
EL2805 Reinforcement Learning

Alice Karnsund
Therese Stålhandske

Abstract

1 This is the Lab Assignment 2 for the course Reinforcement Learning EL2805.

2 1 Formulation of the RL problem

3 Our goal in this lab is to train an agent that finds an optimal policy that prevents the pole on the cart
4 from falling over, see Figure 1. To begin with, we formulate the RL problem, i.e. the state-space,
5 action-space etc.

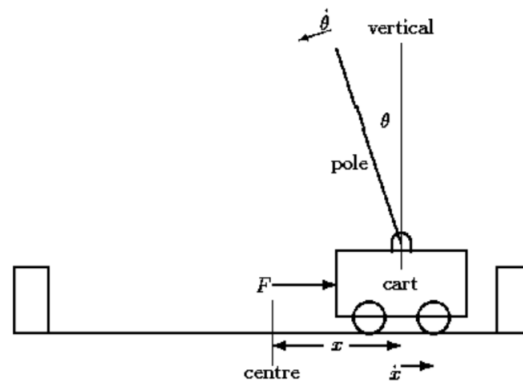


Figure 1: The cart and pole system.

6 Each state is described by four observations made at each time step, the cart position, cart velocity,
7 the pole angle and the pole velocity at the tip of the pole. In Table 1, these observations are described
8 along with their min and max values.

Observation	Min	Max
Cart position, x	-4.8	4.8
Cart velocity, \dot{x}	-inf	inf
Pole Angle, θ	-24°	24°
Pole velocity, $\dot{\theta}$	-inf	inf

Table 1: States

9 The action space constitutes of two actions, either push the cart to the left or to the right, with a
10 magnitude of 1N, to balance the pole.

$$A = \{L = \text{push to the left}, R = \text{push to the right}\}$$

11 The agent can receive two different rewards, which depends on the pole angle θ and the cart position
 12 x . These are represented in Table 2.

State s_t	Reward r_t
$\theta < 12.5^\circ$ and $x < 2,4m$	1
Otherwise	0

Table 2: Rewards

13 Furthermore, we are considering an episodic task. An episode terminates after 200 time steps or once
 14 the agent receives 0 rewards for the first time.

15 For this problem, we cannot use standard methods to solve the problem efficiently. This is due to the
 16 fact that we are dealing with an infinite state space, as can be seen from Table 1.

17 2 Code description

18 The file `cartpole_dqn.py` contains one class, `DQNAgent`, that has seven functions in it, and a
 19 `__main__` method.

20 `__main__`

21 First, we create the cart-pole environment with the use of the `gym`-library. Thereafter we get the
 22 state and action sizes defined by the environment. We create a new agent, an instance of the class
 23 `DQNAgent`. The first for-loop collects a set of test states for plotting the Q values using uniform
 24 random policy. For the first iteration, `done = True`, and we will enter the if-statement and there set
 25 the state to the initial one. Thereafter we will enter the else-statement, until `done` is change, where
 26 we will select a random action and then observe the new state reward, `done` and `info`, after taking this
 27 action. The `done` variable will be set to `True` again if we reach the end of an episode or if the pole
 28 has tipped too far for example.

29 The second for-loop goes over all the episodes, collecting cumulative scores for each episode. For
 30 every episode, we start in the initial state. As long as the episode is not over, in each time step we
 31 select an action, line 3 in Algorithm 1, and collect a sample (s_t, a_t, r_t, s'_t) that is added to the replay
 32 buffer, lines 4-5 in Algorithm 1. Thereafter we train the agent, using `train_model` and add the
 33 reward to this episodes total reward, lines 6-12 in Algorithm 1. If the episode is over, `done` is set to
 34 `True` and we update the target network and collect this episode's cumulative reward to the score list,
 35 line 14 in Algorithm 1. Finally, if we have reached a satisfying solution we plot the solution using
 36 `plot_data` and then terminate the algorithm. This corresponds to lines 15-17 in Algorithm 1. If we
 37 do not reach a satisfying solution we will simply plot what we have.

38 Behavior of each function

39 The class `DQNAgent` has seven functions in it.

40 `__init__(self, state_size, action_size)`

41 First of all the constructor takes as arguments the state size, i.e. the number of elements representing
 42 a state, and the number of actions, in this case, 2. The class holds 16 instance variables. The
 43 first two are boolean variables, where the first one is `True` if we have reached the solution condi-
 44 tion. The other one simply allows us to see the simulation of the learning. The next two variables
 45 hold the state and action size, as given by the arguments. Thereafter we have seven instance vari-
 46 ables that set the hyperparameters for the DQN. `discount_factor`, `learning_rate`, `epsilon`,
 47 `batch_size`, `memory_size`, `train_start` and `target_update`. Here `memory_size` is the
 48 size of the replay buffer and `memory_start` simply assures that we do not start training if we do
 49 not have enough memory. The last five, `test_state_no`, `memory`, `model` and `target_model`
 50 together with `update_target_model()`, these sets the the number of test states for the Q-value
 51 plots, creates a replay buffer and creates two different ANNs for the two Q-functions, Q_ϕ and $Q_{\phi'}$.

