

---

# EL282805 - Reinforcement Learning - Lab 2 - Report

---

Adrian Chmielewski-Anders, Leo Zeitler  
9512115537, 9509289634  
{amca3, llze}@kth.se

## 1

- (a) The state space can be modeled by a vector  $s = (x, \theta, \dot{x}, \dot{\theta})$  where  $x$  is the position from the center of the track and  $x \in [-X, X]$  where  $X \geq 2.4\text{m}$  are the maximum lengths from the middle of the track. Furthermore,  $\theta \in [-90, 90]$  which is measured in degrees, is the angle of the pole.  $\dot{x}$  the linear velocity of the cart which is in  $\mathbb{R}$ . Similarly,  $\dot{\theta}$  is the angular velocity at the tip of the pole which is also in  $\mathbb{R}$ . All state variables take values on  $\mathbb{R}$  (or interval) or a discretized real line in software. Nonetheless, this makes the state space very large. The action space can be modeled as  $\mathcal{A} = \{-1, 1\}$  where the agent can push from the left or right with force 1N. The reward function is defined, though unknown to the agent:

$$r_t(s, a) = \begin{cases} 1 & -12.5 \leq \pi(s, a)_\theta \leq 12.5 \text{ and } -2.4 \leq \pi(s, a)_x \leq 2.4 \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

The issue is the state space is very large for tabular RL methods. Due to this, one might consider deep RL methods. The problem is modeled as an infinite time horizon. Though, the environment ends the episode after 200 steps or when or when the agent encounters a zero reward.

- (b) `main` First it creates the environment, and instantiates the agent. It then initializes some test states for plotting and checking convergence. Then we go through a single training loop with  $N$  episodes, saving values to plot convergence, and going through normal training steps for each episode. That is, take an action, observe rewards and next states and update the two networks. The two plots generated are a plot of the total undiscounted reward of each episode that the agent goes through which are chosen via the  $\epsilon$ -greedy policy. The next graph is for each episode the maximum of the actions averaged over the same random states. In the training loop we are saving the state, action, reward, next state, and done variables into the replay buffer.

`build_model` Creates a simple neural network for approximating the  $Q$  function. The architecture is shared over  $\theta$  and  $\phi$ .

`update_target_model` Sets the target models weights  $\phi$  to the models weights  $\theta$ .

`get_action` Gets an  $\epsilon$ -greedy action, initially it is completely random. However, in later parts of the assignment it is changed to be a  $\epsilon$  greedy action policy with respect to  $Q_\theta$ .

`append_sample` Adds a state, action, reward and next state to the replay list.

`train_model` Trains the model  $\theta$ , going through each step of the Deep Q-learning. In the function `batch_size` items from the replay buffer are sampled and used as the inputs to update the network.

`plot_data` Is used for plotting the behavior over time of the training. Specifically for sample  $Q_\theta$  values and total reward over each episode.

- (c) The pseudo code in Algorithm 1 gives a outline of the deep Q-learning algorithm with experience replay. The lines 15 to 18 are equivalent to lines 192 to 199 in the given code file. Our function `Get-Action` is defined by the function `get_action`. The update of the weights of the target network is called in line 206, whereas the function is implemented

---

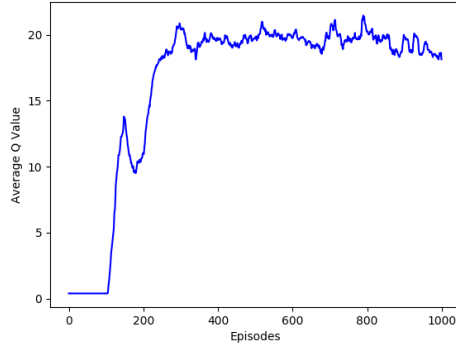
**Algorithm 1** The Deep Q-Learning algorithm

---

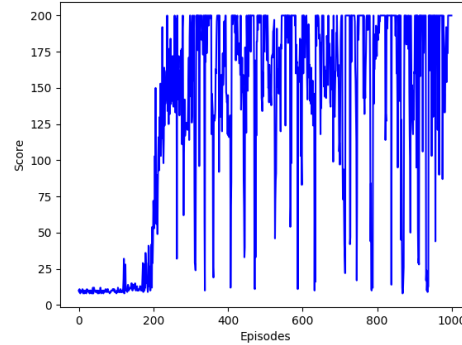
```
1: function UPDATE-STEP( $\theta, \phi, M, \text{batch\_size}$ )
2:   for  $i = 1 \dots \text{batch\_size}$  do
3:     Sample  $s_i, a_i, r_i, s'_i, d$  from  $M$   $\triangleright$  Sample (state, action, reward, next state, done) tuple
       from memory buffer
4:     if  $d$  is True then  $\triangleright$  Check whether it was a final step
5:        $y \leftarrow r_i$   $\triangleright$  If so, assign the reward to the target
6:     else
7:        $y \leftarrow r_i + \lambda \max_b Q_\phi(s'_i, b)$   $\triangleright$  If not, assign current target value prediction to target
       variable
8:      $\theta \leftarrow \theta + \alpha(y - Q_\theta(s_i, a_i)) \nabla_\theta Q_\theta(s_i, a_i)$   $\triangleright$  Update weights of the network according to
       stochastic gradient descent
9: function DQN( $s_1, N, \text{env}, C, \text{batch\_size}$ )
10:  Initialize  $\theta, \phi$   $\triangleright$  Initialize weights for target network and prediction network
11:  Initialize  $M$  as empty array  $\triangleright$  Initialize memory buffer for experience replay
12:  Initialize total_steps as 0  $\triangleright$  Initialize counter that is used for updating the weights of the
       target network every  $C$  steps
13:  for  $e = 1 \dots N$  do  $\triangleright$  Loop over all episodes
14:    while  $t \in \text{episode}$  do  $\triangleright$  Loop over all time steps in episode
15:       $a_t = \text{GET-ACTION}(s_t, Q_\theta)$   $\triangleright$  Get  $\epsilon$ -greedy action
16:       $(r_t, s_{t+1}, d) = \text{env.STEP}(s_t, a_t)$   $\triangleright$  Observe reward, next state and whether the task
       came to an end
17:      Append  $(s_t, a_t, r_t, s_{t+1}, d)$  to  $M$   $\triangleright$  Add experience to memory buffer
18:      UPDATE-STEP( $\theta, \phi, M, \text{batch\_size}$ )  $\triangleright$  Update weights  $\theta$  of the prediction network
19:      total_steps  $\leftarrow$  total_steps + 1  $\triangleright$  Update the step counter
20:      if total_steps %  $C == 0$  then
21:         $\phi \leftarrow \theta$   $\triangleright$  Update target network weights every  $C$  iterations
22:  return  $\theta$ 
```

