

Exercise 4

Gengcong Yan - 1009903
ELEC-E8125 - Reinforcement Learning

October 22, 2021

1 Task 1

```
1 def featurize(self, state):
2     if len(state.shape) == 1:
3         state = state.reshape(1, -1)
4         # Task 1a: TODO: Implement handcrafted features as (s,
5             abs(s))
6         #(x1, x2, x3, x4, abs(x1), abs(x2), ...)
7         # return np.concatenate((state,abs(state)),axis=1)
8
9         ## Task 1b: Use the RBF features using the featurizer
10        return self.featurizer.transform(self.scaler.transform(
11            state))
12
13 def get_action(self, state, epsilon=0.0):
14     # TODO Task 1: Implement epsilon-greedy
15     # Hint: See exercise 3
16     if np.random.uniform(0,1)<epsilon:
17         # random action
18         return env.action_space.sample()
19     else:
20         # Choose current best action
21         feas=self.featurize(state)
22         qvalues=[ q.predict(feas)[0] for q in self.q_functions
23             ]
24         re=np.argmax(np.array(qvalues))
25         return re
```

```

1 def single_update(self, state, action, next_state, reward, done
  ):
2     # Calculate feature representations of the
3     # Task 1: TODO: Set the feature state and feature next
       state
4     featurized_state = self.featurize(state)
5     featurized_next_state = self.featurize(next_state)
6
7     # Task 1: TODO Get Q(s', a) for the next state
8     qss= np.array([q.predict(featurized_next_state)[0] for
       q in self.q_functions])
9     next_qs = np.amax(qss,axis=0)
10
11    # Calculate the updated target Q- values
12    # Task 1: TODO: Calculate target based on rewards and
       next_qs
13    if done: # No next state
14        target=reward
15    else:
16        target = reward+self.gamma*next_qs
17    self.q_functions[action].partial_fit(featurized_state ,
       (target,))

```

The rewards with 2 different representations are shown in Fig.1 and Fig.2.

2 Question 1

NO, Cartpole using linear features failed to learn accurately Q-values. because a linear regressor may be too simple to represent the value approximation. Coaction between different features is not adequately expressed. the position of cartpole may vary differently based on different velocity, which leads to good or bad direction. We need some more complex expressions to reflect the relationship between the different features