52 `build_model(self)`

53 The function `build_model` returns an Artificial Neural Network with one hidden and one output

54 layer, which makes use of a ReLu activation function and a linear activation function, respectively.
55 The ANN is compiled using the mean square error and the stochastic optimization method Adam.

56 `update_target_model(self)`
57 This function updates the weights for the target model, $Q_{\phi'}$, to be the same as those for the model
58 Q_{ϕ} . Q_{ϕ} is the function approximation of the (state, action) value function of a given policy π , Q^{π} .

59 `get_action(self, state)`
60 In this function we provide an ϵ -greedy policy that will, given a state, return an action. This
61 corresponds to line 3 in Algorithm 1.

62 `append_sample(self, state, action, reward, next_state, done)`
63 This function takes care of adding samples to the replay buffer. Where each sample holds the
64 observations, (s_t, a_t, r_t, s'_t) . The replay buffer is of fixed size, therefore old experience will be
65 removed as new is observed. This corresponds to line 5 in Algorithm 1.

66 `train_model(self)`
67 This function is responsible for the training part and is also where we will provide our Q-learning
68 code. First the function makes sure that we have enough memory to start training, i.e. we must have
69 a full replay buffer. We decide upon a `batch_size` that is at most equal to the size of the replay
70 buffer. Then we sample `batch_size` samples from our replay buffer into `mini_batch`, which will
71 be used for training. Now `mini_batch` contains samples of the form (s_i, a_i, r_i, s'_i) . We divide this
72 information into two arrays, `update_input` that holds all s_i and `update_target` that holds all s'_i .
73 And two lists, `action` that holds all a_i and `reward` that holds all r_i . Then we predict an output
74 for the Q_{ϕ} network, using `update_input` as input, and another output for the $Q_{\phi'}$ network, using
75 `update_target` as input. Then after providing our Q-learning code with the target values, we will
76 train the Q_{ϕ} . This corresponds to lines 6-11 in Algorithm 1.

77 `plot_data(self, episodes, scores, max_q_mean)`
78 This function creates two plots, one showing the average Q-value over episodes and the other one
79 showing the score over episodes.

80 3 Neural Network description

81 Keras is a neural network API that is able to run together with TensorFlow, CNTK, or Theano. It
82 offers easy and fast modeling with Artificial Neural Networks.

83 For this task, we are using a sequential model to build our network. This means that we are using a
84 linear stack of layers. To stack layers one simply uses the `.add()` function. In our case, we are using
85 two layers, one hidden layer, and one output layer. Note that the counting index starts with the first
86 hidden layer up to the output layer, thus the "input layer" is not counted.

87 The hidden layer is a Dense layer, i.e. a fully connected layer, and takes as input a vector of dimension
88 4, i.e. the `state_size`. The number 16 is the dimensionality of the output space and means that the
89 layer has 16 hidden units. Furthermore, the hidden layer makes use of an ReLu activation function,
90 $f(x) = \max(0, x)$ when computing the output of each neuron. Lastly, the initial weight values
91 are set according to "he_uniform". That is, weights are drawn from a uniform distribution within
92 $[-limit, limit]$ where the limit is $\sqrt{(6/fan_in)}$, where `fan_in` is the number of input units in the
93 weight tensor¹.

94 The output layer is also a fully connected layer, which takes as inputs the outputs of the hidden layer,
95 and therefore has an input dimension of 16. The output layer consists of 2 units, i.e. `action_size`.
96 Thus the output of the network is an array holding the Q-values for each action of the state that was
97 given as input to the network. This layer makes use of a linear activation function when calculating
98 the output of each neuron, $f(x) = x$. And the initial weights are set in the same way as those for the
99 hidden layer.

100 The whole model makes use of the stochastic optimization method Adam and aims to minimize the
101 Mean Least Square error, both given as an argument to `.compiler()`. Where the `.compiler()`
102 configures the model for training.

¹<https://keras.io/initializers/>

103 4 Modification of code

104 `get_action(self, state)`
 105 DQN is an off-policy algorithm. This means that we evaluate or improve a policy different from that
 106 used to generate the data. Here we generate data w.r.t Q_ϕ , but the target policy is w.r.t $Q_{\phi'}$. Due to
 107 this ϵ can be fixed and do not have to decrease over time, in this case, $\epsilon = 0.02$, always. We make use
 108 of an ϵ -greedy behavior policy,

$$a = \begin{cases} \operatorname{argmax}_b Q_\phi^{(t)}(s, b), & \text{w.p } 1 - \epsilon \\ \operatorname{uniform}(A_s), & \text{w.p } \epsilon \end{cases} \quad (1)$$

109 `train_model(self)`
 110 In this part we provided the targets, y_i for the Q_ϕ network.

$$y_i = \begin{cases} r_i & \text{if episode stops in } s'_i \\ r_i + \lambda \max_a Q_{\phi'}(s'_i, a) & \text{otherwise} \end{cases} \quad (2)$$

111 These are then inserted into the `target` vector, where y_i is set as target for $Q_\phi(s_i, a_i)$. Here a_i is the
 112 action taken in s_i .

