

Chess Assignment

Introduction to Reinforcement Learning (Spring 2022). Source code available at: https://github.com/TwoDigitsOneNumber/IntroRL_ChessAssignment

van den Bergh Laurin

Institute of Informatics

University of Zurich

Zurich, Switzerland

laurin.vandenbergh@uzh.ch — Matriculation number: 16-744-401

Abstract—We explored the use of three deep reinforcement learning methods for training agents to play a simplified version of chess. All algorithms, SARSA, Q-Learning and DQN were able to learn successful strategies. DQN was able to overcome the shortcomings of Q-Learning and provides an off-policy method with performance comparable to the on-policy method SARSA.

Index Terms—deep reinforcement learning, chess, temporal-difference methods

I. INTRODUCTION

In this assignment, we explore three different deep reinforcement learning algorithms to learn how to play a simplified version of chess, which can be thought of as a special instance of an endgame. These three algorithms are SARSA, Q-Learning and DQN¹. We will first provide a general look over our methodology (see Section II) and later discuss the result obtained in our experiments (see Section III).

The focus of this report lies on comparing these three algorithms in theory and in practise on the chess endgame environment. We further explored the impact of the hyper-parameters β and γ , which represent the speed of the decaying trend for the learning rate and the discount factor respectively.

Throughout the report we indicate in footnotes which task a particular section is referring to in terms of answering the task. We do this as the solutions to certain tasks are spread throughout multiple sections, e.g. task 3 is answered in Section II and Section III. Even though this assignment was not solved in a group, we decided to also answer some of the “group only” and we stick to the numbering of the assignment in order to avoid confusion.

II. METHODS

A. Environment

This version of chess takes place on a 4 by 4 board and can be thought of as a specific version of an endgame where the agent has a king and a queen, and the opponent has only a king. Since this game can only end in a win for the agent or in a draw, it is the agent’s goal to learn how to win the game and avoid draws. For all experiments considered, the agent will be given a reward of 1 for winning, 0 for drawing, and 0 for all intermediate steps.

¹SARSA serves as answer to task 3 and DQN serves as answer to task 5.

This chess setting, and chess in general, fulfills the Markov property and therefore justifies the use of the temporal difference methods used in this assignment.

B. SARSA and Q-Learning²

1) *Temporal-Difference Algorithms*: SARSA and Q-Learning are two very related model-free types of temporal-difference (TD) algorithms for learning expected rewards, also known as Q-values, when rewards are not immediate and possibly sparse. The learning takes place via interaction with an environment through trial and error. These Q-values are in general represented by an action-value function Q and, for finitely many state-action pairs (s, a) , can be considered as a Q-table where each state-action pair, (s, a) , maps to a single Q-value, thus providing an estimate of the quality of any given state-action pair (s, a) . In this assignment however we use neural networks to approximate the action-value function, which outputs the Q-values for all possible actions for any given state. This helps to avoid computing large Q-tables. All algorithms explored in this assignment, including DQN, require the environment to fulfill the Markov property.

2) *On-policy vs. Off-policy*: SARSA and Q-Learning address the temporal-credit-assignment problem [1], that is, trying to attribute future rewards to previous actions. These future rewards get discounted with the hyper-parameter γ (see Section III-C). Both algorithms repeatedly choose and take actions in the environment according to some policy π , e.g. an ϵ -greedy policy.

However, this is where they differ. SARSA is an on-policy algorithm, which means that consecutive actions are chosen according to the same policy π , even during the update step of the Q-values, which leads to the update rule:

$$Q_{\pi}(s, a) \leftarrow Q_{\pi}(s, a) + \eta(r + \gamma Q_{\pi}(s_{t+1}, a') - Q_{\pi}(s, a))$$

for some future action a' chosen according to policy π .

Q-learning, on the other hand, is an off-policy algorithm, which means that it takes its actions a according to its policy π , but during the update steps it assumes a greedy policy, i.e. optimal play, for future actions a' . Q-Learning has the update rule:

$$Q_{\pi}(s, a) \leftarrow Q_{\pi}(s, a) + \eta(r + \gamma \max_{a'} Q_{\pi}(s_{t+1}, a') - Q_{\pi}(s, a)).$$

²Answer to task 1.

