

Exercise 4

Tommaso Brumani - 100481325
ELEC-E8125 - Reinforcement Learning

October 10, 2022

Task 1

To implement Q learning using function approximation, the following code was used in the 'update()' and 'get_action()' functions of the 'rbf_agent.py' file:

- get_action():

```
1 ##### Your code starts here #####
2
3 # determine whether to choose greedy policy
4 choose_rand = np.random.random() <= epsilon
5
6 if choose_rand:
7     # take random action
8     return np.random.randint(0, self.num_actions)
9 else:
10    # featurize
11    f_state = self.featurize(state)
12
13    # take greedy action
14    q_values = [func.predict(f_state) for func in self.
15                q_functions]
16    return np.argmax(q_values)
17 ##### Your code ends here #####
```

- update():

```

1 ##### You code starts here #####
2
3 # featurize the state and next_state
4 f_state = self.featurize(batch.state)
5 f_next_state = self.featurize(batch.next_state)
6 q_tar = np.zeros(self.batch_size)
7
8 # variable for storing maximum q for each batch element
9 max_q = np.zeros(self.batch_size)
10
11 # extract next state batch elements in which the
    terminating state is not reached
12 fns_not_done = f_next_state[batch.not_done.astype(bool), :]
13
14 # predict q value for each possible action
15 q_actions = np.array([func.predict(fns_not_done) for func
    in self.q_functions])
16
17 # set maximum q for non terminating batch elements
18 max_q[batch.not_done.astype(bool)] = np.max(q_actions, axis
    =0)
19
20 # calculate q_target
21 q_tar = batch.reward + self.gamma*max_q
22
23 ##### You code ends here #####

```

The model was trained using two different features for state representation:

- (a) A handcrafted feature vector $\Phi(s) = [s, |s|]^T$, the training performance plot for which is reported in Figure 1.

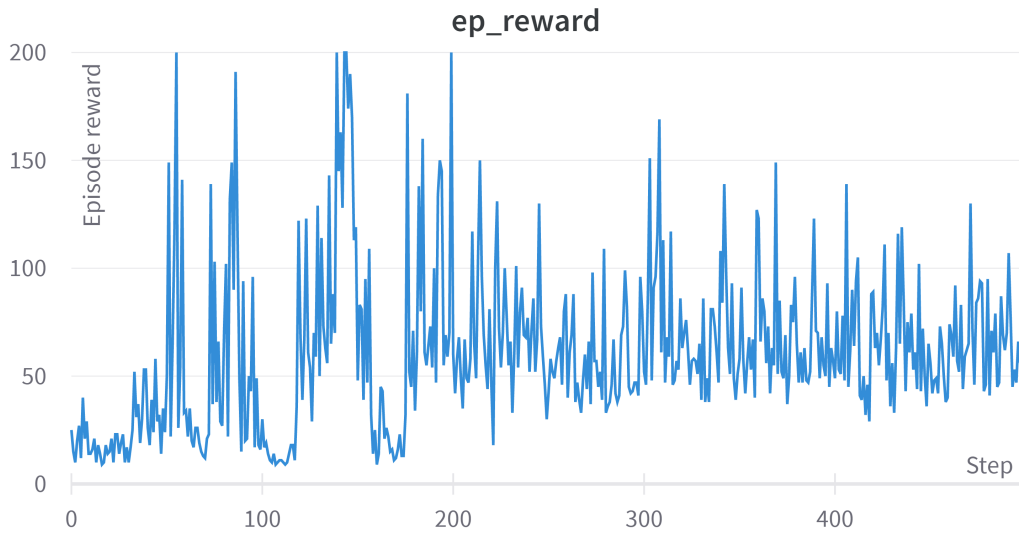


Figure 1: training performance using handcrafted feature vector.

- (b) A radial basis function representation, the training performance plot of which is reported in Figure 2.

