

# Exercise 3

Eugeniu Vezeteu - 886240  
ELEC-E8125 - Reinforcement Learning

October 7, 2020

## 1 Task 1

In this task Q-learning algorithm was implemented.  
The following code implements the epsilon-greedy policy.

```
1 def get_action(state, q_values, greedy=False):
2     state = get_cell_index(state) #convert state from
        continuous to discrete
3     # TODO: Implement epsilon-greedy
4     if greedy or np.random.random() >= epsilon:
5         #greedy is true -> we test the agent, or prob 1-epsilon
6         action = np.argmax(q_values[state])
7     else: #random action with prob epsilon
8         action = int(np.random.random() * num_of_actions)
9
10    return action
```

The next code implements the Q table update

```

1 def update_q_value(old_state, action, new_state, reward, done,
2   q_array, steps, on_policy=True):
3     # TODO: Implement Q-value update
4     s = get_cell_index(old_state)    #current state
5     s_ = get_cell_index(new_state)   #next state
6
7     if done and steps < 199: #Q(terminal,:) is 0
8         q_grid[s][:] = 0.
9     else: #Q update formula
10         if on_policy:
11             a_ = get_action(new_state, q_grid, greedy=True)
12             target = reward + (gamma * q_grid[s_][a_]) - q_grid
13                 [s][action]
14         else: #off_policy
15             target = reward + (gamma * max(q_grid[s_])) -
16                 q_grid[s][action]
17
18     q_grid[s][action] += alpha * target

```

Figure 1 show the result of training for 20000 episodes using epsilon-greedy method, with constant  $\epsilon = 0.2$ .

According to lecture notes, the Greedy in limit with infinite exploration (GLIE) formula is:  $\epsilon = \frac{a}{a+k}$ .

So, in our case, desired  $\epsilon$  is 0.1 and  $k = 20000$ .

The value of the constant  $a$  was computed with the below formula

$$\epsilon = \frac{a}{a+k} \Rightarrow a = \epsilon * a + \epsilon * k \Rightarrow a - \epsilon * a = \epsilon * k \Rightarrow a * (1 - \epsilon) = \epsilon * k \Rightarrow a = \frac{\epsilon * k}{1 - \epsilon} = 199$$

Figure 2 show the result of training for 20000 episodes using GLIE method, with desired  $\epsilon = 0.1$ . after 20000 episodes

Files *q\_values\_task1.npy* and *value\_func\_task1* contains the Q-value and Value function for CartPole with GLIE.

## 2 Task 2

On this task the optimal value function was computed from Q table. The heatmap was plotted in terms of  $x$  and  $\theta$ , and averaging the values over  $\dot{x}$  and  $\dot{\theta}$