113 5 Assessment of model

N hidden layers	1
N neurons	16
Discount factor	0.95
Learning rate	0.005
Memory size	1000

Table 3: Default Hyper-parameter settings

114 The default hyper-parameters, shown in Table 3, are used for the learning procedure demonstrated in
 115 Figure 2. It seems like the agent are able to learn how to keep the pole up, but it is very inconsistent
 116 having large score fluctuations over the run episodes. Clearly, these set of hyper-parameters does not
 117 solve the problem. It was now our task to extrapolate these hyper-parameters in order to improve the
 118 learning performance.

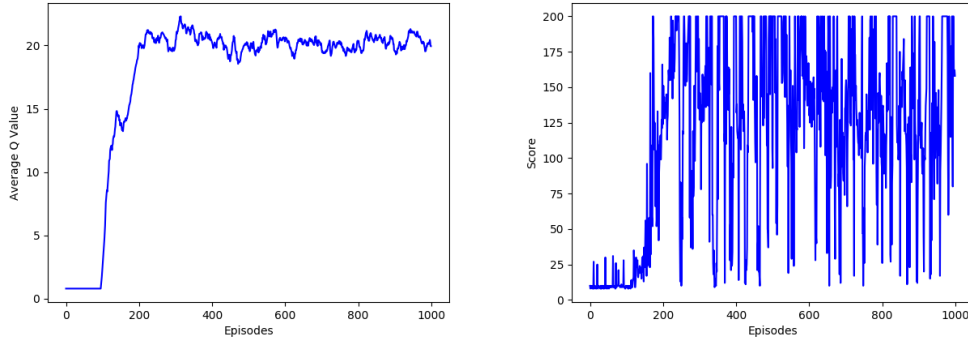


Figure 2: Plot of the average Q-value over episodes (left) and plot of the score over episodes (right), using the default hyper-parameters.

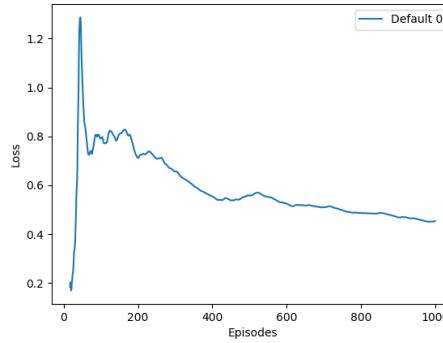


Figure 3: Plot over the average loss function using the default hyper-parameters.

First, we looked at increasing the number of hidden layers to 2 and 3, with different number of hidden neurons in each layer. We also tried different sizes of the replay buffer together with this. We noticed that an overall general solution was best found using one layer with an increased amount of neurons. We found that 84 neurons gave good result after a good amount of trial and error, thus we decided to use this and base further investigations on this setting. The number of neurons is strongly dependent on the learning rate and the discount factor of the learning rate.

6 Investigation

We investigate the three hyper-parameters learning rate, discount factor and memory size and their respective effect on the learning performance. This is done with three representative simulations, varying one parameter at the time, while keeping the other two fixed. Further, we also investigate how the target update frequency effect the training. All the following investigations are based on a model with one hidden layer and 84 neurons.

6.1 Learning Rate

The learning rate is a measure of how far to move the weights in the direction of the gradient for a mini batch. Our intuition then tells us that a low learning rate will give a more reliable, but also slower, learning process. If the learning rate is too high, the training might not even converge as it can overshoot the minimum as the steps are to large.

The loss functions for three different learning rates are displayed below in figure 5. In this graph, the differences between a low and a high learning rate are clearly shown. When the learning rate $= 0.0005$, we see a bad performance in early episodes, where the loss function dramatically increases. But as the number of episodes increases, the loss function starts to steadily decrease. For learning rate $= 0.05$, we see the opposite behavior, the loss function starts low, but when about 400 episodes have elapsed, the loss functions start increasing, indicating a poor learning performance. Learning rate 0.005 , have a behavior somewhere in between and it shows a steady loss function over the first 1000 episodes.

In Figure 4, we see the score and the average q values for each of the learning rates. Considering the score plot, we see that the graph representing learning rate $= 0.0005$, shows the most stable behavior, even if it takes a little more than 400 episodes getting to a score of 200, which in comparisons is quite slow.

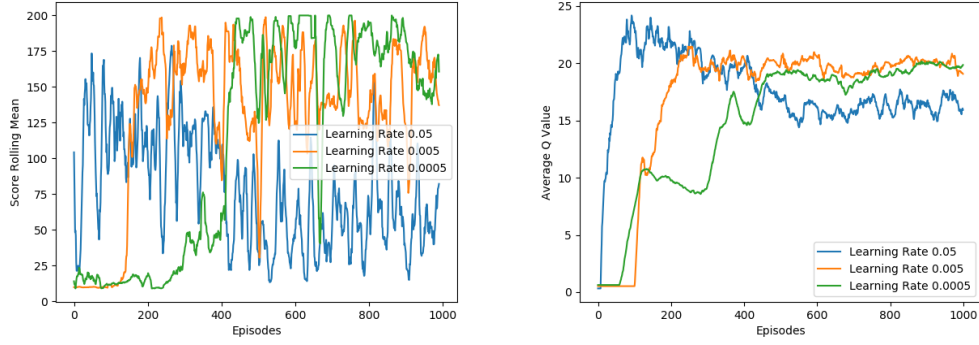


Figure 4: (Left) Rolling mean of the scores per episode for different learning rates. (Right) The Average Q-value per episode for different learning rates.