---

Figure 1: The two plots for the initial parameters after doing parts (e) and (f)



(a) The Q-plot.

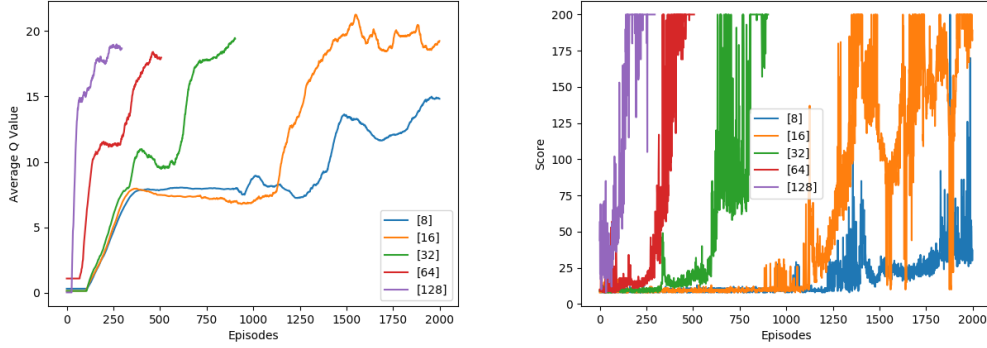


(b) The total reward plot.

in `update_target_model` in line 68. The functionality of the pseudo-code function *Update-Step* is implemented in `train_model`. The calculation of the target value should be implemented between lines 114 and 124. The update of the weights  $\theta$  is defined in the line 127 with the function call `self.model.fit`.

- (d) The network in the `build_model` function is a two layered feed-forward neural network with 16 units in the first hidden layer and 2 in the output layer. It uses ADAM as the optimizer and uses the default activation which is linear for the output layer and the ReLU activation function for the first hidden layer. The loss function is MSE. It uses constant learning rate.

Figure 2: Experimentation with the effect of different parameters on the network complexity.



(a) The Q-value plot for a shallow network with a single hidden layer and different number of neurons (b) The score plot for a shallow network with a single hidden layer and different number of neurons

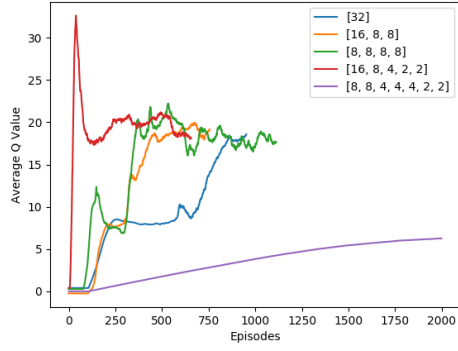
(e) See the attached code.

(f) The code for parts (e) and (f) is appended below. The two plots with all the default parameters are shown in Figure 1

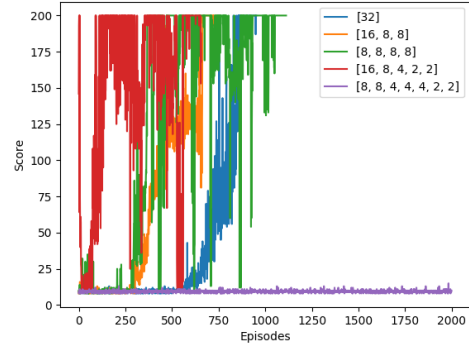
(g) We investigated the impact of the number of hidden neurons and the depth of the network on the performance of the system. As mentioned in the task description, we did some parameter fiddling in order to make the agent solve the task. We discovered that a buffer size of 5000, 2000 epochs and a learning rate set to 0.0005 shows some decent results. We used these parameter settings in all following experiments conducted within the task (g). Firstly, we examined the effect of the number of hidden neurons organised in shallow network with a single hidden layer on the learning. We stopped the learning procedure as soon as the task was considered to be solved, namely when the average reward over 100 consecutive episodes is greater than or equal to 195 to have a measure for the time to converge. As it is depicted in figure 2 the algorithm converges quicker the more hidden neurons are involved. Although the value that the Q-function approaches does not differ for different architectures, it can be seen in the figure 2b that the network with 128 hidden neurons reaches quickly to a score of 200, whereas the network with only 8 is not able to solve the task within 2000 learning epochs. This suggests that a higher number of hidden neurons makes it easier to learn the parameters need to find an accurate mapping from the state space to a corresponding value. Yet it should be noted that the time it takes for the network to learn parameters that are able to solve the task is heavily dependent on the random samples that are drawn from the memory buffer. Thus, the plots are slightly different in different trials, and it might happen in some runs the network with 8 neurons does converge in a lower number of iterations. Nevertheless, even within the trial-to-trial variance we could discover that the number of hidden neurons is (to a certain extent) proportional to the time it takes the network to solve the task.