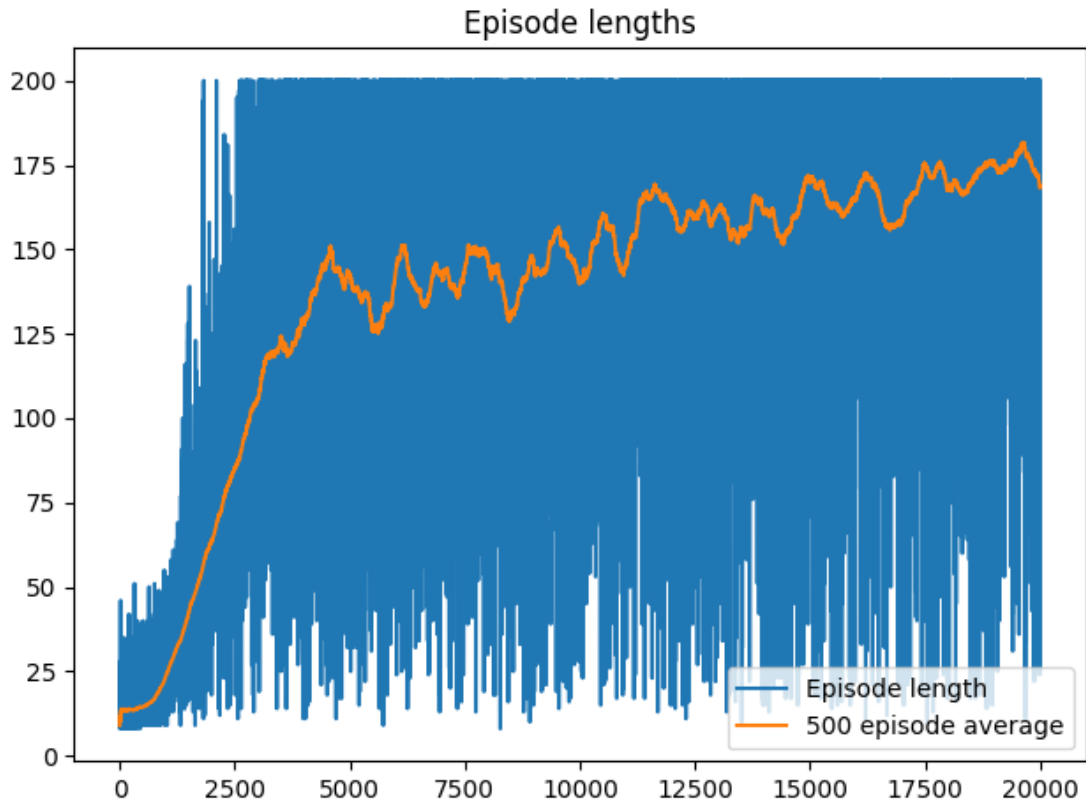


Figure 1: Task 1 Q-learning with epsilon greedy,  $\epsilon = 0.2$ , 20000 episodes

```

1 values = np.amax(q_grid, axis=4)
2 .
3 .
4 .
5 seaborn.heatmap(np.mean(values, axis=(1, 3)),xticklabels=True,
6                  yticklabels=True)
7 plt.xlabel("X")
8 plt.ylabel("theta")
9 plt.title('heatmap')
10 plt.show()

```

For the plots for task 2 please see Figures 3, 4, 5, 6

### 3 Question 1

1. **before the training** - Figure 3 shows that in the beginning the value function is zero for all states, because the environment is completely unexplored.
2. **after a single episode** - In Figure 4 we can see that after one episode the agent