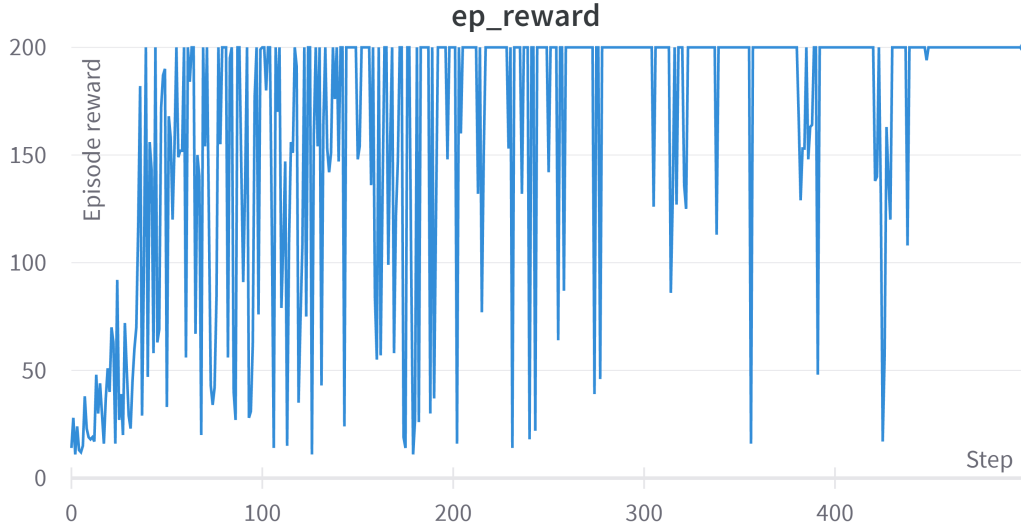


Figure 2: training performance using radial basis function representation.

Question 1.1

It would not be possible to accurately learn Q values for the Cartpole problem using linear features, since we can expect the real value function to take a symmetric shape around the center of the space (the cart should try to keep the pole going away from the two sides of the

screen), and such a function cannot be achieved using linear features with a linear regressor, which could only produce values proportional to the state.

It then becomes necessary to map the state to a space with a higher number of features, so that a linear regressor may find a fitting linear function over a higher number of dimensions, resulting in a non-linear function in the original state space.

From the fact that using the handcrafted feature vector $\Phi(s) = [s, |s|]^T$ still fails to converge to a good policy, we can deduce that the number of features in that instance is still insufficient to learn the problem.

Question 1.2

The use of an experience replay buffer makes it possible to learn from past experiences, which achieves a number of positive effects.

For one, it allows to increase sample efficiency by making it possible to reuse samples multiple times, and the ability to store samples in memory allows to take mini-batches to speed up computation.

Furthermore, it helps stabilize the learning process by breaking the temporal correlations of subsequent updates, reducing bias: by sampling the batches randomly from memory, the algorithm learns from samples that have low correlation and are thus more close to the assumption of being i.i.d.

Question 1.3

Tabular methods are generally sample-inefficient compared to RBF function approximation, as the number of states is usually greater than that of parameters used in function approximation.

Since modifying a single parameter has a cascading effect on the value of multiple states, function approximation makes it possible to train a model using a lower number of samples compared to tabular methods, as the variables to be learned are fewer.

Task 2

The 2D plot of policy learned with RBF in terms of x and θ for $\dot{x} = 0$ and $\dot{\theta} = 0$ is reported in Figure 3.