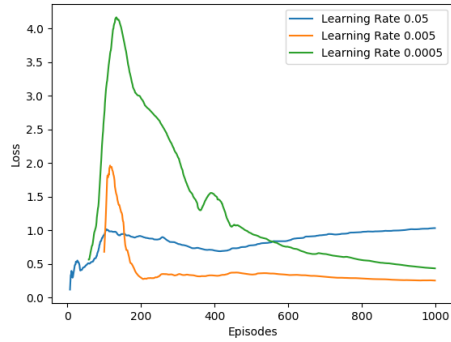


Figure 5: Mean loss function for different learning rates.

148 6.2 Discount Factor

149 The discount factor regulates the length of the time horizon in which that the agent should consider
 150 the rewards. A low discount factor premiers immediate rewards, while a higher discount factor also
 151 considers rewards further into the future.

152 Different discount factors are demonstrated in the figures below. With a discount factor closer to 1,
 153 0.99, we can see that the average Q-value, in Figure 6 to the right, strongly increases and do not seem
 154 to converge at the same pace as the lower discount factors 0.90 and 0.95. This also leads to larger
 155 fluctuation in the score plot, shown in Figure 6 to the left. Looking at loss function, in figure 7, the
 156 loss for a discount rate of 0.99 increases together with the number of episodes. This is clearly not a
 157 desirable behavior.

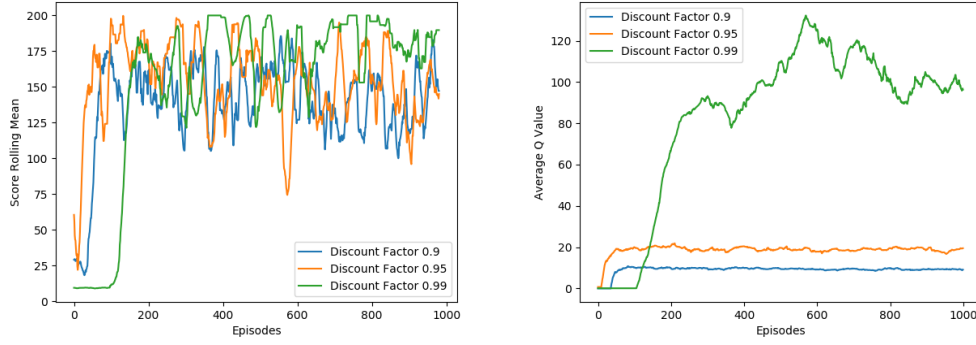


Figure 6: (Left) Rolling mean of the scores per episode for different discount factors. (Right) The Average Q-value per episode for different discount factors.

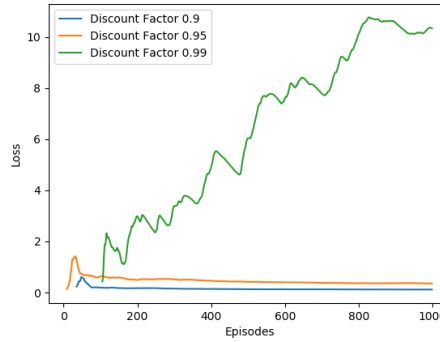


Figure 7: Mean loss function for different discount factors.

158 6.3 Memory Size

159 The size of the memory determines the number of experiences that the agent will base its next action
 160 on. The size of the optimal memory depends, among other things, on the variability of the training
 161 samples. For simulation purposes, we try to increase the memory with an exponential of 10. The
 162 most distinguishable difference between the different memory sizes, is that when the memory size
 163 is smaller. Then the learning scores, shown in 8, have higher fluctuation. The scores have a higher
 164 tendency to drop, even in later episodes. This might indicate that, in this case, a larger memory size
 165 is preferred.

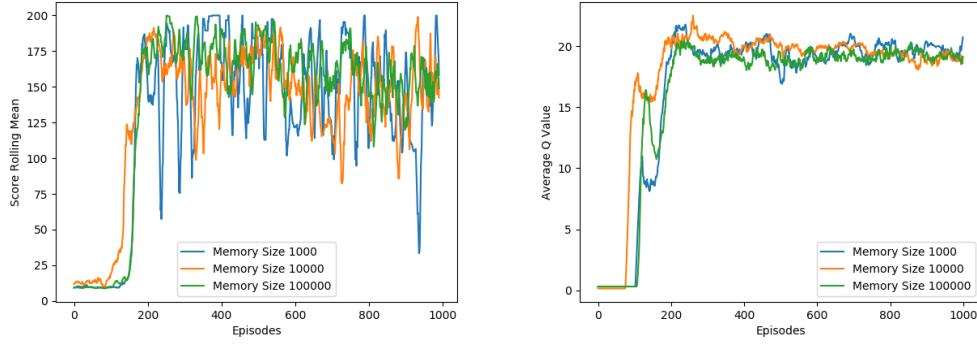


Figure 8: (Left) Rolling mean of the scores per episode for different memory sizes. (Right) The Average Q-value per episode for different memory sizes.