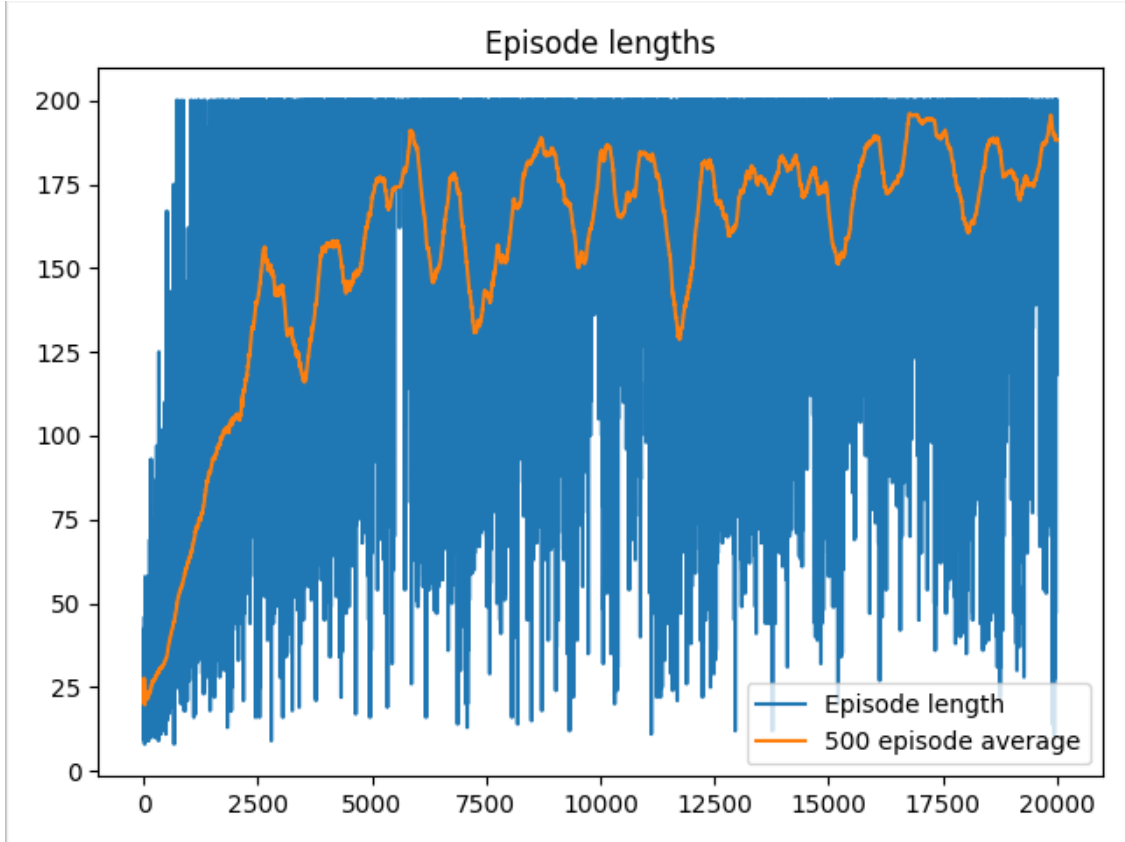


Figure 2: Task 1 Q-learning GLIE with desired epsilon,  $\epsilon = 0.1$ , 20000 episodes

started to explore the environment, and there are some non zero values cells.

3. **halfway through the training** - After the halfway of training we can see from Figure 5 that agent explored more states of the environment, which values are non zero, however there are black pixels (unexplored states), because the agent hasn't been in Q-table cell corresponding to those state and action.

Figure 6 shows the map of states after the 20000 episodes of training.

## 4 Task 3

Figures 7 and 8 shows performance of training starting from initial Q value 0 and 50.

## 5 Question 2.1

The model has a better performance starting with the initial Q values 50, for each state-action pair.