3) *Advantages and Disadvantages*: This leads to one of Q-Learning’s major advantages: Because of Bellman’s optimality equation, Q-Learning is guaranteed to learn the values for the optimal policy, i.e. $Q_*(s, a) = \max_{\pi} Q_{\pi}(s, a)$, regardless of the policy used to train it, and in a greedy setting will take the optimal actions, at least if it was trained sufficiently. However, this can in certain cases mean that the online performance of Q-Learning will be worse than the one from SARSA, as Sutton et al. [2] demonstrate with their “gridworld” example “Cliff Walking”. Our chess game is a similar situation, because a win and a draw can be very close, thus during exploration Q-Learning can accidentally create a draw because it is going for the optimum when exploiting. Q-Learning is however relatively unstable and the parameters can even diverge when it is combined with non-linear function approximators [3], making the guarantee to learn the optimal policy irrelevant.

SARSA will learn to take a safer path, because it keeps its policy in mind when updating the Q-values, i.e. it keeps in mind that it will explore in future actions. This has the advantage that SARSA in general tends to explore more than Q-Learning.

C. Experience Replay³

Experience replay is a technique proposed by Lin [4] to speed up the training process for reinforcement learning algorithms by reusing past experiences for future training. This is analogous to the human ability to remember past experiences and learn from them even after the fact. The past experiences are stored in a replay memory of fixed size at each time step t as a tuple $e_t = (s_t, a_t, r_t, s_{t+1})$. This essentially allows us to transform the learning process from online learning to mini-batch learning, where a batch of experiences e_j is randomly sampled for each update step. Experience replay can only be used in combination with off-policy algorithms, because otherwise the current parameters determine the next sample and create unwanted feedback loops [3], [5].

Experience replay provides many benefits over online Q-Learning, especially when neural networks are used to approximate the action-value function. First, it enables the agent to learn from past experiences more than once, leading to increased data efficiency and faster convergence [4], [5]. Second, since for each update step past experiences are sampled randomly, the correlations between the individual actions are reduced, which then reduces the variance of the updates [5]. This leads to the experience samples e_j being closer to i.i.d. and thus guaranteeing better convergence when using optimization algorithms such as stochastic gradient descent as most convergence proofs assume i.i.d. data.

D. Deep Q-Networks (DQN)⁴

A first version of the DQN algorithm was proposed by Mnih et al. [5] and combined experience replay with Q-learning, where a neural network was used as a non-linear function approximator for the action-value function. Mnih et al. [3] later

improved upon the method and presented the DQN algorithm, as it is known today, where they address the problem of the Q-values $Q_{\pi}(s, a)$ being correlated to the target values $y = r + \gamma \max_{a'} Q_{\pi}(s', a')$ because they are generated using the same neural network. In the DQN algorithm they separated the Q-network from the target network and only update the target network every C steps, which helps to break this correlation and combat diverging network parameters.

Since DQN uses experience replay, we essentially transform the reinforcement learning task to a supervised learning task. Therefore a suitable loss function for the neural network is needed. Mnih et al. [3] used a squared loss of the temporal-difference error, also known as delta: $\delta = y - Q_{\pi}(s, a)$.

E. Experiments

In order to address all tasks, we divided the tasks into several independent experiments. First, we conducted seeded runs⁵ for all three algorithms using seed 21 for reproducibility, which was chosen a-priori. These seeded runs serve as examples to compare the algorithm’s online performance qualitatively. The seeds are used such that the weights of all neural networks are instantiated identically for all algorithms and they subsequently serve as seeds for any random number used during training. This makes sure that all agents start with the same initial conditions and that the results are reproducible (see Section V-A). All algorithms were run for 100000 episodes using identical model architecture and hyper-parameters (see Section II-F).

Since the seeded runs are heavily influenced by the choice of the seed, we could end up with anything between a very lucky and well performing seed, or with a very unlucky one. Also the interpretation of the seeded runs is more difficult as we just have one run for each algorithm. Therefore, we decided to perform a simulation study and complete 30 non-seeded runs for each algorithm in order to get a better idea of how the algorithms perform on average. For computational reasons we limited these runs to 40000 episodes as we realized with test runs that by then most of the training progress has already taken place.