Leaving the number of hidden units constant while increasing the depth of the network does not give such clear results. In a second experimental setup we kept the number of hidden neurons constant at 32, but increased the depth of the network to the following architectures: [32], [16, 8, 8], [8, 8, 8, 8], [16, 8, 4, 2, 2], and [8, 8, 4, 4, 4, 2, 2], where the number in the lists represents the number of units in the particular hidden layer. The Q-values of the initial state as well as the scores are shown in the graphs in figure 3. The results obtained from the experiments are a mixed bag and do not allow a general statement. While the deep network with 5 hidden layers performed better than the shallow network, it converged quicker than the networks with 4 and 7 hidden layers. It should be stressed that these particular results are random and that we face a trial-to-trial variance. Nonetheless, we were not able to deduce any particular trend from the experiments. Hence, we will use a shallow network in the following tasks.

Figure 3: Experimentation with the effect of different parameters on the network complexity.

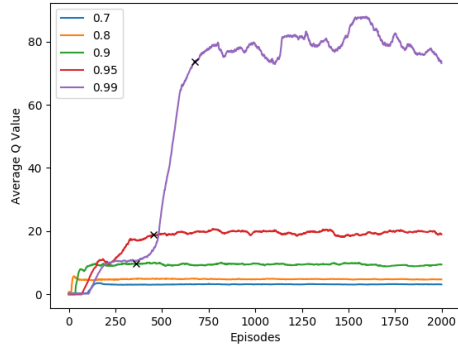


(a) The Q-value plot for different deep architectures.

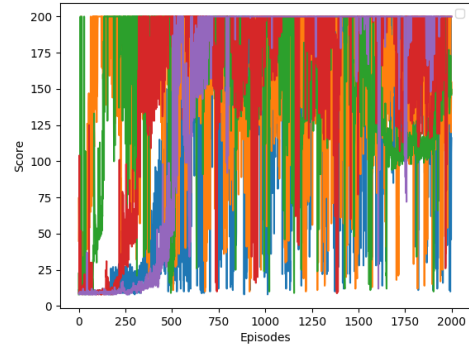


(b) The score plot for different deep architectures.

Figure 4: Effects of discount factor on training

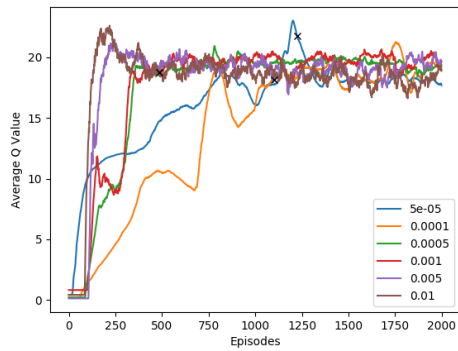


(a) The Q-value plot

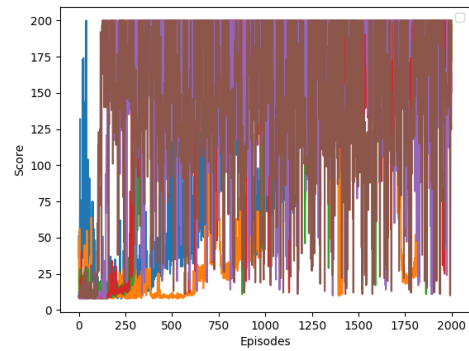


(b) The score plot

Figure 5: Effects of learning rate on training

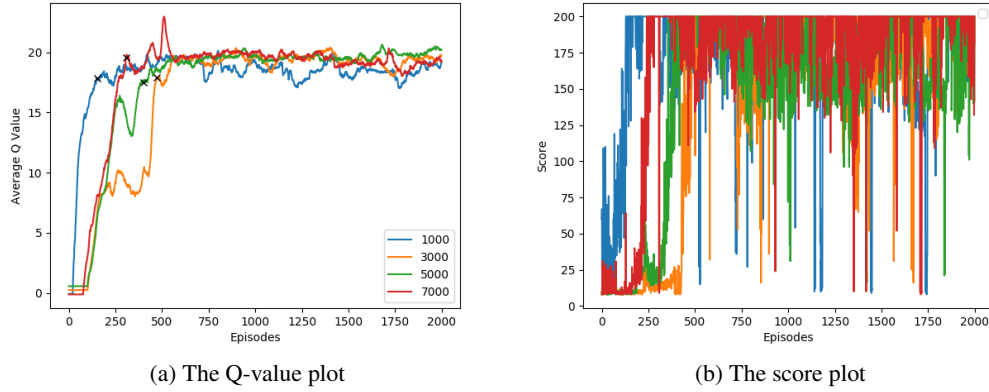


(a) The Q-value plot



(b) The score plot

Figure 6: Effects of memory size on training



- (h) Figures 4, 5, and 6 show experiments examining the effects on training when changing the discount factor, learning rate, and memory size, respectively. For all three experiments, a variety of different parameter values were tried while keeping all other variables constant. The model that was selected had one hidden layer with 128 units, discount factor of 0.95, learning rate of 0.0005, memory size of 5000 and update frequency of 1. Put in other words, each experiment changed one of these variables but left all others unchanged. Each training run lasted 2000 episodes. The check solve flag was turned on, and a black “X” on each of the graphs marks the point where this configuration had solved the problem.