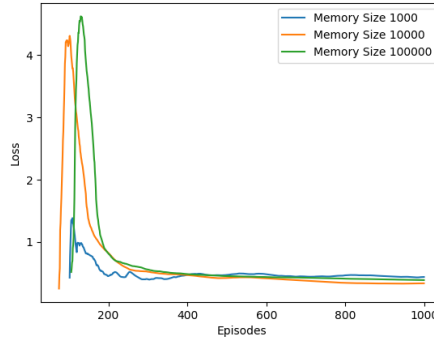


Figure 9: Average loss function for each episode for different memory sizes.

166 6.4 Target Update Frequency

167 The target update frequency controls how often the target network is updated. In Figure 10, to the
 168 left, we see the scores over the episodes. Naturally, when the target update frequency = 100, it will
 169 take a larger number of episodes to learn the true behavior, as the updates are not as frequent. An
 170 interesting thing though, can be seen in figure 11, where the loss function decreases rather fast for a
 171 frequency of = 100.

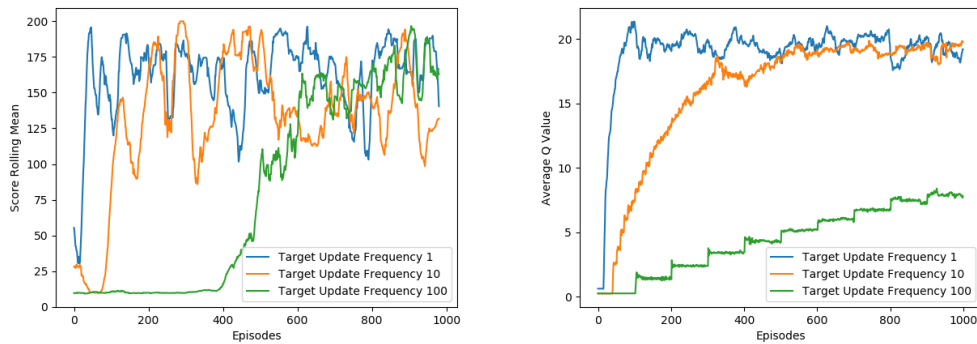


Figure 10: (Left) Rolling mean of the scores per episode for different target updates frequencies. (Right) The Average Q-value per episode for different target update frequencies.

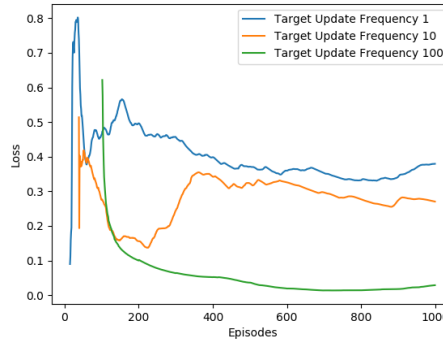


Figure 11: Average loss function for each episode for different target update frequencies.

6.5 Extension

In the extension, we try to further improve the learning performance of the model. First, this is done by introducing a decaying learning rate. This means that in the start of the learning process, when we are supposedly far from a solution, the agent will take larger steps towards the solution, decreasing the learning time. When we are getting closer to the minimum, the steps become smaller, increasing the precision of the algorithm. When introducing a decay parameter, we will alter the learning rate, shifting it upwards to counter the decay.

Further, we investigate different activation function for the hidden layer in the neural network. If several of the activation's gets below zero, many of the neurons in the network will "die" and prohibit learning, a phenomenon commonly referred to as "the dying ReLU problem"². To cope with this we tried using a sigmoid and tanh activation function instead. Since we are not using multiple layers, the problem with vanishing gradients won't be significant.

7 Conclusion

Finding an optimal model that solves the CartPole problem has constitute partly of a trail-and-error procedure. This combined with the experience gathered during the investigation section of this lab, has allowed us to choose an appropriate set of hyper-parameters and model.

N hidden layers	1
N neurons	84
Discount factor	0.95
Learning rate	0.03
Decay Rate	0.0008
Memory size	10000
Activation function	Sigmoid
Update frequency	1

Table 4: Chosen hyper-parameter settings

Our best model solved the problem after 127 episodes with the chosen hyper-parameters shown in Table 4.