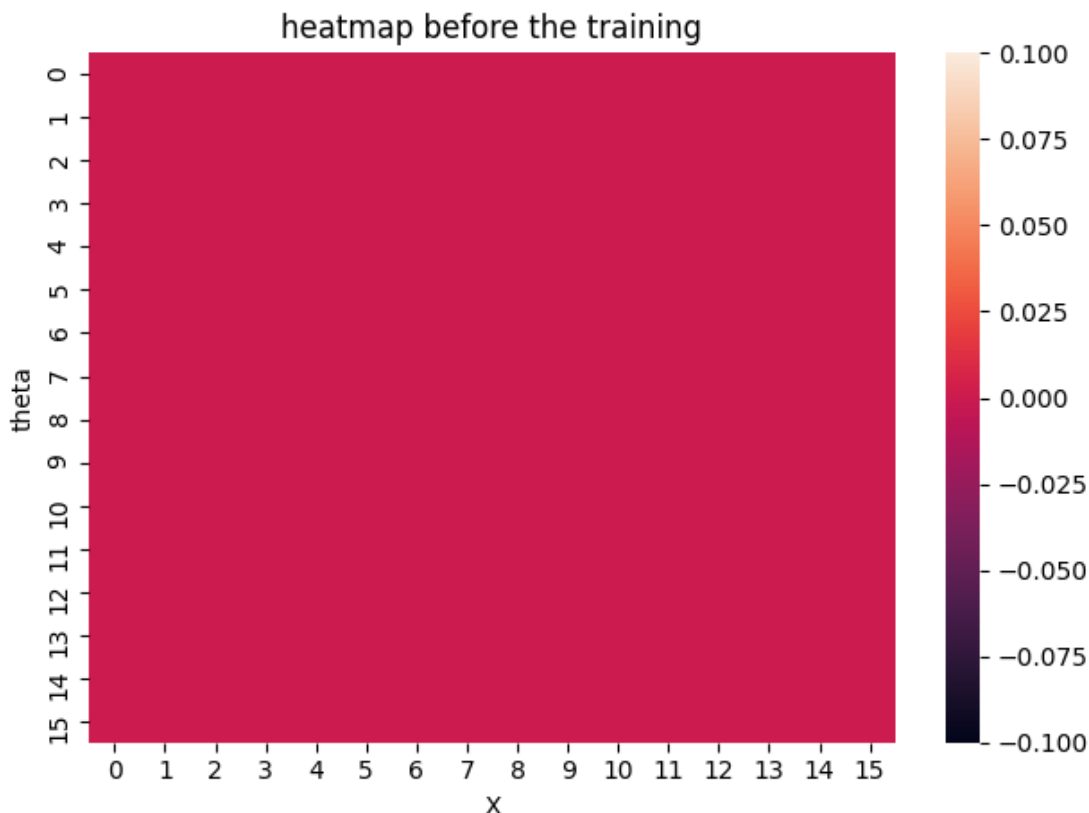


Figure 3: Task 2 a) Optimal value function before the training

## 6 Question 2.2

The agent performs better with higher initial Q values because each state-action pair will be chosen by agent at least once, in the beginning. For example, when the agent will take an action  $action = np.argmax(q\_values[state])$ , each action will have the same chance, but after some iterations, the Q value for state-action pair that did not lead to reward, will decrease, and the agent will take another action being in the same state. This encourages the exploration, the agent discover the best action-value pairs faster.

## 7 Task 4. Lunar lander

On this task the code from task 1 was modified.