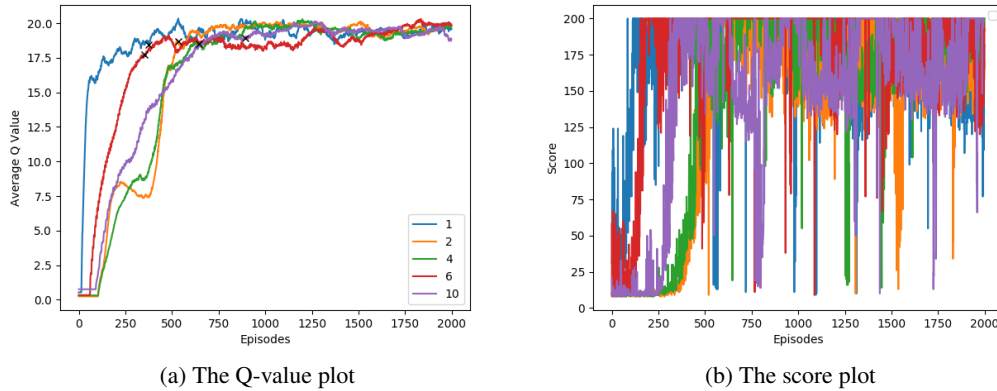
**Discount factor:** For lower values of  $\lambda$  the problem never was solved, in our experiments these values were  $\lambda = 0.7, 0.8$ . The results show that for discount factors which did result in a solved game, the higher the discount factor, the more episodes it took to solve. It is possible that this is because higher discount factor leads to policies which have more weight on future reward. Larger discount factors had better scores on average than smaller discount factors. Also the Q-values for larger values of  $\lambda$  are larger, which intuitively makes sense because the Q-function grows with increasing  $\lambda$ .

**Learning rate:** As one would expect, with higher learning rate, there is much more noise in the average Q values since  $\theta$  is changed more each update. Lower learning rates solved the task later. That is, after experiencing more episodes. On the other hand, high learning rates did not solve the problem in some cases at all.

**Memory size:** The results can be seen in Figure 6. Generally all values used for memory size solved the problem, in fact after around the same number of episodes. The main difference was the swings in the average Q-values where it was observed that smaller memory size had larger swings in average Q-value. One reason for this might be that for larger memory size, the sampling is more “uniform” and thus in the long term average won’t cause the Q-values to fluctuate so much, whereas smaller memory size have more correlated randomness since they are fewer values to sample from and they all happened more recently. This is the case because the memory buffer is like a sliding window. For all the experiments a value of 5000 was chosen because it was in the middle of values that did work, and it worked consistently, and was not as noisy as smaller memory sizes.

- (i) Using the default values for the model as specified in part (h) the update frequency was changed and results/effects on training can be seen in Figure 7. Interestingly, all solve the problem, just after different number of observed episodes. The best seemed to be a frequency of 6. Generally, the lower the value (i.e. higher frequency) the faster change and more noise in the average Q-values, and scores. Oddly, this is not the case for the highest value we tried of 10.
- (j) The parameters chosen for the previous parts do solve the problem, as do many others as made evident by the plots shown. For example, the values which we used as baselines for parts (h) and (i) solve the problem.

Figure 7: Effects of update frequency on training



## 2 Code

The code after completing parts (e) and (f) is appended below.

---

```

1  import sys
2  import gym
3  import pylab
4  import random
5  import numpy as np
6  from collections import deque
7  from keras.layers import Dense
8  from keras.optimizers import Adam
9  from keras.models import Sequential
10
11  EPISODES = 2000 #Maximum number of episodes
12
13  #DQN Agent for the Cartpole
14  #Q function approximation with NN, experience replay, and target
15  ↪ network
16  class DQNAgent:
17      #Constructor for the agent (invoked when DQN is first called
18      ↪ in main)
19      def __init__(self, state_size, action_size):
20          self.check_solve = True #If True, stop if you
21          ↪ satisfy solution condition
22          self.render = False #If you want to see Cartpole
23          ↪ learning, then change to True
24
25          #Get size of state and action
26          self.state_size = state_size
27          self.action_size = action_size
28
29          # Modify here
30
31          #Set hyper parameters for the DQN. Do not adjust those
32          ↪ labeled as Fixed.
33          self.discount_factor = 0.95
34          self.learning_rate = 0.0005
35          self.epsilon = 0.02 #Fixed