²[https://en.wikipedia.org/wiki/Rectifier_\(neural_networks\)](https://en.wikipedia.org/wiki/Rectifier_(neural_networks))

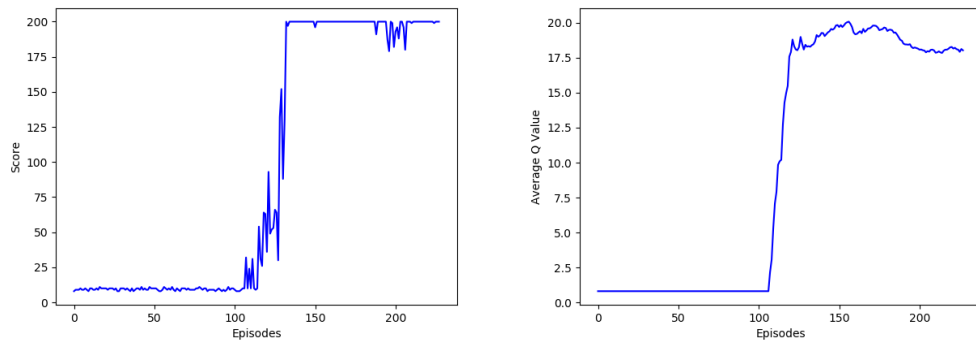


Figure 12: (Left) Scores per episode for our chosen model. (Right) The Average Q-value per episode for our chosen model

190 8 Appendix A

```

191 import sys
192 import gym
193 import pylab
194 import random
195 import numpy as np
196 from collections import deque
197 import keras
198 from keras.layers import Dense
199 from keras.optimizers import Adam
200 from keras.models import Sequential
201 from keras.callbacks import History
202
203 history = History()
204
205
206 EPISODES = 1000 # Maximum number of episodes
207
208 # PLOT VARIABLES
209
210 nr_layers = 1
211
212 #Name_of_plot = "Learning Rate"
213 #Name_of_plot = "Discount Factor"
214 #Name_of_plot = "Memory Size"
215 #Name_of_plot = "Target Update Frequency"
216 #Name_of_plot = "Default"
217 #Name_of_plot = "Record"
218 Name_of_plot = "Number_of_Neurons"
219
220 learning_rate = 0.03
221 discount_factor = 0.95
222 memory_size = 10000
223 target_update_frequency = 1
224 nodes = 100
225
226 # Iteration parameters
227 #parameter_variable = [0.05, 0.005, 0.0005]
228 #parameter_variable = [0.90, 0.95, 0.99]
229 #parameter_variable = [1000,10000,100000]
230 #parameter_variable = [1,10,100]
231 parameter_variable = [24,48,84,100]
232
233 parameter_variable = [100]
234
235 # DQN Agent for the Cartpole
236
237 # Q function approximation with NN, experience replay, and target network
238 class DQNAgent:
239     # Constructor for the agent (invoked when DQN is first called in main)
240     def __init__(self, state_size, action_size, nodes_first, learning_rate, discount):
241         self.check_solve = True # If True, stop if you satisfy solution condition
242         self.render = False # If you want to see Cartpole learning, then change to True
243
244         # Get size of state and action
245         self.state_size = state_size
246         self.action_size = action_size
247

```

```

248 #####
249 #####
250 # Set hyper parameters for the DQN. Do not adjust those labeled as Fixed.
251
252 self.discount_factor = discount_factor # lambda
253 self.learning_rate = learning_rate
254 self.memory_size = memory_size
255
256 self.epsilon = 0.02 # Fixed
257 self.batch_size = 32 # Fixed
258 self.train_start = 1000 # Fixed
259
260 self.target_update_frequency = target_update_frequency
261
262 # New added hyper parameteres
263 self.nodes_first_hidden = nodes_first
264 #####
265 #####
266
267 # Number of test states for Q value plots
268 self.test_state_no = 10000
269
270 # Create memory buffer using deque
271 self.memory = deque(maxlen=self.memory_size)
272
273 # Create main network and target network (using build_model defined below)
274 self.model = self.build_model()
275 self.target_model = self.build_model()
276
277 self.losses = []
278
279 # Initialize target network
280 self.update_target_model()
281
282 # Approximate Q function using Neural Network
283 # State is the input and the Q Values are the output.
284 #####
285 #####
286 # Edit the Neural Network model here
287 # Tip: Consult https://keras.io/getting-started/sequential-model-guide/
288 def build_model(self):
289     model = Sequential()
290     model.add(Dense(self.nodes_first_hidden, input_dim=self.state_size, activation='relu',
291                     kernel_initializer='he_uniform'))
292     #model.add(Dense(24, activation='relu',
293     #                kernel_initializer='he_uniform'))
294     model.add(Dense(self.action_size, activation='linear',
295                     kernel_initializer='he_uniform'))
296     model.summary()
297     model.compile(loss='mse', optimizer=Adam(lr=self.learning_rate, decay=0.0008))
298     return model
299
300 #####
301 #####
302
303 # After some time interval update the target model to be same with model
304 def update_target_model(self):
305     self.target_model.set_weights(self.model.get_weights())
306