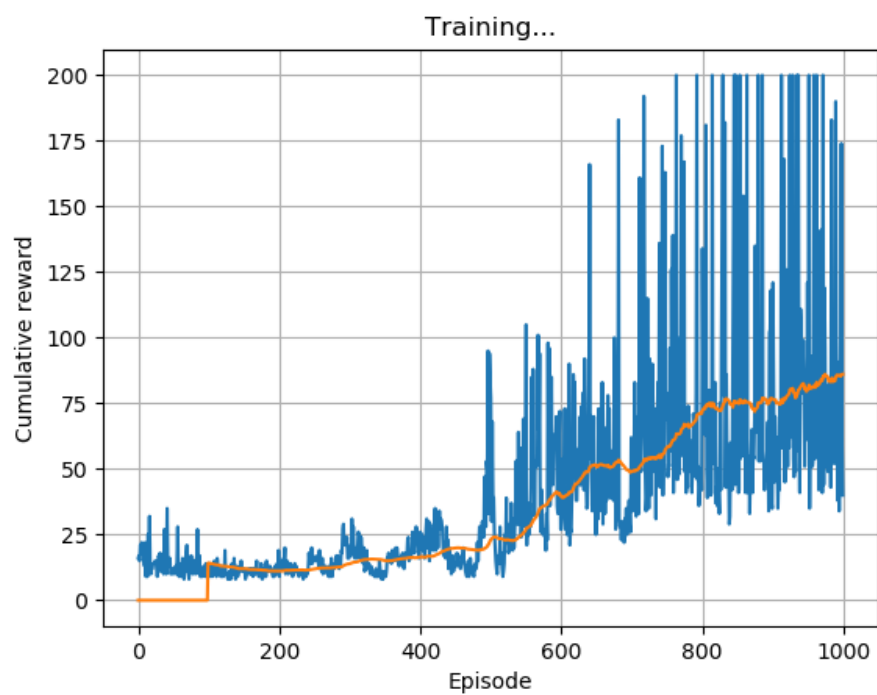


Figure 1: Rewards with handcrafted feature

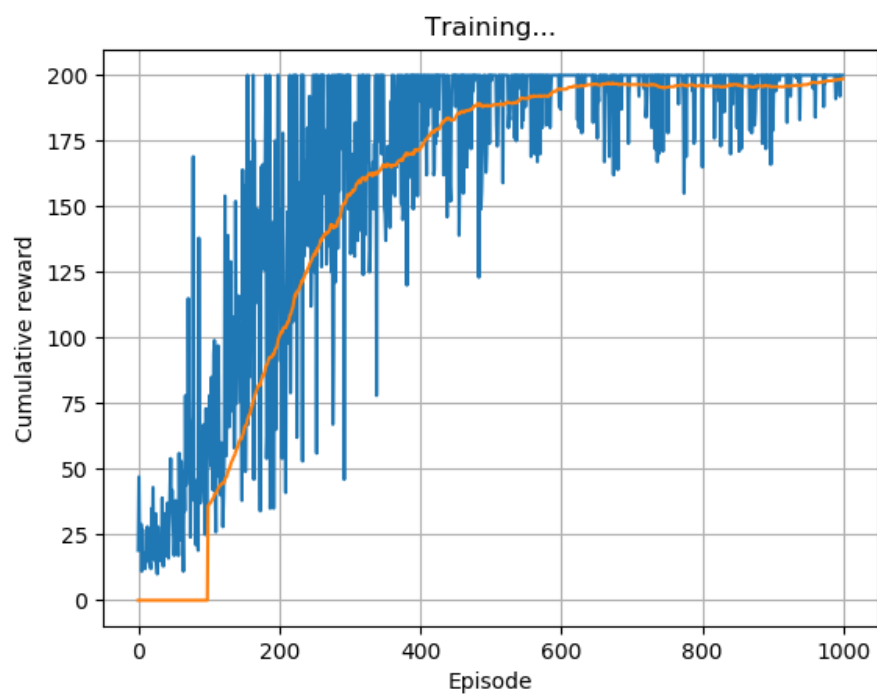


Figure 2: Rewards with RBF

3 Task 2

```
1  ##### Code in function update_estimator, not all
2  # Task 2: TODO: Reformat data in the minibatch
3      states = []
4      action = []
5      next_states = []
6      rewards = []
7      dones = []
8      for i in range(len(samples)):
9          states.append(samples[i].state)
10         action.append(samples[i].action)
11         next_states.append(samples[i].next_state)
12         rewards.append(samples[i].reward)
13         dones.append(samples[i].done)
14     states = np.array(states)
15     action = np.array(action)
16     next_states = np.array(next_states)
17     rewards = np.array(rewards)
18     dones = np.array(dones)
19
20     # Task 2: TODO: Calculate  $Q(s', a)$ 
21     featurized_next_states = self.featurize(next_states)
22     next_qs = []
23     for fns in featurized_next_states:
24         fns=fns.reshape(1,-1)
25         next_qs.append(np.amax(np.array([q.predict(fns)[0] for
26             q in self.q_functions]),axis=0))
27     next_qs=np.array(next_qs)
28
29     # Calculate the updated target values
30     # Task 2: TODO: Calculate target based on rewards and
31         next_qs
32
33     targets = rewards+self.gamma*next_qs* (1-dones)
```

The rewards using minibatch update and experience replay are shown in Fig.3 and Fig.4.

4 Question 2

4.1 Question2.1

With experience replay, the training process became faster, the rewards increased a lot just in the beginning. This strategy will save computation power and converge faster, but ultimately the final performance won't get any better reward.