```

```

31     self.batch_size = 32 #Fixed
32     self.memory_size = 5000
33     self.train_start = 1000 #Fixed
34     self.target_update_frequency = 1
35
36     #Number of test states for Q value plots
37     self.test_state_no = 10000
38
39     #Create memory buffer using deque
40     self.memory = deque(maxlen=self.memory_size)
41
42     #Create main network and target network (using build_model
43     ↪ defined below)
44     self.model = self.build_model()
45     self.target_model = self.build_model()
46
47     #Initialize target network
48     self.update_target_model()
49
50     #Approximate Q function using Neural Network
51     #State is the input and the Q Values are the output.
52     #####
53     #Edit the Neural Network model here
54     #Tip: Consult
55     ↪ https://keras.io/getting-started/sequential-model-guide/
56     def build_model(self):
57         model = Sequential()
58         model.add(Dense(128, input_dim=self.state_size,
59             ↪ activation='relu',
60                 kernel_initializer='he_uniform'))
61         model.add(Dense(self.action_size, activation='linear',
62             ↪ kernel_initializer='he_uniform'))
63         model.summary()
64         model.compile(loss='mse',
65             ↪ optimizer=Adam(lr=self.learning_rate))
66         return model
67     #####
68
69     #After some time interval update the target model to be same
70     ↪ with model
71     def update_target_model(self):
72         self.target_model.set_weights(self.model.get_weights())
73
74     #Get action from model using epsilon-greedy policy
75     def get_action(self, state):
76     #####
77     #Insert your e-greedy policy code here
78     #Tip 1: Use the random package to generate a random
79     ↪ action.
80     #Tip 2: Use keras.model.predict() to compute Q-values from
81     ↪ the state.
82     action = random.randrange(self.action_size)
83     if random.random() < self.epsilon:
84         return action
85     else:
86         Q_actions = self.model.predict(state)

```

```

83         return np.argmax(Q_actions, axis=1)[0]
84
85     ↪ #####
86 #####
87     #Save sample <s,a,r,s'> to the replay memory
88     def append_sample(self, state, action, reward, next_state,
89         ↪ done):
90         self.memory.append((state, action, reward, next_state,
91             ↪ done)) #Add sample to the end of the list
92
93     #Sample <s,a,r,s'> from replay memory
94     def train_model(self):
95         if len(self.memory) < self.train_start: #Do not train if
96             ↪ not enough memory
97             return
98         batch_size = min(self.batch_size, len(self.memory)) #Train
99             ↪ on at most as many samples as you have in memory
100         mini_batch = random.sample(self.memory, batch_size)
101             ↪ #Uniformly sample the memory buffer
102         #Preallocate network and target network input matrices.
103         update_input = np.zeros((batch_size, self.state_size))
104             ↪ #batch_size by state_size two-dimensional array (not
105             ↪ matrix!)
106         update_target = np.zeros((batch_size, self.state_size))
107             ↪ #Same as above, but used for the target network
108         action, reward, done = [], [], [] #Empty arrays that will
109             ↪ grow dynamically
110
111         for i in range(self.batch_size):
112             update_input[i] = mini_batch[i][0] #Allocate s(i) to
113                 ↪ the network input array from iteration i in the
114                 ↪ batch
115             action.append(mini_batch[i][1]) #Store a(i)
116             reward.append(mini_batch[i][2]) #Store r(i)
117             update_target[i] = mini_batch[i][3] #Allocate s'(i)
118                 ↪ for the target network array from iteration i in
119                 ↪ the batch
120             done.append(mini_batch[i][4]) #Store done(i)
121
122         target = self.model.predict(update_input) #Generate target
123             ↪ values for training the inner loop network using the
124             ↪ network model
125         target_val = self.target_model.predict(update_target)
126             ↪ #Generate the target values for training the outer
127             ↪ loop target network
128
129         #Q Learning: get maximum Q value at s' from target network
130         #####
131         #####
132         #Insert your Q-learning code here
133         #Tip 1: Observe that the Q-values are stored in the
134             ↪ variable target
135         #Tip 2: What is the Q-value of the action taken at the
136             ↪ last state of the episode?
137         for i in range(self.batch_size): #For every batch
138             if done[i]:
139                 target[i][action[i]] = reward[i]
140             else:

```



```

122         target[i][action[i]] = self.discount_factor *
           ↳ np.max(target_val[i]) + reward[i]
123 #####
124 #####
125
126         #Train the inner loop network
127         self.model.fit(update_input, target,
           ↳ batch_size=self.batch_size,
128                   epochs=1, verbose=0)
129         return
130         #Plots the score per episode as well as the maximum q value
           ↳ per episode, averaged over precollected states.
131     def plot_data(self, episodes, scores, max_q_mean):
132         pylab.figure(0)
133         pylab.plot(episodes, max_q_mean, 'b')
134         pylab.xlabel("Episodes")
135         pylab.ylabel("Average Q Value")
136         pylab.savefig("qvalues.png")
137
138         pylab.figure(1)
139         pylab.plot(episodes, scores, 'b')
140         pylab.xlabel("Episodes")
141         pylab.ylabel("Score")
142         pylab.savefig("scores.png")
143
144 #####
145 #####
146
147     if __name__ == "__main__":
148         #For CartPole-v0, maximum episode length is 200
149         env = gym.make('CartPole-v0') #Generate Cartpole-v0
           ↳ environment object from the gym library
150         #Get state and action sizes from the environment
151         state_size = env.observation_space.shape[0]
152         action_size = env.action_space.n
153
154         #Create agent, see the DQNAgent __init__ method for details
155         agent = DQNAgent(state_size, action_size)
156
157         #Collect test states for plotting Q values using uniform
           ↳ random policy
158         test_states = np.zeros((agent.test_state_no, state_size))
159         max_q = np.zeros((EPISODES, agent.test_state_no))
160         max_q_mean = np.zeros((EPISODES, 1))
161
162         done = True
163         for i in range(agent.test_state_no):
164             if done:
165                 done = False
166                 state = env.reset()
167                 state = np.reshape(state, [1, state_size])
168                 test_states[i] = state
169             else:
170                 action = random.randrange(action_size)
171                 next_state, reward, done, info = env.step(action)
172                 next_state = np.reshape(next_state, [1, state_size])
173                 test_states[i] = state
174                 state = next_state
175