In the following code, the limits of the state space are defined

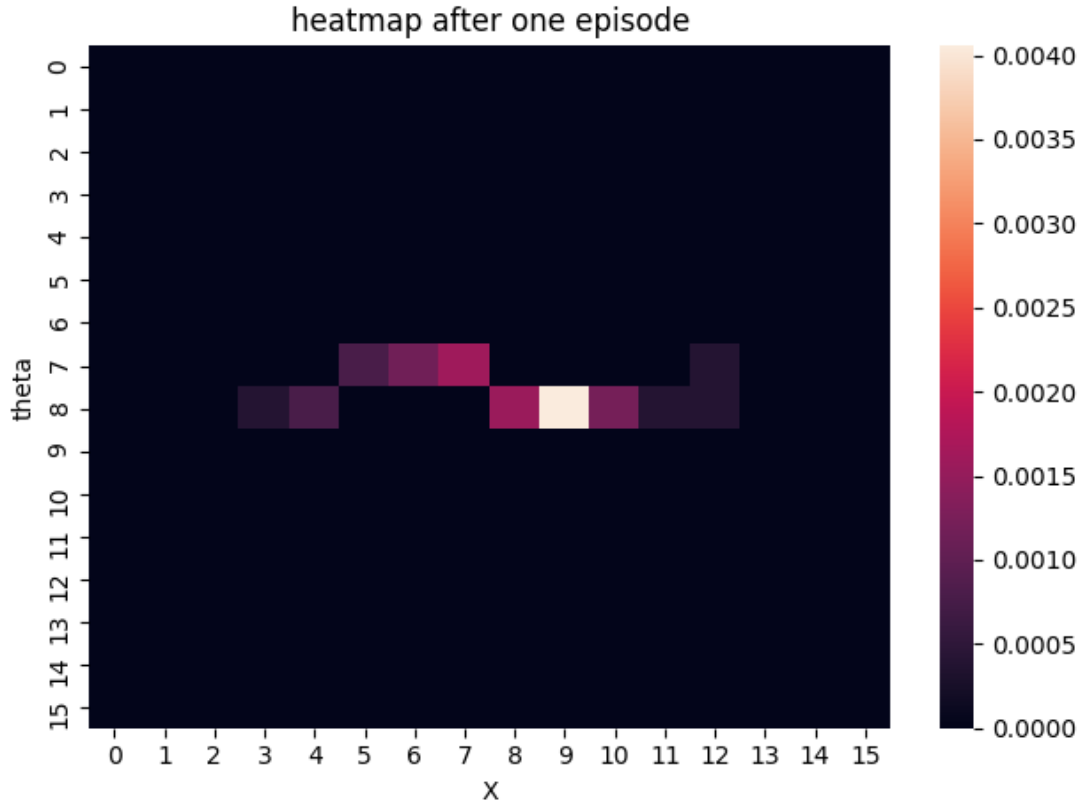


Figure 4: Task 2 b) Optimal value function after 1 episode

```

1 x_min, x_max = -1.2, 1.2
2 y_min, y_max = -.3, 1.2
3 xdot_min, xdot_max = -2.4, 2.4
4 ydot_min, ydot_max = -2., 2.
5 theta_min, theta_max = -6.28, 6.28
6 thetadot_min, thetadot_max = -8., 8.
7 cl_min, cl_max = 0, 1
8 cr_min, cr_max = 0, 1
9 .
10 x_grid = np.linspace(x_min, x_max, discr)
11 y_grid = np.linspace(y_min, y_max, discr)
12 xdot_grid = np.linspace(xdot_min, xdot_max, discr)
13 ydot_grid = np.linspace(ydot_min, ydot_max, discr)
14 theta_grid = np.linspace(theta_min, theta_max, discr)
15 thetadot_grid = np.linspace(thetadot_min, thetadot_max, discr)
16 cl_grid = np.linspace(cl_min, cl_max, 2)
17 cr_grid = np.linspace(cr_min, cr_max, 2)
18 .
19 q_grid = np.zeros((discr, discr, discr, discr, discr, discr, 2,
    2, num_of_actions)) + initial_q #LunarLander