```

```

307 # Get action from model using epsilon-greedy policy
308 def get_action(self, state):
309     #####
310     #####
311     # Insert your e-greedy policy code here
312     # Tip 1: Use the random package to generate a random action.
313     # Tip 2: Use keras.model.predict() to compute Q-values from the state.
314     if random.uniform(0, 1) < self.epsilon:
315         action = random.randrange(self.action_size)
316     else:
317         Q_state = self.model.predict(state)
318         action = np.argmax(Q_state)
319     return action
320
321     #####
322     #####
323     # Save sample <s,a,r,s'> to the replay memory
324     def append_sample(self, state, action, reward, next_state, done):
325         self.memory.append((state, action, reward, next_state, done))
326 # Add sample to the end of the list
327
328 # Sample <s,a,r,s'> from replay memory
329 def train_model(self):
330     if len(self.memory) < self.train_start:
331 # Do not train if not enough memory
332         return
333         batch_size = min(self.batch_size, len(self.memory))
334 # Train on at most as many samples as you have in memory
335         mini_batch = random.sample(self.memory, batch_size)
336 # Uniformly sample the memory buffer
337         # Preallocate network and target network input matrices.
338         update_input = np.zeros(
339             (batch_size, self.state_size)) # batch_size by state_size two-dimension
340         update_target = np.zeros((batch_size, self.state_size))
341 # Same as above, but used for the target network
342         reward = np.zeros((batch_size, 1))
343         action, done = [], [] # Empty arrays that will grow dynamically
344
345         for i in range(self.batch_size):
346             update_input[i] = mini_batch[i][0]
347 # Allocate s(i) to the network input array from iteration i in the batch
348             action.append(mini_batch[i][1]) # Store a(i)
349             # reward.append(mini_batch[i][2]) #Store r(i)
350             reward[i] = mini_batch[i][2]
351             update_target[i] = mini_batch[i][
352                 3] # Allocate s'(i) for the target network array from iteration i in
353             done.append(mini_batch[i][4]) # Store done(i)
354
355         target = self.model.predict(
356             update_input) # Generate target values for training the inner loop network
357         target_val = self.target_model.predict(
358             update_target) # Generate the target values for training the outer loop
359
360         # Q Learning: get maximum Q value at s' from target network
361         #####
362         #####
363         # Insert your Q-learning code here
364         # Tip 1: Observe that the Q-values are stored in the variable target
365         # Tip 2: What is the Q-value of the action taken at the last state of the episode

```

```

366
367     max_Qtarget = np.reshape(np.max(target_val , axis=1), (batch_size , 1))
368     max_Qtarget[done] = 0 # Q-value of the action taken at the last state of the
369     y = reward + self.discount_factor * max_Qtarget
370
371     # Line 10 in Algo. 1
372
373     for i in range(self.batch_size): # For every batch
374         target[i][action[i]] = y[i] # random.randint(0,1)
375     # print(target)
376     #####
377     #####
378
379     # Train the inner loop network
380     history = self.model.fit(update_input , target , batch_size=self.batch_size ,
381                             epochs=1, verbose=0)
382
383     self.losses.append(history.history[ 'loss' ][0])
384
385     return
386
387
388
389     # Plots the score per episode as well as the maximum q value per episode , averaged
390     def plot_data(self , episodes , scores , max_q_mean):
391         pylab.figure(0)
392         pylab.plot(episodes , max_q_mean , 'b')
393         pylab.xlabel("Episodes")
394         pylab.ylabel("Average_Q_Value")
395         # pylab.savefig("qvalues.png")
396
397         pylab.figure(1)
398         pylab.plot(episodes , scores , 'b')
399         pylab.xlabel("Episodes")
400         pylab.ylabel("Score")
401         # pylab.savefig("scores.png")
402         pylab.show()
403
404
405     def multi_plot(episodes , all_scores , all_max_q_mean , all_loss , agent):
406
407         pylab.figure(0)
408
409         discount_factor = agent.discount_factor
410         learning_rate = agent.learning_rate
411         memory_size = agent.memory_size
412         nodes = agent.nodes_first_hidden
413
414         # Rolling Average
415         N = 20
416
417
418         for i in range(np.shape(all_scores)[0]):
419             cumsum = np.cumsum(np.insert(all_scores[i] , 0, 0))
420             moving_avg = (cumsum[N:] - cumsum[:-N]) / float(N)
421             pylab.plot(moving_avg , label= Name_of_plot + "_" + str(parameter_variable[i]
422             pylab.xlabel("Episodes")
423             pylab.ylabel("Score_Rolling_Mean")
424             pylab.legend()