```

```

176     scores, episodes = [], [] #Create dynamically growing score
    ↪ and episode counters
177     for e in range(EPISODES):
178         done = False
179         score = 0
180         state = env.reset() #Initialize/reset the environment
181         state = np.reshape(state, [1, state_size]) #Reshape state
    ↪ so that to a 1 by state_size two-dimensional array ie.
    ↪ [x_1,x_2] to [[x_1,x_2]]
182         #Compute Q values for plotting
183         tmp = agent.model.predict(test_states)
184         max_q[e][:] = np.max(tmp, axis=1)
185         max_q_mean[e] = np.mean(max_q[e][:])
186
187         while not done:
188             if agent.render:
189                 env.render() #Show cartpole animation
190
191                 #Get action for the current state and go one step in
    ↪ environment
192                 action = agent.get_action(state)
193                 next_state, reward, done, info = env.step(action)
194                 next_state = np.reshape(next_state, [1, state_size])
    ↪ #Reshape next_state similarly to state
195
196                 #Save sample <s, a, r, s'> to the replay memory
197                 agent.append_sample(state, action, reward, next_state,
    ↪ done)
198                 #Training step
199                 agent.train_model()
200                 score += reward #Store episodic reward
201                 state = next_state #Propagate state
202
203             if done:
204                 #At the end of every episode, update the target
    ↪ network
205                 if e % agent.target_update_frequency == 0:
206                     agent.update_target_model()
207                 #Plot the play time for every episode
208                 scores.append(score)
209                 episodes.append(e)
210
211                 print("episode:", e, " score:", score, "
    ↪ q_value:", max_q_mean[e], " memory length:",
    ↪ len(agent.memory))
212
213                 # if the mean of scores of last 100 episodes is
    ↪ bigger than 195
214                 # stop training
215                 if agent.check_solve:
216                     if np.mean(scores[-min(100, len(scores)):]) >=
    ↪ 195:
217                         print("solved after", e-100, "episodes")
218
    ↪ agent.plot_data(episodes,scores,max_q_mean[:e+1])
219
220                 sys.exit()
221 agent.plot_data(episodes,scores,max_q_mean)

```

---

The same code with additional functions to graph and with our experiments follows:

---

```

1  import os
2  import sys
3  import gym
4  import pylab
5  import random
6  import numpy as np
7  from collections import deque
8  from keras.layers import Dense
9  from keras.optimizers import Adam
10 from keras.models import Sequential
11 from gym import wrappers
12
13 EPISODES = 2000 #Maximum number of episodes
14
15 #DQN Agent for the Cartpole
16 #Q function approximation with NN, experience replay, and target
   ↪ network
17 class DQNAgent:
18     #Constructor for the agent (invoked when DQN is first called
   ↪ in main)
19     def __init__(self, state_size, action_size,
   ↪ target_update_frequency=1, arch=[16],
   ↪ discount_factor=0.95, learning_rate=0.0005,
   ↪ mem_size=5000):
20         self.check_solve = True #If True, stop if you
   ↪ satisfy solution condition
21         self.render = False #If you want to see Cartpole
   ↪ learning, then change to True
22
23         #Get size of state and action
24         self.state_size = state_size
25         self.action_size = action_size
26
27         # Modify here
28
29         #Set hyper parameters for the DQN. Do not adjust those
   ↪ labeled as Fixed.
30         self.discount_factor = discount_factor
31         self.learning_rate = learning_rate
32         self.epsilon = 0.02 #Fixed
33         self.batch_size = 32 #Fixed
34         self.memory_size = mem_size
35         self.train_start = 1000 #Fixed
36         self.target_update_frequency = target_update_frequency
37
38         #Number of test states for Q value plots
39         self.test_state_no = 10000
40
41         #Create memory buffer using deque
42         self.memory = deque(maxlen=self.memory_size)
43
44         self.arch = arch
45         #Create main network and target network (using build_model
   ↪ defined below)
46         self.model = self.build_model()
47         self.target_model = self.build_model()
48
49         #Initialize target network

```

```

50         self.update_target_model()
51
52         #Approximate Q function using Neural Network
53         #State is the input and the Q Values are the output.
54         #####
55         #####
56         #Edit the Neural Network model here
57         #Tip: Consult
58         ↪ https://keras.io/getting-started/sequential-model-guide/
59     def build_model(self):
60         model = Sequential()
61         for num_units in self.arch:
62             model.add(Dense(num_units, input_dim=self.state_size,
63                             ↪ activation='relu',
64                               kernel_initializer='he_uniform'))
65         model.add(Dense(self.action_size, activation='linear',
66                         ↪ kernel_initializer='he_uniform'))
67         model.summary()
68         model.compile(loss='mse',
69                       ↪ optimizer=Adam(lr=self.learning_rate))
70         return model
71         #####
72         #####
73         #After some time interval update the target model to be same
74         ↪ with model
75     def update_target_model(self):
76         self.target_model.set_weights(self.model.get_weights())
77
78         #Get action from model using epsilon-greedy policy
79     def get_action(self, state):
80         #####
81         #####
82         #Insert your e-greedy policy code here
83         #Tip 1: Use the random package to generate a random
84         ↪ action.
85         #Tip 2: Use keras.model.predict() to compute Q-values from
86         ↪ the state.
87         action = random.randrange(self.action_size)
88         if random.random() < self.epsilon:
89             return action
90         else:
91             Q_actions = self.model.predict(state)
92             return np.argmax(Q_actions, axis=1)[0]
93
94         ↪ #####
95         #####
96         #Save sample <s,a,r,s'> to the replay memory
97     def append_sample(self, state, action, reward, next_state,
98                     ↪ done):
99         self.memory.append((state, action, reward, next_state,
100                             ↪ done)) #Add sample to the end of the list
101
102         #Sample <s,a,r,s'> from replay memory
103     def train_model(self):
104         if len(self.memory) < self.train_start: #Do not train if
105             ↪ not enough memory
106             return