```

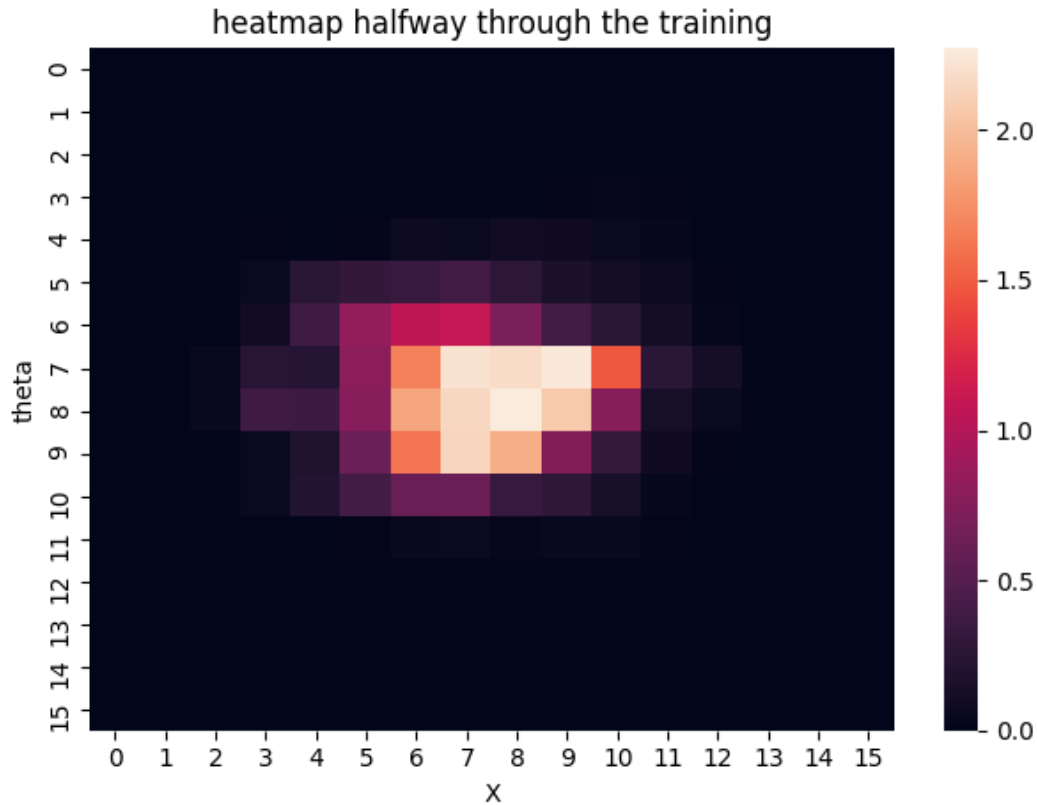


Figure 5: Task 2 c) Optimal value function halfway through the training

The function *get\_cell\_index(state)*: was changed to:

```

1  #used for LunarLander
2  def get_cell_index(state):
3      x = find_nearest(x_grid, state[0])
4      y = find_nearest(y_grid, state[1])
5      xdot = find_nearest(xdot_grid, state[2])
6      ydot = find_nearest(ydot_grid, state[3])
7      theta = find_nearest(theta_grid, state[4])
8      thetadot = find_nearest(thetadot_grid, state[5])
9      cl = find_nearest(cl_grid, state[6])
10     cr = find_nearest(cr_grid, state[7])
11     return x, y, xdot, ydot, theta, thetadot, cl, cr

```

The policy and update Q table are the same.

The performance of the agent on Luna Lander is presented in Figure 9



Figure 6: Task 2 Optimal value function after the training

## 8 Question 3.1

The agent was not able to learn the Luna Lander environment. This environment is more complicated than the CartPole env, and the previous setup is simple for this problem. The world discretization used in the CartPole env, is very poor for Luna Lander, the discretization of the Luna Lander cannot sufficiently cover the entire state space.

In our case, we discretize the state space into 16 bins for each feature, we may increase the number of bins for each feature for a better performance. However, it is not recommended, since we will have a lot of states, a lot of memory will be consumed, and in addition we need a lot of computational power to explore each state.

A better idea is to use function approximators for state space (Neural Networks, RBF).