To analyze the impact of the hyper-parameters β and γ ⁶ we trained 49 agents with different combinations for β and γ but keeping all other hyper-parameters and model architecture identical. We chose SARSA for this experiment as we found it to have very low variance between its unseeded runs, which makes it an ideal candidate for comparing individual runs. These runs are seeded identically to the seeded runs mentioned above.

F. Implementation and Hyper-parameters

We implemented all algorithms from scratch according to Sutton et al. [2] (SARSA and Q-Learning⁷) and Mnih et al. [3] (DQN⁸). For the implementation see file `neural_net.py`

³Answer to “group only” task 2.

⁴Answer to task 5: Describing the used method.

⁵Answers to task 3 and 5.

⁶Answer to task 4.

⁷SARSA as answer to task 3 and Q-Learning as additional algorithm.

⁸Answer to task 5.

on GitHub or Listing 1. All algorithms use a neural network with 58 input neurons, 200 hidden neurons and 32 output neurons, not including the biases for the input and hidden layer. The neural network automatically adds a constant input for the bias and the hidden layer. The implementation treats the biases like any other weights and thus they are part of any matrix multiplication. We used a ReLU activation function for the hidden layer and no activation on the output layer. The weights were initialized using Glorot initialization [6], such that the weights are sampled from a normal distribution with mean 0 and variance $\frac{2}{n_{in} + n_{out}}$, where n_{in} and n_{out} denote to the number of input and output neurons of the respective layer. This helped preventing exploding gradients for the most part. [check again with results](#)

For all experiments we used the default hyperparameters provided in the `Assignment.ipynb` file unless otherwise noted (see Table I). For DQN we updated the weights of the target network after every $C = 10$ steps, as most games take fewer steps than that. We used a replay memory of size 100000 and a batch size of 32.

Parameter	Value
Nr. input neurons	58+1
Nr. hidden neurons	200+1
Nr. output neurons	32
Initial exploration probability ϵ_0	0.2
Learning rate η	0.035
Decay rate of ϵ , β	0.00005
Discount factor γ	0.85

TABLE I: Common hyper-parameters shared by all algorithms.

III. RESULTS

A. Seeded Runs⁹

The rewards and number of moves for the seeded runs are depicted in Figures 1 and 2 respectively. Since the curves are very noisy, we smoothed them using an exponential moving average (EMA) with a weight on the most recent observation of $\alpha = 0.001$.

As expected, the online performance of Q-Learning in terms of the rewards is generally lower than the rewards for SARSA but they converge slowly as ϵ decreases (Figure 1). Also in Figures 1 and 2 we can see that Q-Learning experiences instable learning behavior as both plots are a lot more noisy and at about 20000 and 90000 episodes the rewards decrease for some period. SARSA and DQN don't show this behavior.

Even though the number of steps is not punished, all agents still learn to reduce the number of steps over time, as they do not give rewards and their goal is to take actions that do. SARSA seems to do the best job at this, which perhaps is caused by its tendency to explore more and find better strategies. Q-Learning however seems to struggle to reduce the number of steps it takes.

⁹Answer to task 3 (SARSA and Q-Learning as additional algorithm) and 5 (DQN).

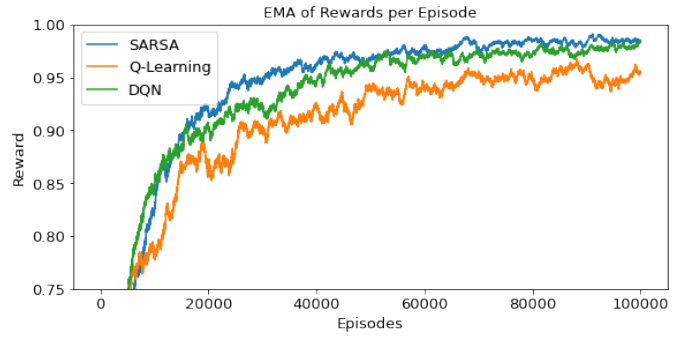


Fig. 1: Exponential moving average of the rewards achieved during training for 100000 episodes with identical hyper-parameters, weight initialization and model architecture.