```

```

99     batch_size = min(self.batch_size, len(self.memory)) #Train
    ↪ on at most as many samples as you have in memory
100     mini_batch = random.sample(self.memory, batch_size)
    ↪ #Uniformly sample the memory buffer
101     #Preallocate network and target network input matrices.
102     update_input = np.zeros((batch_size, self.state_size))
    ↪ #batch_size by state_size two-dimensional array (not
    ↪ matrix!)
103     update_target = np.zeros((batch_size, self.state_size))
    ↪ #Same as above, but used for the target network
104     action, reward, done = [], [], [] #Empty arrays that will
    ↪ grow dynamically
105
106     for i in range(self.batch_size):
107         update_input[i] = mini_batch[i][0] #Allocate s(i) to
    ↪ the network input array from iteration i in the
    ↪ batch
108         action.append(mini_batch[i][1]) #Store a(i)
109         reward.append(mini_batch[i][2]) #Store r(i)
110         update_target[i] = mini_batch[i][3] #Allocate s'(i)
    ↪ for the target network array from iteration i in
    ↪ the batch
111         done.append(mini_batch[i][4]) #Store done(i)
112
113     target = self.model.predict(update_input) #Generate target
    ↪ values for training the inner loop network using the
    ↪ network model
114     target_val = self.target_model.predict(update_target)
    ↪ #Generate the target values for training the outer
    ↪ loop target network
115
116     #Q Learning: get maximum Q value at s' from target network
117     #####
118     #####
119     #Insert your Q-learning code here
120     #Tip 1: Observe that the Q-values are stored in the
    ↪ variable target
121     #Tip 2: What is the Q-value of the action taken at the
    ↪ last state of the episode?
122     for i in range(self.batch_size): #For every batch
123         if done[i]:
124             target[i][action[i]] = reward[i]
125         else:
126             target[i][action[i]] = self.discount_factor *
    ↪ np.max(target_val[i]) + reward[i]
127     #####
128     #####
129
130     #Train the inner loop network
131     self.model.fit(update_input, target,
    ↪ batch_size=self.batch_size,
132                     epochs=1, verbose=0)
133     return
134     #Plots the score per episode as well as the maximum q value
    ↪ per episode, averaged over precollected states.
135     def plot_data(self, episodes, scores, max_q_mean, dir_name,
    ↪ arch=None):
136
137     if not os.path.exists(dir_name):

```

```

138         os.makedirs(dir_name)
139
140     if arch is None:
141         subname = ''
142     else:
143         subname = str(arch)
144     pylab.figure(0)
145     pylab.plot(episodes, max_q_mean, 'b')
146     pylab.xlabel("Episodes")
147     pylab.ylabel("Average Q Value")
148     pylab.savefig("%s/qvalues-%s.png" % (dir_name, subname))
149
150     pylab.figure(1)
151     pylab.plot(episodes, scores, 'b')
152     pylab.xlabel("Episodes")
153     pylab.ylabel("Score")
154     pylab.savefig("%s/scores-%s.png" % (dir_name, subname))
155
156 def plot_data_multiple(episodes, scores, max_q_mean, solved_times,
157     ↪ dir_name, names):
158
159     if not os.path.exists(dir_name):
160         os.makedirs(dir_name)
161
162     pylab.figure(0)
163     pylab.clf()
164     for i in range(len(episodes)):
165         pylab.plot(episodes[i], max_q_mean[i], label=names[i])
166     actually_solved_times = [time for time in solved_times if time
167     ↪ != -1]
168     solved_values = [max_q_mean[i][solved_times[i]] for i in
169     ↪ range(len(solved_times)) if solved_times[i] != -1]
170     pylab.plot(actually_solved_times, solved_values, 'kx',
171     ↪ label='solved')
172     pylab.xlabel("Episodes")
173     pylab.ylabel("Average Q Value")
174     pylab.legend(names)
175     pylab.savefig("%s/qvalues.png" % dir_name)
176
177     pylab.figure(1)
178     pylab.clf()
179     for i in range(len(episodes)):
180         pylab.plot(episodes[i], scores[i])
181     pylab.xlabel("Episodes")
182     pylab.ylabel("Score")
183     pylab.legend()
184     pylab.savefig("%s/scores.png" % dir_name)
185
186 #####
187 #####
188
189 def simulate(agent, times):
190     env = gym.make('CartPole-v0')
191     env = wrappers.Monitor(env, directory='sims/trained',
192     ↪ force=True)
193
194     for run in range(times):
195         state = env.reset()
196         state_size = env.observation_space.shape[0]

