         # Update total reward
64         total_reward += reward
65
66         # Update Q-values for each option
67         for O in range(N_options):
68             opt_act_option, opt_done_option =
69                 ↪ Deliver_policy(Q_values_options, goal_pos, O,
70                 ↪ state, eps_options[O])
71             if opt_act_option == opt_act:
72                 if opt_done_option:
73                     Q_values_IOQL[sub_state, O] += alpha *
74                     ↪ (reward + gamma *
75                     ↪ np.max(Q_values_IOQL[next_sub_state]) -
76                     ↪ Q_values_IOQL[sub_state, O])
77                 else:
78                     Q_values_IOQL[sub_state, O] += alpha * (
79                     ↪ reward + gamma *
80                     ↪ (Q_values_IOQL[next_sub_state, O])
81                     ↪ -Q_values_IOQL[sub_state, O])
82
83         # Update the current state
84         state = next_state
85
86         # Append total reward for the episode
87         rewards.append(total_reward)
88
89     return rewards

```

IOQL Reward Plot:

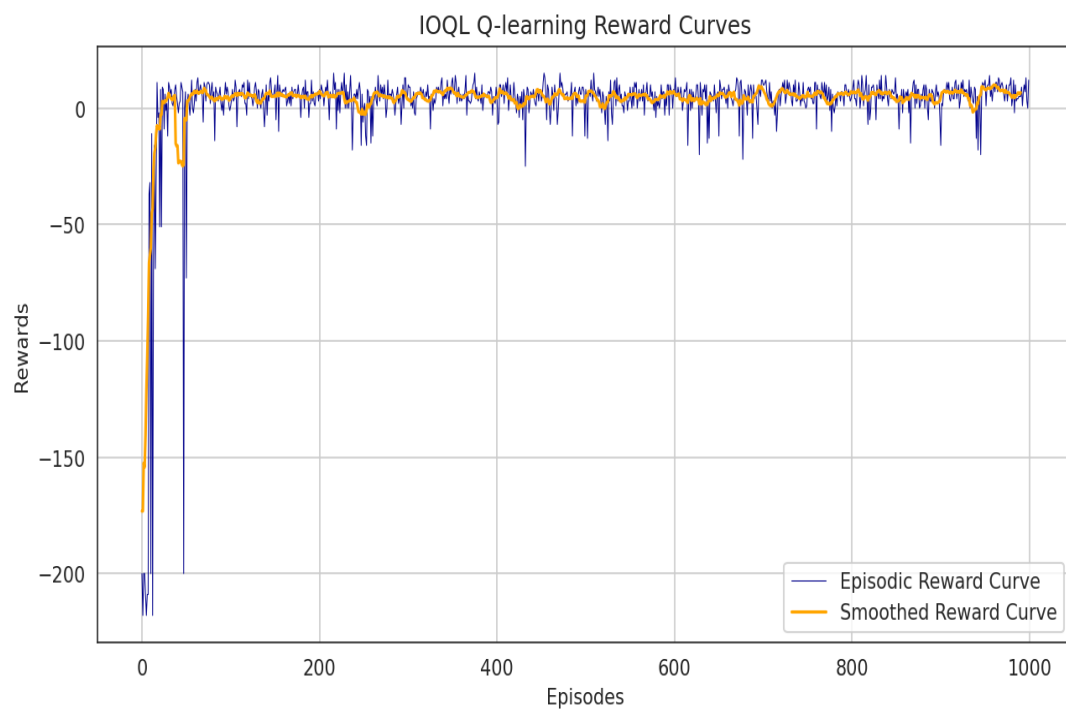


Figure 4: IOQL Q learning Reward Curve

Q values option[Move Red (option 0) is 25 x 4 Qmap denoting the Qmap to choose primitive action to move the taxi towards that option goal.