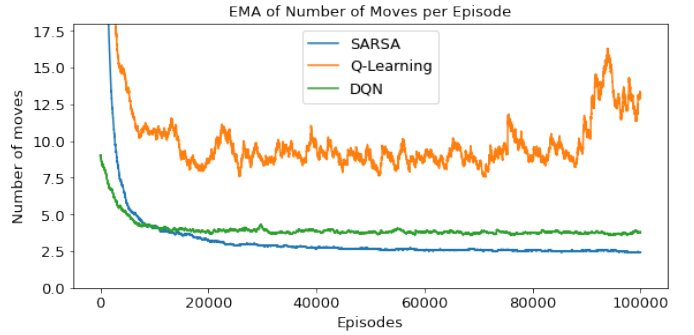


Fig. 2: Exponential moving average of the number of moves per episode achieved during training for 100000 episodes with identical hyper-parameters, weight initialization and model architecture.

As suggested by Mnih et al. [3], [5], DQN¹⁰ was able to overcome the downsides of Q-Learning which lead to an online performance which is comparable to that of SARSA in terms of reward and number of moves it achieved, and also in terms of the stability during training. It did however not learn to reduce the amounts of steps as much as SARSA.

B. Simulation Study (Non-seeded Runs)

We can confirm that the qualitative results from the seeded runs reasonably well represent the average case. The only notable exception being Q-Learning, for which most runs performed equally well to the seeded run, but some runs experienced huge increases in the number of steps, which influenced the average run dramatically, leading to an average of about 23 moves per episode after 40000 episodes.

We observed that DQN and SARSA show very comparable learning curves with DQN showing slightly faster convergence in the first 5000 episodes. SARSA showed the lowest variance in all runs and seems to be a very stable algorithm. Q-Learning on the other hand showed clear signs of divergence as for some runs the rewards consistently dropped while the number of moves consistently increased. This shows that the measures

¹⁰Answer to task 5 for comparing DQN to SARSA and Q-Learning.

taken by Mnih et al. [3] to combat the disadvantages of Q-Learning worked and increased the stability as well as the convergence speed. We were able to verify that the gradients of the Q-Learning agents were a lot less stable than the gradients of the other agents. However, using the Glorot initialization [6] helped prevent exploding gradients from occurring.

We also found out that, unsurprisingly, the effective training time is mainly dependent on the number of steps an algorithm takes per episode. This leads to Q-Learning having by far the longest training time, especially when the parameters diverge and the number of steps increase. DQN and SARSA have relatively short training times, with SARSA being the fastest.

We can conclude that the seeded runs in our initial experiment truthfully represent the average run and therefore some level of inference is justified.

C. Hyper-parameters¹¹

Figure 3(a) depicts the rewards and number of moves per episode as a function of β and γ . We can see that the reward increases monotonically as γ is increased, suggesting that a value of $\gamma \in [0.80, 1)$ should be chosen for almost all values of β . This intuitively makes sense, as we have very sparse rewards and want the agent to “backpropagate” this reward through its sequence of actions. The left plot of Figure 3 suggests that reducing γ to a value in $[0.5, 0.8]$ can teach the SARSA agent to not reduce the number of steps. Intuitively this makes sense, as the only reward will be “backpropagated” less to earlier states and thus the agent will move faster towards setting the opponent’s king checkmate.

We can not see a clear relationship between β and the rewards, apart from $\beta = 0$ being an inferior choice for all values of γ . In Figure 3(b) we can see that the number of steps taken by the agent decreases drastically when increasing γ from very low levels, but this effect seems larger for larger values of β . We can however see that there is a slight, but possibly insignificant, peak in the rewards around $\beta = 5 \cdot 10^{-3}$. In summary, for reasonably chosen values of γ the choice of β seems to not have much of an influence for training periods of around 40000 episodes.

IV. CONCLUSION

We are aware that the performance of the individual algorithms could be improved by tuning the hyper-parameters, however, this was not explicitly asked for and the focus on this assignment lies on the comparison of these algorithms from a theoretical and practical perspective.