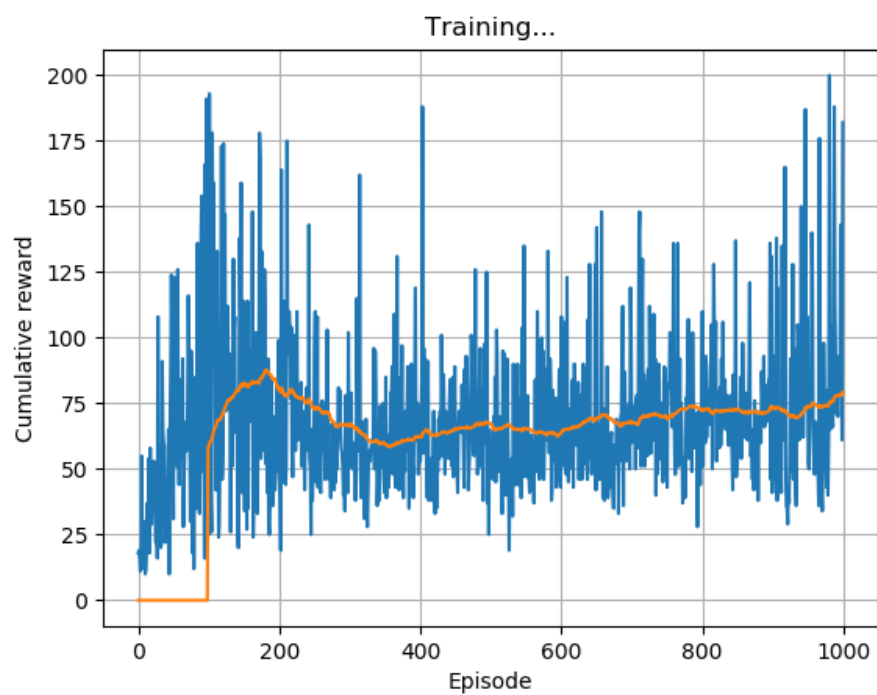


Figure 3: Rewards with handcrafted feature (experience replay)

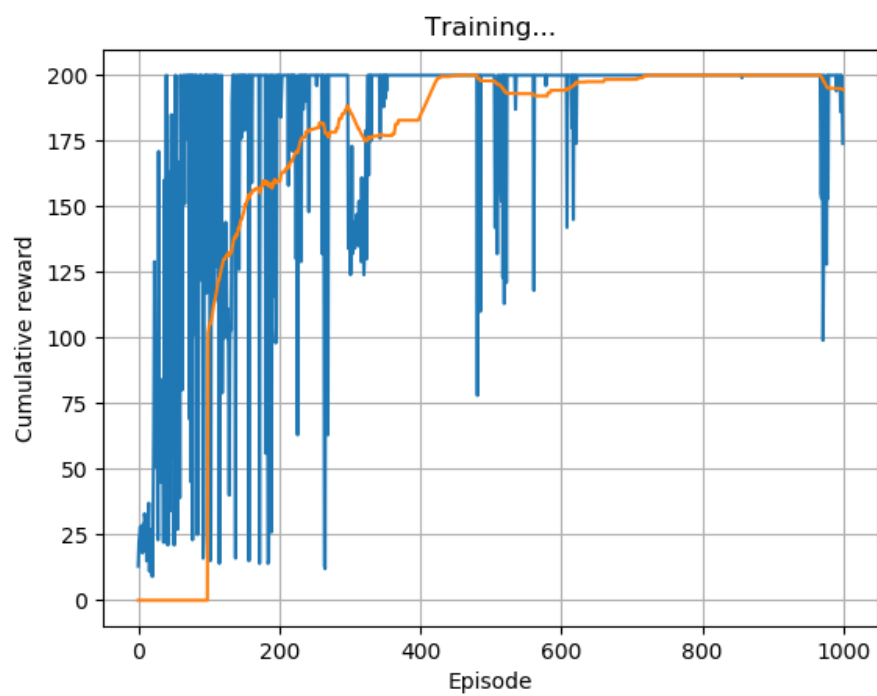


Figure 4: Rewards with RBF (experience replay)

4.2 Question2.2

- **Pros:** It saves computing power and train fast. It will work well if you can effectively construct representations of the important features in problem. But this is generally very difficult.
- **Cons:** It's too simple in some complex problems, resulting in bad performance. There is a lack of interaction between the different features. Their interconnection with each other can influence the training effect.

4.3 Question2.3

No, the grid based methods need much more episodes to train for convergence. grid base methods have to learn each behavior with huge episodes. The function approximation methods can generate different behaviors more efficiently.