Visulaization of Polciy Qmap -IOQL:

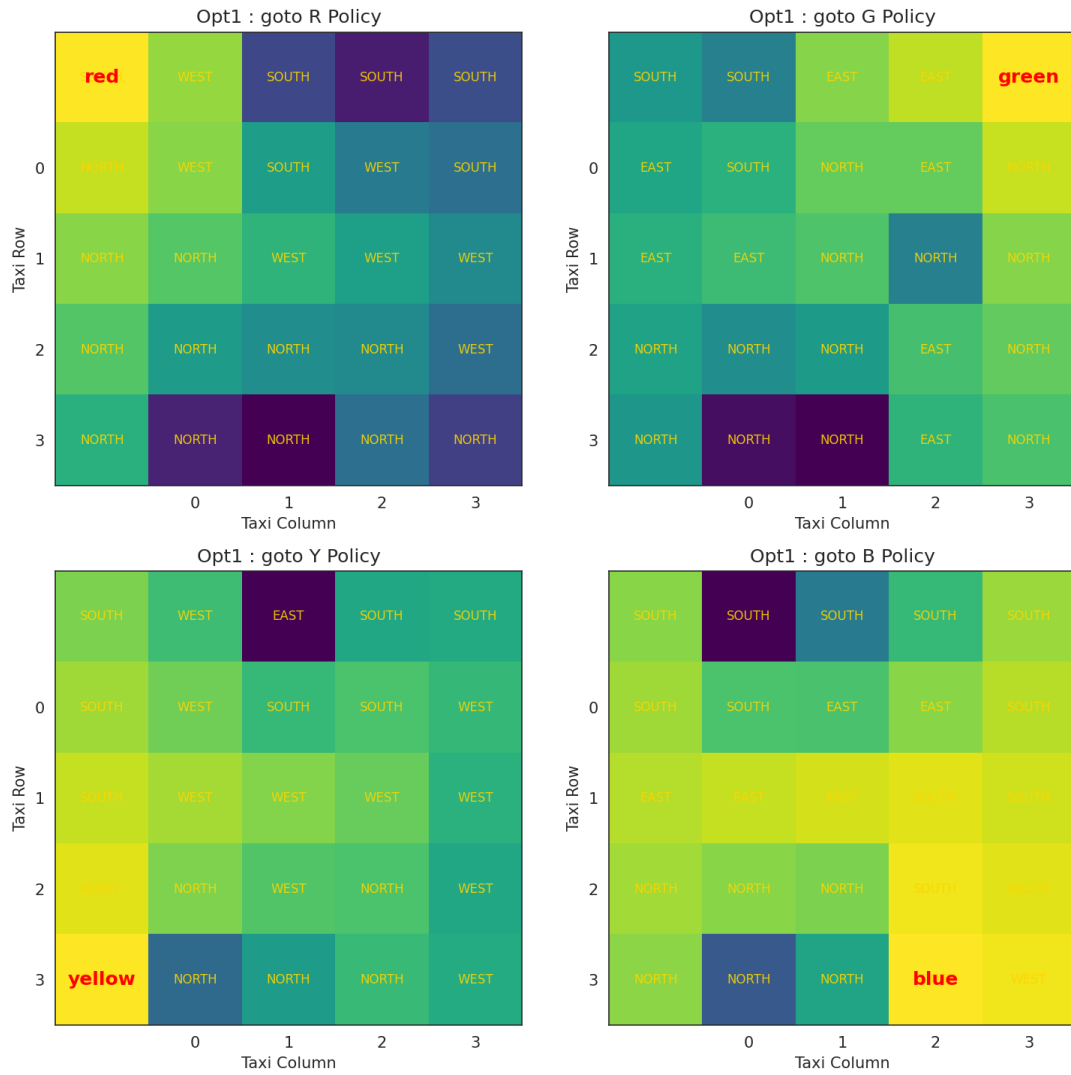


Figure 5: Qmap of each policy

For each sub-state we choose which one of the 4 options is best.

Visulaization of IOQL Polciy Qmap:



Figure 6: IOQL Qmap

4 Inferences:

4.1 Reasons for the Working of Hierarchical Learning Algorithms:

Hierarchical learning enables agents to think in terms of larger actions, such as "go to the pickup location in Red" or "drop off the passenger in Blue," rather than focusing on individual steps. By considering these higher-level actions, agents can learn and plan more efficiently, ultimately leading to better performance in complex tasks like taxi navigation.

SMDP Q-learning Policy Reasoning:

SMDP Q-learning involves two levels of policies: the main policy for choosing options and the option policies for completing the option goal.

Main Policy to Choose Option:

1. Figure 3 illustrates the decoding of 500 states into 20 sub-states, representing pickup and drop locations.
2. Each sub-state corresponds to selecting an option to reach the goal, simplifying the learning process by focusing on a single action rather than multiple movements.
3. Figure 3 depicts the optimal option for each sub-state.

Option Policy:

1. Within the main policy, option policies dictate primitive actions to complete a taxi maneuver.
2. Figure 2 displays the Q-map indicating the shortest path for each option goal.

SMDP Q-learning learns this policy by allowing the agent to operate at a higher level, focusing on actions like "go to location Y" or "go to location Blue." By simplifying the decision-making process to selecting one option for each sub-state, the agent can efficiently achieve its objectives.

IOQL (Intra-Option Q-learning) Policy:

The main and option policies used in IOQL are analogous to those in SMDP (see Figure 6 and Figure 5).

IOQL optimizes Q-values within each option using an intra-option learning approach. By emphasizing the value of staying within an option until completion, IOQL enhances decision-making efficiency. Incorporating 20 sub-states enables IOQL to effectively explore and exploit the environment, resulting in improved performance compared to traditional Q-learning methods.