For any deep reinforcement learning method the choice of suitable hyper-parameters for the task is crucial and can have large impacts on the training outcome. In our case, the default parameters provided to us performed very well so no need for much further consideration was necessary.

All three algorithms were able to learn to play the simplified version of chess to a very high degree even without hyper-parameter tuning. SARSA proved to be the most stable

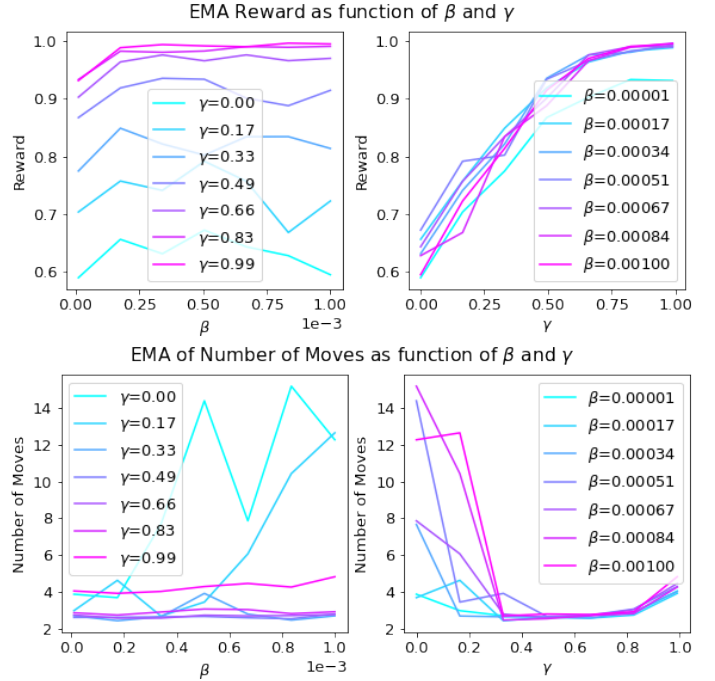


Fig. 3: Rewards and number of moves as functions of the speed of the decaying trend β and the discount factor γ after training a SARSA agent for 40000 episodes.

algorithm, which was confirmed to be the general case with 30 non-seeded runs. Q-Learning suffers from some instabilities when training, but DQN was able to overcome all of the problems of Q-Learning and provides an off-policy method that can learn with high stability, fast convergence and a low training time comparable to SARSA. Since DQN is an off-policy method, it comes with the added advantage that it will learn an optimal policy, similar to Q-Learning.

REFERENCES

- [1] R. S. Sutton, “Temporal credit assignment in reinforcement learning,” 1984.
- [2] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. The MIT Press, 2018. [Online]. Available: <http://incompleteideas.net/book/the-book-2nd.html>
- [3] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. A. Riedmiller, A. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, “Human-level control through deep reinforcement learning,” *Nat.*, vol. 518, no. 7540, pp. 529–533, 2015. [Online]. Available: <https://doi.org/10.1038/nature14236>
- [4] L.-J. Lin, “Self-improving reactive agents based on reinforcement learning, planning and teaching,” *Mach. Learn.*, vol. 8, no. 3–4, p. 293–321, may 1992. [Online]. Available: <https://doi.org/10.1007/BF00992699>
- [5] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. A. Riedmiller, “Playing atari with deep reinforcement learning,” *CoRR*, vol. abs/1312.5602, 2013. [Online]. Available: <http://arxiv.org/abs/1312.5602>
- [6] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *In Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS’10). Society for Artificial Intelligence and Statistics*, 2010.

¹¹ Answer to task 4.

A. Reproducibility

In order to reproduce the results presented in this report we provide the code on GitHub in the repository https://github.com/TwoDigitsOneNumber/IntroRL_ChessAssignment. Along this we provide a conda environment file `environment.yaml` which can be used to recreate the exact environment we used. We recommend to run the file `Assignment_Train_Algorithms.ipynb` before the files `Assignment_Compare_Algorithms.ipynb` and `Assignment_Hyperparameter_Influence.ipynb` as the latter use files generated by the former. However, even on fast hardware running the former file takes between 5-6 hours, so we provide all necessary intermediate outputs in the repository as well.