5 Task 3

```
1  # There are also code in EX3 used here
2
3  def get_cell_index2(x,th):
4      x = find_nearest(x_grid, x)
5      # v = find_nearest(v_grid, state[1])
6      th = find_nearest(th_grid, th)
7      # av = find_nearest(av_grid, state[3])
8      return x, th,
9
10 x_grid=x_grid.round(2)
11 th_grid=th_grid.round(2)
12
13 for i,x in enumerate(x_grid):
14     for j,th in enumerate(th_grid):
15
16         x,th=get_cell_index2(x,th)
17         state=np.array([x,0,th,0])
18         action=agent.get_action(state)
19         q_grid[(i,j)]=action
20
21 # print(q_grid,x_grid,th_grid)
22
23 print("Plotting Policy with RBF experience replay")
24 sb.heatmap(q_grid,xticklabels=x_grid,yticklabels=th_grid,cbar=
25     False)
26 plt.title("Policy with RBF experience replay")
27 plt.xlabel("X")
28 plt.ylabel("Theta")
29 plt.show()
```

Base on the code in Exercise 3, after some change we can plot the heatmap of policy learned with RBF and experience replay in Fig.5. Black and white represent 0 and 1 in plot, showing the current action in different situation. We can see theta of cartpole Basically no impact on action, but the position decides different reactions preventing cartpole falling down.

6 Task 4

The rewards in Cartpole and Lunar landing environment with DQN implementation are shown in Fig.6 and Fig.7.

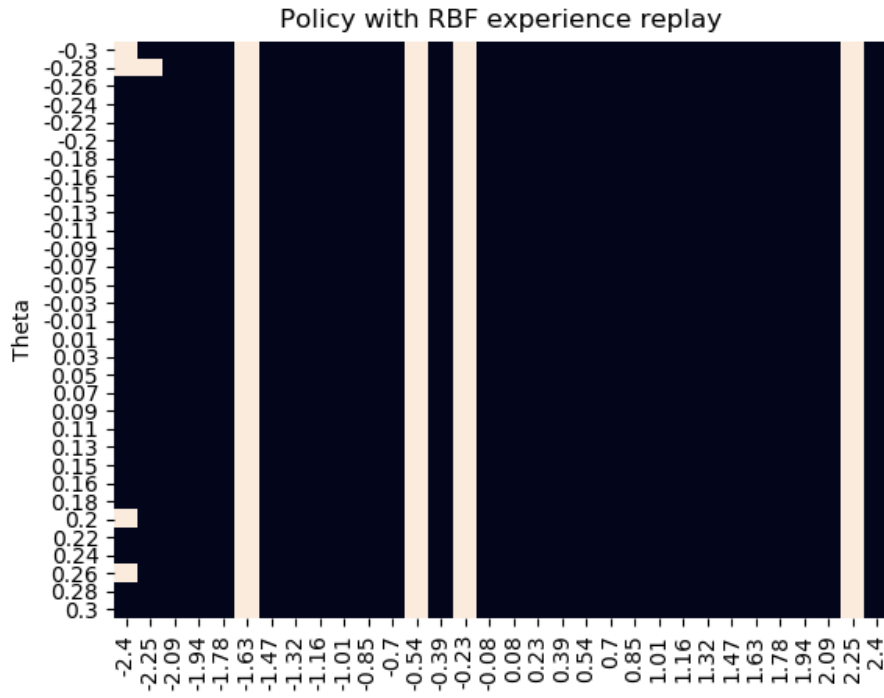


Figure 5: Heatmap of policy learned with RBF and experience replay

```

1  # get_action in DQN_agent
2  def get_action(self, state, epsilon=0.05):
3      # TODO: Implement epsilon-greedy
4      # Hint:
5      # If greedy
6      #   convert state to torch float tensor
7      #   get q values from policy net (consider using with
8      #   torch.no_grad())
9      #   return action
10     # If not greedy
11     #   return random action
12     if random.random() < epsilon:
13         # random action
14         return random.randrange(self.n_actions)
15     else:
16         # find best action
17         with torch.no_grad():
18             state = torch.from_numpy(state).float()
19             qvalues = self.policy_net(state)
20             # print("qvalues", qvalues)
21             qmax_action = np.argmax(qvalues).item()
22             # print("qmax_action", qmax_action)
23             return qmax_action