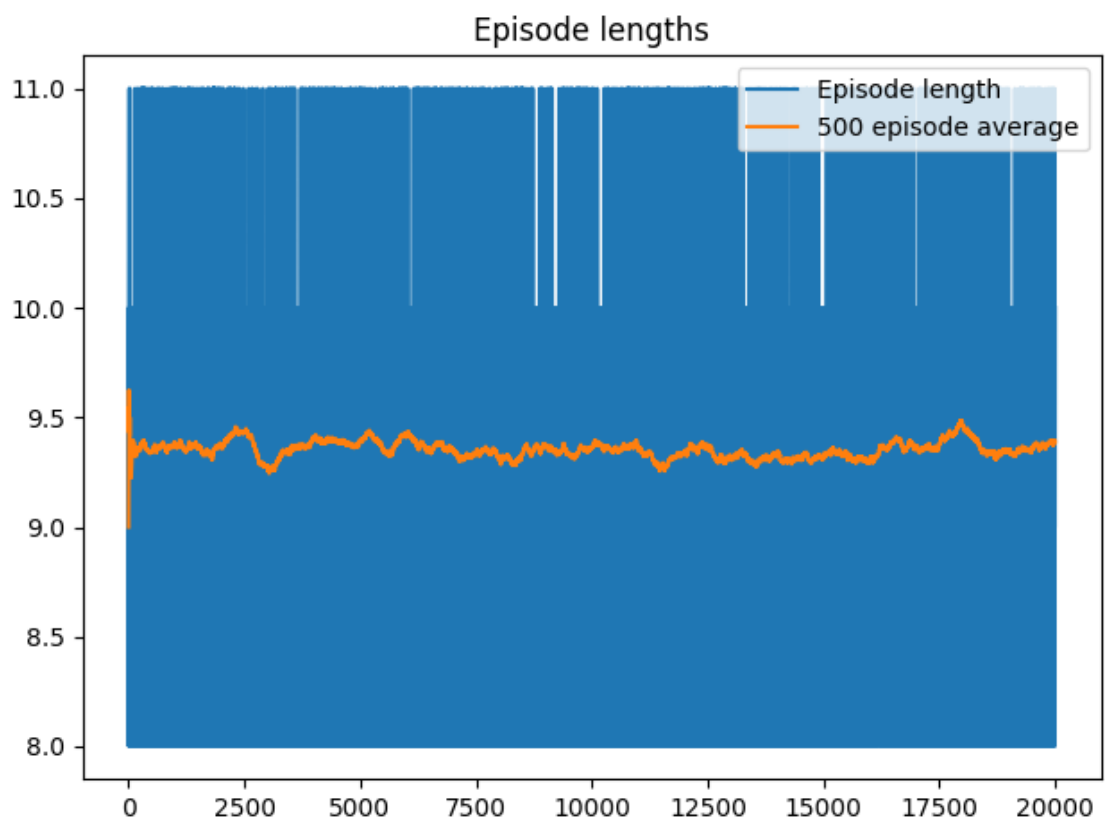


Figure 7: Task 3 Q-learning, greedy policy with Q values starting from 0

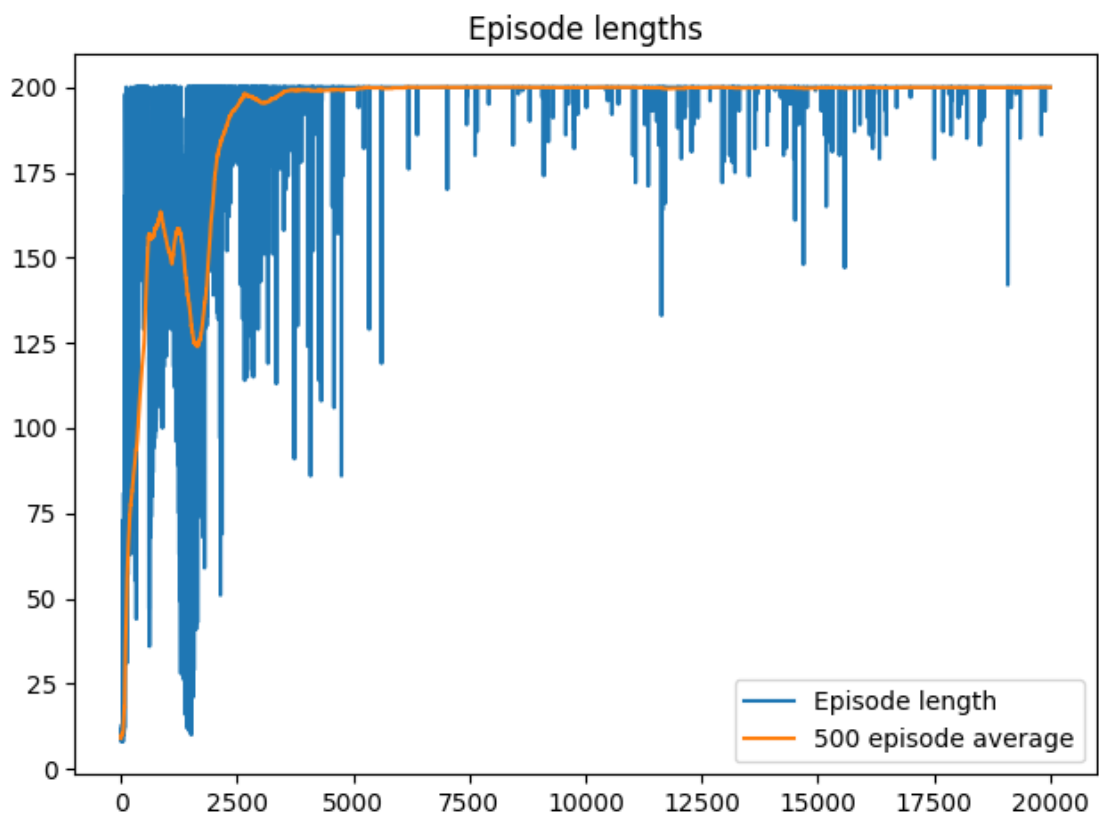