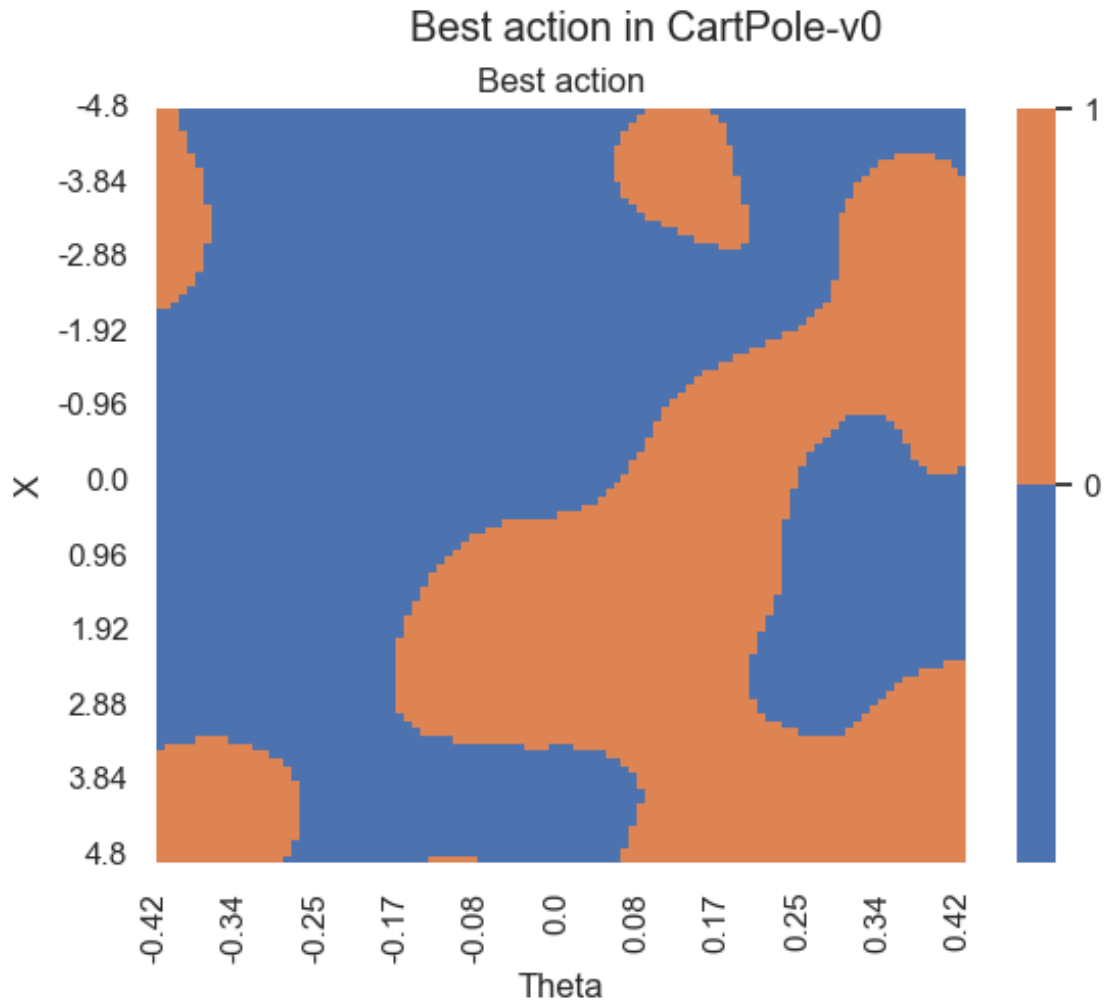


Figure 3: learned policy with RBF.

Task 3

To implement the DQN agent in the 'dqn_agent.py' file, the code in the 'update()' and 'get_action()' functions were modified as follows:

- update():

```
1  ##### You code starts here #####
2
3  # calculate the q(s,a)
4  q_actions_policy = self.policy_net(batch.state)
5  qs = torch.gather(input=q_actions_policy, dim=1, index=
    batch.action.type(torch.int64))
6
7  # calculate q target
8  q_actions_target = self.target_net(batch.next_state)
9  max_q = torch.max(input=q_actions_target, dim=1).values
10
11 # set to zero where final state is reached
12 max_q = torch.reshape(max_q, batch.not_done.size())*batch.
    not_done
13
14 # calculate q target
15 q_tar = batch.reward + self.gamma*max_q
16
17 #calculate the loss
18 loss = (qs - q_tar.detach())**2 * 0.5
19 loss = loss.sum()
20
21 # clip grad norm and perform the optimization step
22 self.optimizer.zero_grad()
23 loss.backward()
24 torch.nn.utils.clip_grad_norm_(self.policy_net.parameters()
    , self.grad_clip_norm)
25 self.optimizer.step()
26
27 ##### You code ends here #####
```