```

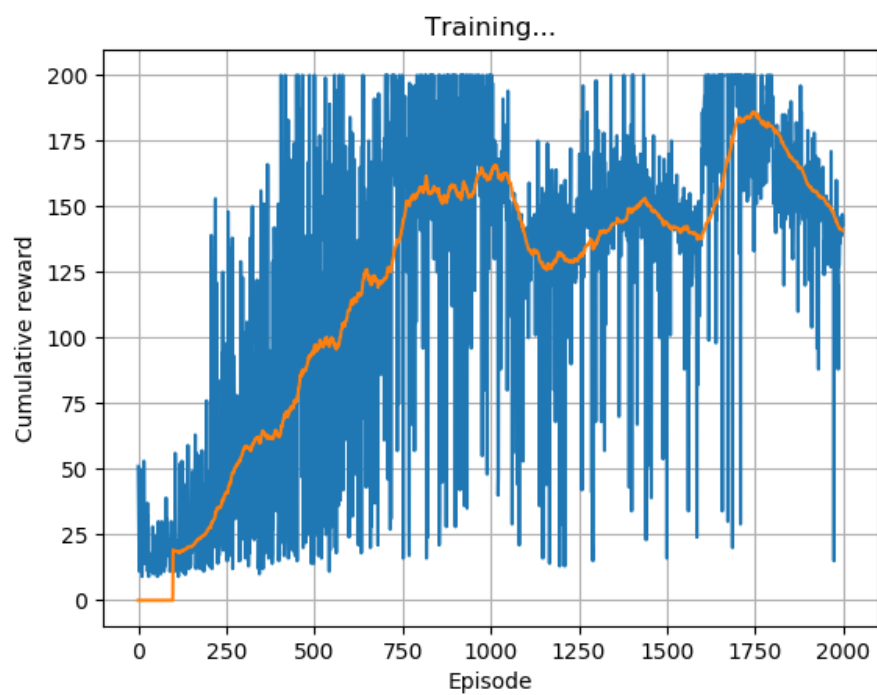


Figure 6: Rewards in Cartpole with DQN

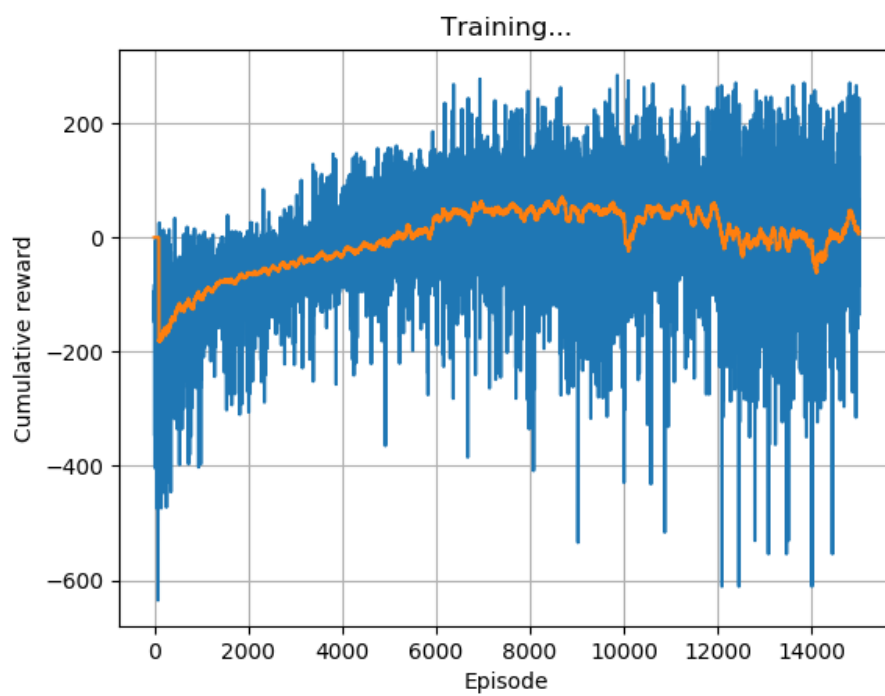


Figure 7: Rewards in Lunar landing with DQN

7 Question 3

7.1 Question 3.1

We **can not** directly use Q-learning in continuous action space. But with some mechanism like discretion and function approximation, we can still apply it in continuous space.

7.2 Question 3.2

In Q-learning, we always need to get $\max Q(S_{T+1}, a)$ for the next state. we need to iterate all actions to obtain maximum value of it. if the environment is a continuous action space, finding the maximum will be time and computation consuming. Even if we discretize the space we may not get a good result. So we can use advanced methods like policy gradient and Actor-Critic.

8 Feedback

1. 10h