5 Comparsion of IOQL vs SMDP:

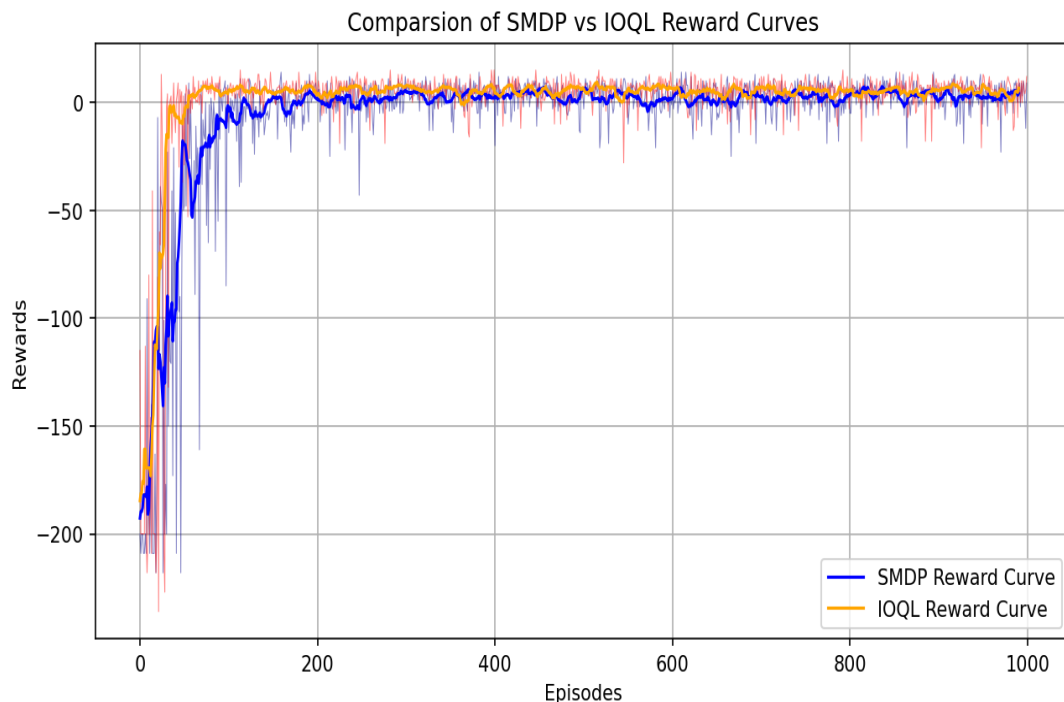


Figure 7: IOQL vs SMDP Reward Curve

IOQL performs better than SMDP Q learning (see Fig 7)

IOQL performs better because it allows for simultaneous updates across multiple options for a single action. This means that when taking a certain action, IOQL can update the Q-values for multiple options if those options were followed. This simultaneous updating enables IOQL to make more efficient use of experience, leading to faster learning and improved performance compared to SMDP Q-learning.

6 New Set of Options:

The set of options are taken such that the taxi tries to move towards the corners of the grid provided the taxi is already not present. So the **total number of actions taken are 8 including the 4 primitive actions (UP, DOWN, RIGHT, LEFT).**

No of Options : 4

Total No of options : 8 (4 Primitive actions + 4 options

Corners : (0, 0), (0, 4), (4, 0), (4, 4)

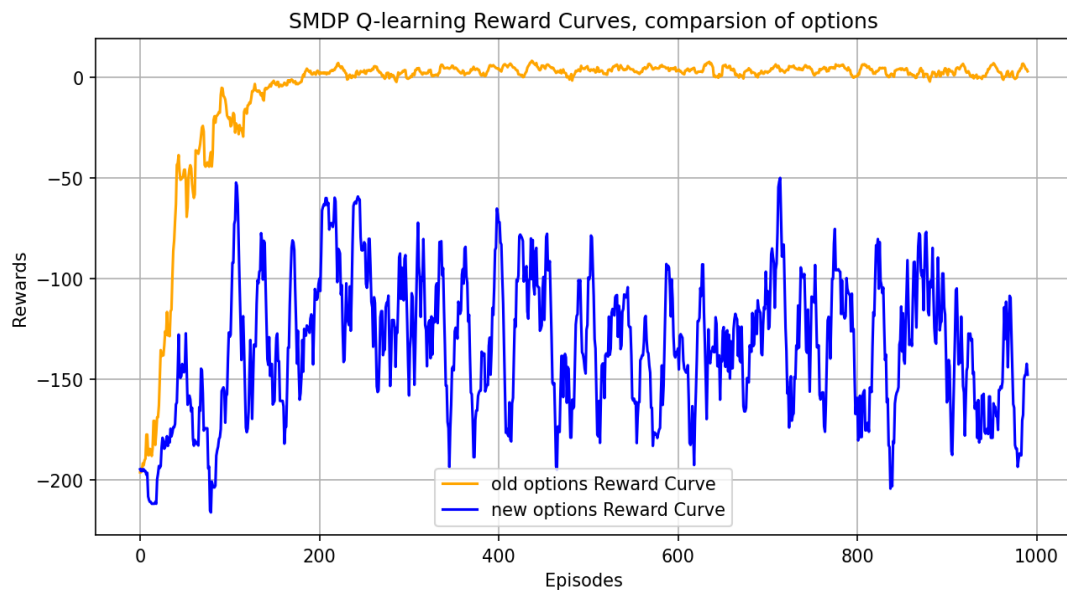
New Option Policy Code:

```

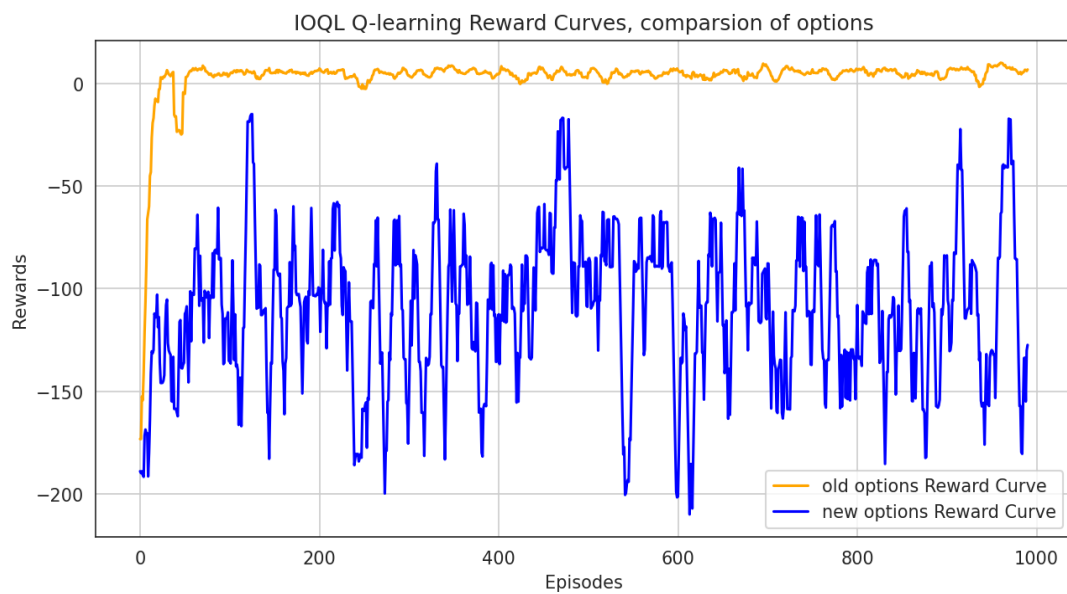
1  num_options = 4
2
3  goal_pos = [(0,0), (0,1),(2,0), (4,4)]
4
5  def new_policy(Q_option, goal_pos, goal, state,epsilon):
6