B. Code Excerpts

```

1 # import libraries
2 from types import MethodDescriptorType
3 import numpy as np
4 from tqdm.notebook import tqdm
5 import os
6 import json
7 import time
8 import random
9 from collections import namedtuple, deque
10
11 # import from files
12 from Chess_env import *
13
14
15
16 # ===== Epsilon-greedy Policy =====
17
18 def EpsilonGreedy_Policy(Qvalues, allowed_a, epsilon):
19     """
20     returns: tuple
21         an action in form of a one-hot encoded
22         vector with the same shape as
23         Qvalues.
24         an action as decimal integer (0-based)
25
26     Assumes only a single state, i.e. online
27     learning and NOT (mini-)batch learning.
28     """
29     # get the Qvalues and the indices (relative to
30     # all Qvalues) for the allowed actions
31     allowed_a_ind = np.where(allowed_a==1)[0]
32     Qvalues_allowed = Qvalues[allowed_a_ind]
33
34     # ----- epsilon greedy -----
35
36     # draw a random number and compare it to epsilon
37     rand_value = np.random.uniform(0, 1, 1)
38
39     if rand_value < epsilon: # if the random number
40         is smaller than epsilon, draw a random
41         action
42         action_taken_ind_of_allwed_only = np.random.
43             randint(0, len(allowed_a_ind))
44     else: # greedy action
45         action_taken_ind_of_allwed_only = np.argmax(
46             Qvalues_allowed)
47
48     # get index of the action that was chosen (
49         relative to all actions, not only allowed)

```

```

ind_of_action_taken = allowed_a_ind[
    action_taken_ind_of_allwed_only]

```

```

# ----- create usable output -----

# get the shapeensions of the Qvalues
N_a, N_samples = np.shape(Qvalues) # N_samples
must be 1

# initialize all actions of binary mask to 0
A_binary_mask = np.zeros((N_a, N_samples))
# set the action that was chosen to 1
A_binary_mask[ind_of_action_taken, :] = 1

return A_binary_mask, ind_of_action_taken

# ===== activation functions and it's derivatives
=====

# relu and its derivative
def relu(x):
    return np.maximum(0, x)

def heaviside(x):
    return np.heaviside(x, 0)

# sigmoid and its derivative
def sigmoid(x):
    return 1 / (1 + np.exp(-x))

def gradient_sigmoid(x):
    return sigmoid(x) * (1 - sigmoid(x))

# tanh and its derivative
def tanh(x):
    return np.tanh(x)

def gradient_tanh(x):
    return 1 - np.tanh(x)**2

# identity and its derivative
def identity(x):
    return x

def const(x):
    return np.ones(x.shape)

def act_f_and_gradient(activation_function="relu"):
    if activation_function == "relu":
        return relu, heaviside
    elif activation_function == "sigmoid":
        return sigmoid, gradient_sigmoid
    elif activation_function == "tanh":
        return tanh, gradient_tanh
    else: # identity and constant 1
        return identity, const

# ===== Replay Memory for Experience Replay (with
DQN) =====

Transition = namedtuple('Transition', ("state", "
    action", "reward", "next_state", "done"))

class ReplayMemory(object):
    def __init__(self, capacity):
        self.memory = deque(maxlen=capacity)

    def push(self, *args):