```

```

425     pylab.savefig("C:/Users/Therese/PycharmProjects/ReinforcementLearningLAB1/plots/%
426     pylab.show()
427
428
429     for i in range(np.shape(all_scores)[0]):
430         pylab.plot(episodes, all_max_q_mean[i][:], label= Name_of_plot + "_" + str(p
431     pylab.xlabel("Episodes")
432     pylab.ylabel("Average_Q_Value")
433     pylab.legend()
434     pylab.savefig("C:/Users/Therese/PycharmProjects/ReinforcementLearningLAB1/plots/%
435     pylab.show()
436
437     pylab.figure(1)
438     for i in range(np.shape(all_scores)[0]):
439         pylab.plot(episodes, all_scores[i][:], label= Name_of_plot + "_" +
440 str(parameter_variable[i]))
441     pylab.xlabel("Episodes")
442     pylab.ylabel("Score")
443     pylab.legend()
444     pylab.savefig("C:/Users/Therese/PycharmProjects/ReinforcementLearningLAB1/plots/%
445     pylab.show()
446
447     pylab.figure(2)
448     for i in range(np.shape(all_scores)[0]):
449         pylab.plot(episodes, all_loss[i][:], label= Name_of_plot + "_" +
450 str(parameter_variable[i]))
451     pylab.xlabel("Episodes")
452     pylab.ylabel("Loss")
453     pylab.legend()
454     pylab.savefig("C:/Users/Therese/PycharmProjects/ReinforcementLearningLAB1/plots/%
455     pylab.show()
456
457
458
459 if __name__ == "__main__":
460     all_scores = np.zeros((len(parameter_variable), EPISODES))
461     all_max_q_mean = np.zeros((len(parameter_variable), EPISODES))
462     all_losses = np.zeros((len(parameter_variable), EPISODES))
463
464     for n in range(len(parameter_variable)):
465
466         # For CartPole-v0, maximum episode length is 200
467         env = gym.make('CartPole-v0') # Generate Cartpole-v0 environment object from
468         # Get state and action sizes from the environment
469         state_size = env.observation_space.shape[0]
470
471         action_size = env.action_space.n
472
473         # Setting iteration variable (Mainly for the plots)
474         if Name_of_plot == "Learning_Rate":
475             learning_rate = parameter_variable[n]
476
477         if Name_of_plot == "Memory_Size":
478             memory_size= parameter_variable[n]
479
480         if Name_of_plot == "Discount_Factor":
481             discount_factor = parameter_variable[n]
482
483         if Name_of_plot == "Target_Update_Frequency":

```

```

484         target_update_frequency = parameter_variable[n]
485
486     if Name_of_plot == "Number_of_Neurons":
487         nodes = parameter_variable[n]
488
489     # Create agent, see the DQNAgent __init__ method for details
490     agent = DQNAgent(state_size, action_size, nodes, learning_rate, discount_factor)
491
492
493     print(agent.learning_rate, agent.discount_factor, agent.memory_size)
494     #print("Learning Rate %f" %agent.learning_rate)
495
496     # agent.nodes_first_hidden = nodes[n]
497
498     # Collect test states for plotting Q values using uniform random policy
499     test_states = np.zeros((agent.test_state_no, state_size))
500     max_q = np.zeros((EPISODES, agent.test_state_no))
501     max_q_mean = np.zeros((EPISODES, 1))
502
503     done = True
504     for i in range(agent.test_state_no):
505         if done:
506             done = False
507             state = env.reset()
508             state = np.reshape(state, [1, state_size])
509             test_states[i] = state
510         else:
511             action = random.randrange(action_size)
512             next_state, reward, done, info = env.step(action)
513             next_state = np.reshape(next_state, [1, state_size])
514             test_states[i] = state
515             state = next_state
516
517     scores, episodes, losses = [], [], []
518     # Create dynamically growing score and episode counters
519     for e in range(EPISODES):
520         done = False
521         score = 0
522         state = env.reset() # Initialize/reset the environment
523         state = np.reshape(state, [1,
524                                state_size])
525     # Reshape state so that to a 1 by state_size two-dimensional array ie. [x_1,x_2] to
526     # Compute Q values for plotting
527     tmp = agent.model.predict(test_states)
528     max_q[e][:] = np.max(tmp, axis=1)
529     max_q_mean[e] = np.mean(max_q[e][:])
530
531     while not done:
532         if agent.render:
533             env.render() # Show cartpole animation
534
535         # Get action for the current state and go one step in environment
536         action = agent.get_action(state)
537         next_state, reward, done, info = env.step(action)
538         next_state = np.reshape(next_state, [1, state_size])
539     # Reshape next_state similarly to state
540
541     # Save sample <s, a, r, s'> to the replay memory
542     agent.append_sample(state, action, reward, next_state, done)

```



```

543     # Training step
544     agent.train_model()
545
546     score += reward # Store episodic reward
547     state = next_state # Propagate state
548
549     if done:
550         # At the end of very episode, update the target network
551         if e % agent.target_update_frequency == 0:
552             agent.update_target_model()
553         # Plot the play time for every episode
554
555         loss_list = agent.losses
556         losses.append(np.mean(loss_list))
557
558         scores.append(score)
559         episodes.append(e)
560
561         print("episode:", e, "_score:", score, "_q_value:", max_q_mean[
562             len(agent.memory))
563
564         # if the mean of scores of last 100 episodes is bigger than 195
565         # stop training
566         if agent.check_solve:
567             if np.mean(scores[-min(100, len(scores)):]) >= 195:
568                 print("solved_after", e - 100, "episodes")
569                 agent.plot_data(episodes, scores, max_q_mean[:e + 1])
570                 sys.exit()
571     # agent.plot_data(episodes, scores, max_q_mean)
572     all_scores[n][:] = np.asarray(scores)
573     all_max_q_mean[n][:] = max_q_mean.ravel()
574     #print(losses)
575     all_losses[n][:] = np.asarray(losses)
576
577 multi_plot(episodes, all_scores, all_max_q_mean, all_losses, agent)

```