Figure 8: Task 3 Q-learning, greedy policy with Q values starting from 50

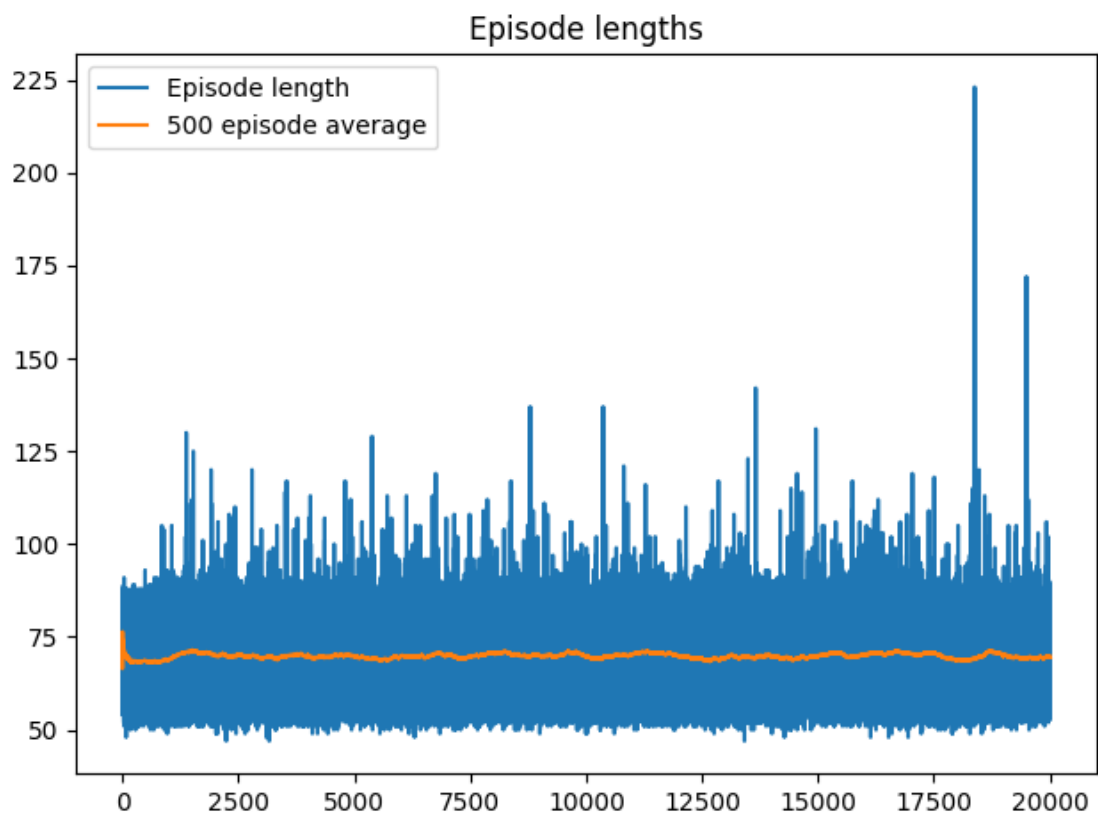


Figure 9: Task 4 Q-learning, GLIE for Luna Lander