- get_action():

```

1 ##### You code starts here #####
2
3 # determine wether to choose greedy policy
4 choose_rand = np.random.random() <= epsilon
5
6 if choose_rand:
7     # take random action
8     return np.random.randint(0, self.n_actions)
9 else:
10    # take greedy action
11    q_actions_target = self.policy_net(torch.from_numpy(
12        state).to(device))
13    return torch.argmax(q_actions_target).item()
14 ##### You code ends here #####

```

The agent was trained in both the Cartpole and Lunar Lander environments, and its performance is presented in Figure 4 and Figure 5, respectively.



Figure 4: training performance of DQN in the Cartpole environment.

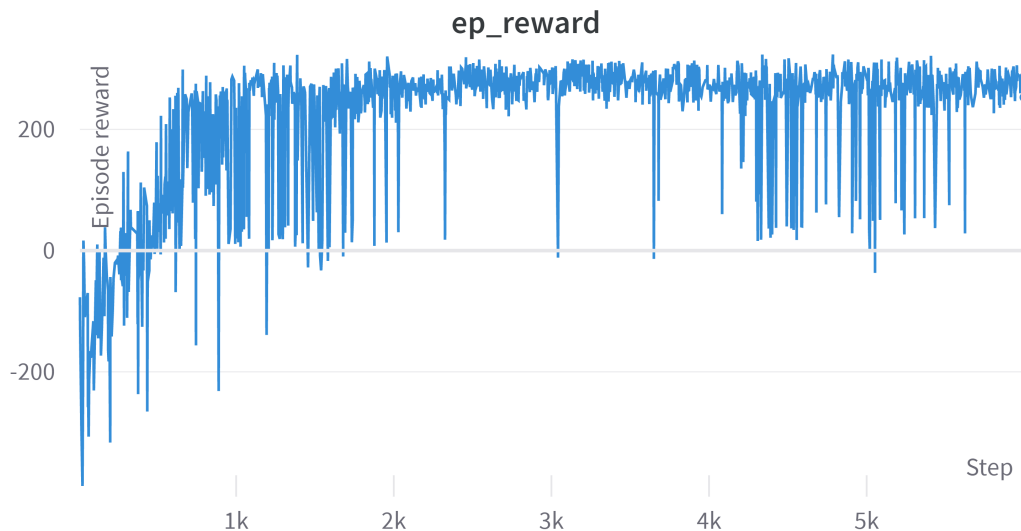


Figure 5: training performance of DQN in the Lunar Lander environment.

Question 3.1

Standard Q-learning cannot be used directly in continuous action space environments because discrete and finite actions and states are necessary for it to calculate the Q function for specific state-action pairs.

Question 3.2

The main difficulty when dealing with continuous action spaces is in calculating the maximum Q value over all actions for the next state. If the action space is continuous this problem becomes very complex, as simply evaluating the value over all possible actions is impossible. A possible (albeit imperfect) solution would be to simply discretize the action space.

Question 3.3

The use of an additional neural network is necessary to stop the bias derived from bootstrapping from causing the estimated Q values to diverge, for example in instances in which an outlier reward has been registered by chance. If a single neural network was used, for example, registering an abnormally high reward value during a transition could cause the algorithm to iterate over that estimate, thus it skewing its estimates.

Using a second network relatively fixed in time, instead, gives the algorithm the opportunity to correct this bias over multiple iterations by identifying its sharp divergence from previous estimates.

In this instance where a single network is used, not stopping the gradient of the target

Q value would result in it being taken into account and propagated when calculating the gradient, which might skew the update.