```

```

192     state = np.reshape(state, [1, state_size]) # Reshape
        ↳ state so that to a 1 by state_size two-dimensional
        ↳ array ie. [x_1,x_2] to [[x_1,x_2]]
193
194     done = False
195     while not done:
196         #env.render() # Show cartpole animation
197
198         # Get action for the current state and go one step in
        ↳ environment
199         action = agent.get_action(state)
200         state, reward, done, info = env.step(action)
201         state = np.reshape(state, [1, state_size]) # Reshape
        ↳ next_state similarly to state
202     env.close()
203
204
205 def train(arch, discount_factor=0.95, learning_rate=0.0005,
        ↳ mem_size=5000, update_freq=1):
206     solved = -1
207     #For CartPole-v0, maximum episode length is 200
208     env = gym.make('CartPole-v0') #Generate Cartpole-v0
        ↳ environment object from the gym library
209     #env = wrappers.Monitor(env, directory='sims/training',
        ↳ force=True)
210
211     #Get state and action sizes from the environment
212     state_size = env.observation_space.shape[0]
213     action_size = env.action_space.n
214
215     #Create agent, see the DQNAgent __init__ method for details
216     agent = DQNAgent(state_size, action_size, arch=arch,
        ↳ discount_factor=discount_factor,
        ↳ learning_rate=learning_rate, mem_size=mem_size,
        ↳ target_update_frequency=update_freq)
217
218     #Collect test states for plotting Q values using uniform
        ↳ random policy
219     test_states = np.zeros((agent.test_state_no, state_size))
220     max_q = np.zeros((EPISODES, agent.test_state_no))
221     max_q_mean = np.zeros((EPISODES,1))
222
223     done = True
224     for i in range(agent.test_state_no):
225         if done:
226             done = False
227             state = env.reset()
228             state = np.reshape(state, [1, state_size])
229             test_states[i] = state
230         else:
231             action = random.randrange(action_size)
232             next_state, reward, done, info = env.step(action)
233             next_state = np.reshape(next_state, [1, state_size])
234             test_states[i] = state
235             state = next_state
236
237     scores, episodes = [], [] #Create dynamically growing score
        ↳ and episode counters
238     for e in range(EPISODES):

```

```

239     done = False
240     score = 0
241     state = env.reset() #Initialize/reset the environment
242     state = np.reshape(state, [1, state_size]) #Reshape state
243     ↳ so that to a 1 by state_size two-dimensional array ie.
244     ↳ [x_1,x_2] to [[x_1,x_2]]
245     #Compute Q values for plotting
246     tmp = agent.model.predict(test_states)
247     max_q[e][:] = np.max(tmp, axis=1)
248     max_q_mean[e] = np.mean(max_q[e][:])
249
250     while not done:
251         if agent.render:
252             env.render() #Show cartpole animation
253
254         #Get action for the current state and go one step in
255         ↳ environment
256         action = agent.get_action(state)
257         next_state, reward, done, info = env.step(action)
258         next_state = np.reshape(next_state, [1, state_size])
259         ↳ #Reshape next_state similarly to state
260
261         #Save sample <s, a, r, s'> to the replay memory
262         agent.append_sample(state, action, reward, next_state,
263             ↳ done)
264         #Training step
265         agent.train_model()
266         score += reward #Store episodic reward
267         state = next_state #Propagate state
268
269         if done:
270             #At the end of every episodesepisode, update the
271             ↳ target network
272             if e % agent.target_update_frequency == 0:
273                 agent.update_target_model()
274             #Plot the play time for every episode
275             scores.append(score)
276             episodes.append(e)
277
278             if e % 100 == 0:
279                 print("episode:", e, " score:", score, "
280                     ↳ q_value:", max_q_mean[e], " memory
281                     ↳ length:",
282                     len(agent.memory))
283
284             # if the mean of scores of last 100 episodes is
285             ↳ bigger than 195
286             # stop training
287             if agent.check_solve and solved == -1:
288                 if np.mean(scores[-min(100, len(scores))]) >=
289                     195:
290                     print("solved after", e-100, "episodes")
291                     solved = e-100
292                     #
293                     ↳ agent.plot_data(episodes,scores,max_q_mean[:e+1], 'part-g
294                     env.close()
295                     simulate(agent, 3)
296                     return episodes, scores, max_q_mean,
297                     ↳ solved

```



```

286     env.close()
287     #agent.plot_data(epsisodes,scores,max_q_mean, '')
288     return episodes, scores, max_q_mean, solved
289
290
291 if __name__ == '__main__':
292     episodes = []
293     scores = []
294     max_q_means = []
295     # num_unit_values = [16, 32, 64]
296     # archs = [[8, 32], [16, 32, 32], [16, 32, 32, 32]]
297     # archs = [[8], [16], [32], [64], [128]]
298     # archs = [[8]]
299     #names = list(map(str, archs))
300
301     # discount factor tests
302     # discount_factors = [0.7, 0.8, 0.9, 0.95, 0.99]
303     # names = list(map(str, discount_factors))
304     arch = [128]
305     # learning_rates = [0.00005, 0.0001, 0.0005, 0.001, 0.005,
306     ↪ 0.01]
307     # names = list(map(str, learning_rates))
308     # mem_sizes = [500, 1000, 5000, 9000]
309     # names = list(map(str, mem_sizes))
310
311     # update_freqs = [1, 2, 4, 6, 10]
312     # names = list(map(str, update_freqs))
313     #
314     # solved_times = []
315     #
316     # for update_freq in update_freqs:
317     #     eps, score, qs, solved = train(arch,
318     ↪ discount_factor=0.95, learning_rate=0.0005, mem_size=5000,
319     ↪ update_freq=update_freq)
320     #     episodes.append(eps)
321     #     scores.append(score)
322     #     max_q_means.append(qs[:len(score)])
323     #     solved_times.append(solved)
324     # plot_data_multiple(episodes, scores, max_q_means,
325     ↪ solved_times, 'update-freq', names)
326
327 train([128], discount_factor=0.95, learning_rate=0.0005,
328     ↪ mem_size=5000, update_freq=1)

```

---