

Exercise 3

Amr Alkhashab - 905833
ELEC-E8125 - Reinforcement Learning

October 5, 2020

1 Task 1

Two functions were created for this task. First the Qlearning (The Cartpole code is presented in qlearning.py):

```
1 def update_q_value(old_state, action, new_state, reward, done,
2   q_array):
3     # TODO: Implement Q-value update
4     if done == False: #for non terminating next state
5         old_cell_index = get_cell_index(old_state)
6         new_cell_index = get_cell_index(new_state)
7         q_old = q_array[old_cell_index[0], old_cell_index[1],
8             old_cell_index[2], old_cell_index[3], action]
9         q_max = np.max(q_array[new_cell_index[0],
10             new_cell_index[1], new_cell_index[2], new_cell_index
11             [3], :])
12         q_old = q_old + alpha * (reward + (gamma * q_max) -
13             q_old)
14         q_array[old_cell_index[0], old_cell_index[1],
15             old_cell_index[2], old_cell_index[3], action] =
16             q_old
17     else: #for terminating next state we set next state state-
18         action value to zero.
19         old_cell_index = get_cell_index(old_state)
20         q_old = q_array[old_cell_index[0], old_cell_index[1],
21             old_cell_index[2], old_cell_index[3], action]
22         q_old = q_old + alpha * (reward - q_old)
23         q_array[old_cell_index[0], old_cell_index[1],
24             old_cell_index[2], old_cell_index[3], action] =
25             q_old
```

The second function is for action selection based on epsilon:

```
1 def get_action(state, q_values, greedy=False):
2     #TODO: Implement epsilon-greedy
3     cell_index = get_cell_index(state)
4     greedyaction = -1
5     min = float('-inf')
6
7     #Value a is determined by round((target_eps * 20000)/(1 -
8       target_eps)), while K is incremented inside the episode
9       loop.
10    #epsilon = a / (a + K)
11    #epsilon = 0
12
13    # Find greedy action
14    for i in range(num_of_actions):
15        if q_grid[cell_index[0], cell_index[1], cell_index[2],
16            cell_index[3], i] > min:
17            min = q_grid[cell_index[0], cell_index[1],
18                cell_index[2], cell_index[3], i]
19            greedyaction = i
20
21    #select between all action with probability epsilon
22    if np.random.random() < epsilon:
23        action = np.random.randint(0, num_of_actions)
24    else:
25        action = greedyaction
26
27    return int(action)
```

The result of constant 0.2:

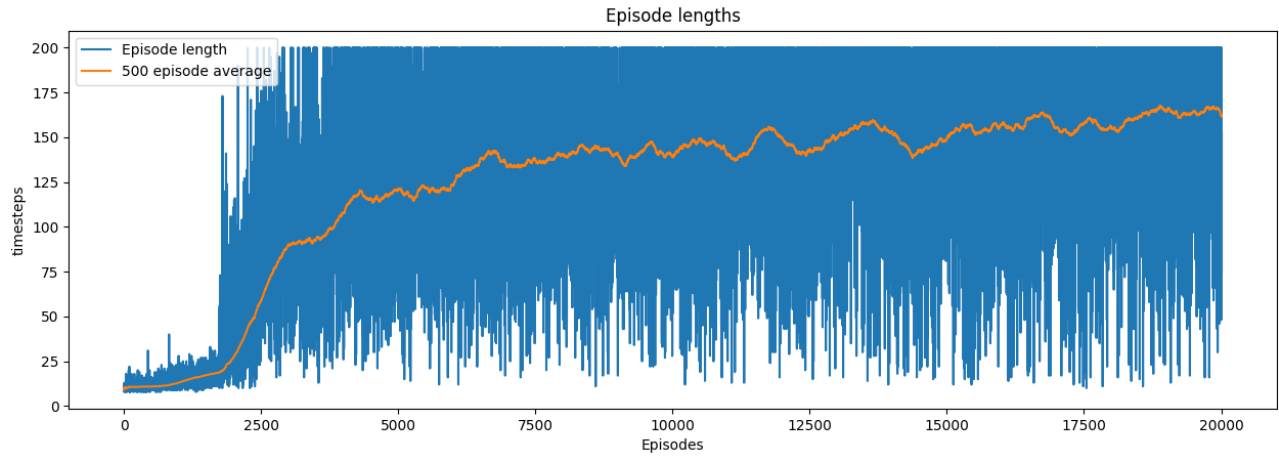


Figure 1: Training Cartpole with constant Epsilon 0.2

The result of constant GLIE:

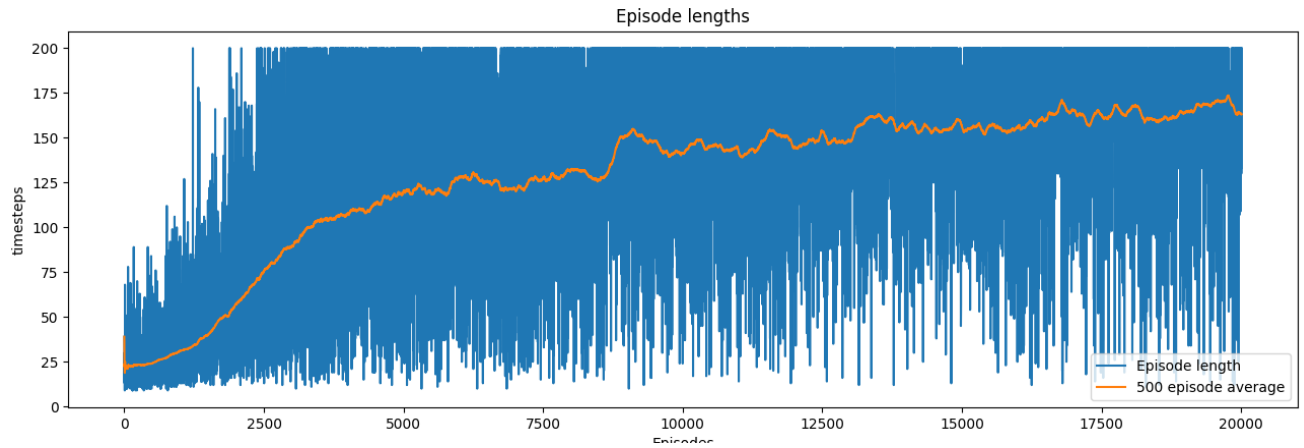


Figure 2: Training Cartpole with GLIE

2 Task 2

The Optimal value function can be easily determine by taking the maximum of state-actions value for each state.

```
1 values = np.max(q_grid, axis=4)
```

For the heatmap a function is created to draw the x and theta, by taking avearge over velocity and angular velocity:

```
1 #heatmap drawing
2 def heatmap(q_grid):
3     value= np.zeros(q_grid.shape[:-1]) # TODO: COMPUTE THE
        VALUE FUNCTION FROM THE Q-GRID
4     #value_function
5     value = np.max(q_grid, axis=4)
6     values_array = np.zeros(q_grid.shape[:-3])
7     #values_array is 2d array with mean values over velocity and
        angular velocity for each position and theta
8     values_array = np.mean(np.mean(value, axis = 3), axis = 1)
9     ax = sns.heatmap(values_array)
10    plt.xlabel("X")
11    plt.ylabel("Theta")
12    plt.show()
```

The heatmap after the training looks like this:

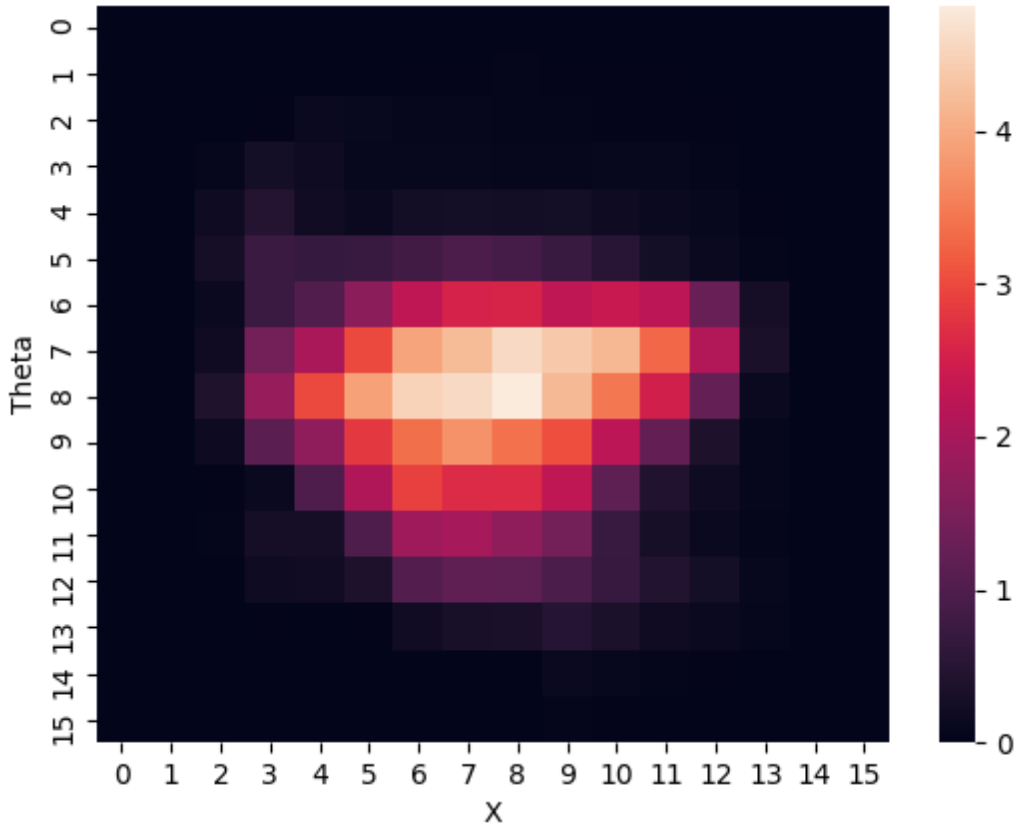


Figure 3: Heatmap after training

3 Question 1

The three plot will be summited with the files, since attaching to report is not required.

3.1 (a)

Before training, the heatmap will all be zeros,since no training has taken place and the values will be the intial intialized values (zeros).

3.2 (b)

After single episode, the area of change is still very small from observation it is between X of grid (7 and 12) and theta grid of around (7 to 9), it is still far away from Figure. 3. The value around the center will change very slightly (very small change) toward the optimal value, due to small step size of 0.1 and the small covearage area is because the agent still didnt explore most of the other states yet and thus the avearge is nearly zero.

3.3 (c)

Halfway through, It will look more like Figure. 3, the agent was able to explore more states and due to multiple episode the values converged more and more toward the true value, as they have updated more than once, and more states has been visited.

4 Task 3

(a)- Running the training for initial 0 and episolon equals 0: (b)- Running the training for

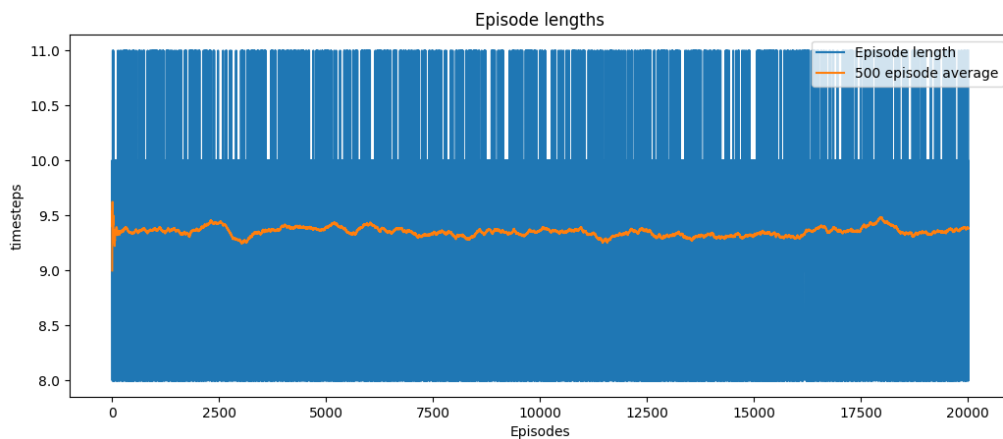


Figure 4: Training for Epislon = 0 and initial $Q = 0$

initial 50 and episolon equals 0:

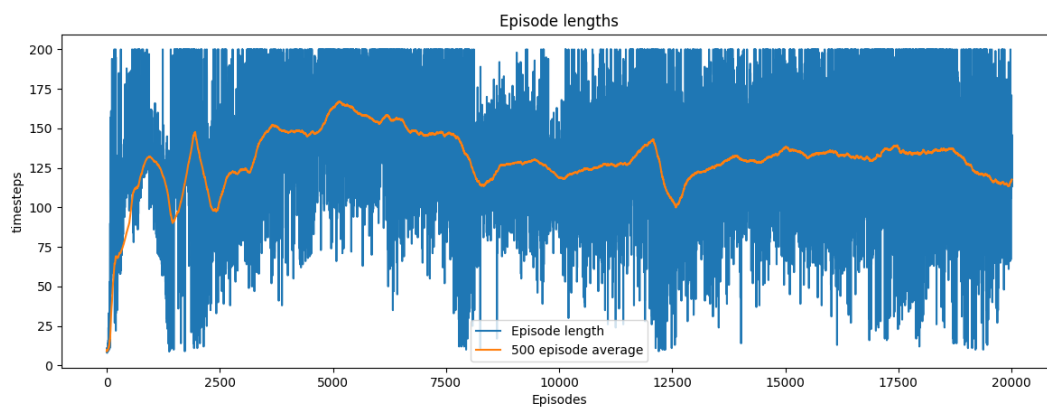


Figure 5: Training for Epislon = 0 and initial $Q = 50$

5 Question 2

5.1 Question 2.1

The model perform better in the second case, where the initial estimate is set to 50.

5.2 Question 2.2

For the first case, since all values are zero which is lower than the optimal value for each state, no exploration will take place and the only greedy action will be selected. The method of the second case is called Optimistic Initial Values. Since the true optimal state-value function is less than the initial value, every time an action is selected the estimate decreases, Thus, all state-action value pair are selected atleast once and tried several time before it converges. Thus, The agent does fair amount of exploration even if greedy actions are selected most of the time. As long the initial values are larger than the optimal values for each state, there will be optimistic exploration, the larger the value, the more the exploration.

6 Task 4

The code is done in new python file (qlarning2). The grid size of the first continous parameter are selected as 16 for each, similar to previous Task1. The results:

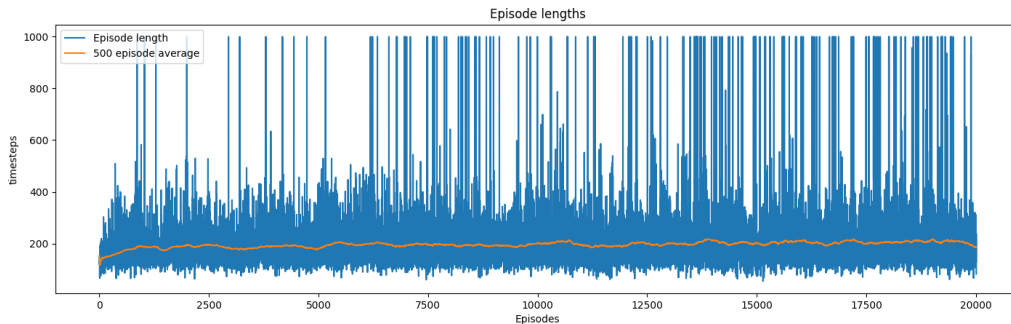


Figure 6: Lunar lander training

7 Question 3.1

The lander doesnot learn much. This is due to two reasons:

1. The range of the observation is much larger than the previous Tasks, and it still is divided across 16 grid, which mean states representation is by far less accurate.
2. The action-state values of Task1 was $16 * 16 * 16 * 16 * 2 = 131072$, for lunar lander, the action-state values are of size $16 * 16 * 16 * 16 * 16 * 16 * 2 * 2 * 4 = 268435456$, which is 2048 more values than previous Task to explore and converge. Thus, it is impossible to train for the same episodes.