7      taxi_X, taxi_Y, Passenger, DropLoc = env.decode(state)
8      optdone = False
9      if(goal >3): # 4, 5, 6, 7 are chosen for UP, DOWN, RIGHT, and LEFT
10         ↪ ↪ respectively as the input
11         optact = goal - 4
12         optdone = True
13         return [optact,optdone]
14     if (taxi_X == goal_pos[goal][0] and taxi_Y == goal_pos[goal][1]):
15         optdone = True
16         if (Passenger == goal):
17             optact = 4 #pick the passenger at goal
18         elif (DropLoc == goal):
19             optact = 5 # Drop the passenger at goal
20         else:
21             optact = epsilon_policy(Q_option[goal],taxi_X*5+taxi_Y,epsilon)
22     else:
23         optact = epsilon_policy(Q_option[goal],taxi_X*5+taxi_Y,epsilon) ##
24         ↪ ↪ epsilon_greedy(Q_option[goal][state]) will choose one among UP,
25         ↪ ↪ DOWN, RIGHT, LEFT
26
27     return [optact,optdone]

```

Comparision of Old options vs New options for SMDP-QLearning



Comparision of Old options vs New options for IOQL-QLearning



It is evident from both plots that old options are better

GITHUB REPO LINK:

GITHUB USERNAME:

ME20B083 : harshavardhan379

ME20B138 : SATWIK-ME20B138

GITHUB REPO LINK : https://github.com/harshavardhan379/CS6700_P43.