```

```

111         self.memory.append(Transition(*args))
112
113     def sample(self, batch_size):
114         # if less data than batch size, return all
115         data
116         if len(self) < batch_size:
117             batch_size = len(self)
118         return random.sample(self.memory, batch_size)
119
120     def __len__(self):
121         return len(self.memory)
122
123
124
125 # ===== Neural Network =====
126
127 class NeuralNetwork(object):
128
129     def __init__(self, N_in, N_h, N_a,
130                 activation_function_1="relu",
131                 activation_function_2=None, method="
132                 qlearning", seed=None, capacity=100_000, C
133                 =100):
134         """
135         activation functions: "relu", "sigmoid", "
136         tanh", None
137         methods: "qlearning", "sarsa", "dqn"
138         """
139         self.D = N_in # input dimension (without
140         bias)
141         self.K = N_h # nr hidden neurons (without
142         bias)
143         self.O = N_a # nr output neurons (letter
144         , not digit 0)
145
146         # store method and seed
147         self.method = method
148         self.seed = seed
149
150         if self.method == "dqn":
151             self.capacity = capacity
152             self.replay_memory = ReplayMemory(
153                 capacity)
154             self.C = C
155
156         # set activation function and gradient
157         function
158         self.act_f_1_name = activation_function_1
159         self.act_f_2_name = activation_function_2
160         self.act_f_1, self.grad_act_f_1 =
161         act_f_and_gradient(activation_function_1)
162         self.act_f_2, self.grad_act_f_2 =
163         act_f_and_gradient(activation_function_2)
164
165         # initialize the weights and biases and seed
166         global seed
167         np.random.seed(self.seed)
168
169         # self.W1 = np.random.randn(self.K+1, self.D
170         +1)/np.sqrt(self.D+1) # standard norm
171         distribution, shape: (K+1, D+1)
172         # gloriot/xavier normal initialization
173         self.W1 = np.random.randn(self.K+1, self.D
174         +1)*np.sqrt(2/ (self.D+1 + self.K+1)) #
175         standard normal distribution, shape: (K
176         +1, D+1)
177         self.W1 = np.random.standard_normal((self.K
178         +1, self.D+1))*np.sqrt(2/ (self.D+1 +
179         self.K+1)) # standard normal
180
181         distribution, shape: (K+1, D+1)
182         # self.W2 = np.random.randn(self.O, self.K
183         +1)/np.sqrt(self.K+1) # standard normal
184         distribution, shape: (O, K+1)
185         # gloriot/xavier normal initialization
186         self.W2 = np.random.standard_normal((self.O,
187         self.K+1))*np.sqrt(2/ (self.K+1 + self.
188         O)) # standard normal distribution,
189         shape: (O, K+1)
190         # self.W2 = np.random.randn(self.O, self.K
191         +1) # standard normal distribution,
192         shape: (O, K+1)
193
194         if self.method == "dqn":
195             self.W1_target = np.copy(self.W1)
196             self.W2_target = np.copy(self.W2)
197
198     def forward(self, x, target=False):
199         """
200         x has shape: (D+1, 1) (constant bias 1 must
201         be added beforehand)
202         target: if True, use the weights of the
203         target network
204
205         returns:
206         last logits (i.e. Qvalues) of shape (O,
207         1)
208         """
209         if target == True:
210             W1 = np.copy(self.W1_target)
211             W2 = np.copy(self.W2_target)
212         else:
213             W1 = np.copy(self.W1)
214             W2 = np.copy(self.W2)
215
216         # forward pass/propagation
217         a1 = W1 @ x
218         h1 = self.act_f_1(a1)
219         h1[0,:] = 1 # set first row (bias to second
220         layer) to 1 (this ignores the weights
221         for the k+1th hidden neuron, because
222         this should not exist; this allows to
223         only use matrix multiplication and
224         simplify the gradients as we only need 2
225         instead of 4)
226         a2 = W2 @ h1
227         h2 = self.act_f_2(a2)
228         return a1, h1, a2, h2
229
230     def backward(self, R, x, Qvalues, Q_prime, a1,
231                 h1, a2, gamma, future_reward,
232                 action_binary_mask):
233         """
234         backward for methods "qlearning" and "sarsa"
235
236         x has shape (D+1, 1) (constant bias 1 must
237         be added beforehand)
238         set future_reward=True for future reward
239         with gamma>0, False for immediate reward
240
241         Q_prime must be chosen according to the
242         method on x_prime (on- or off-policy)
243         """
244
245         # backward pass/backpropagation
246         # compute the gradient of the square loss
247         with respect to the parameters

```



```

209                                     270
210 # ===== compute TD error (aka delta) ===== 271
211                                     272
212 # make reward of shape (O, 1) 273
213 R_rep = np.tile(R, (self.O, 1)) 274
214 if future_reward: # future reward 275
215     delta = R_rep + gamma*Q_prime - Qvalues
216         # -> shape (O, 1) 276
217 else: # immediate reward 277
218     delta = R_rep - Qvalues # -> shape (O
219         1) 278
220                                     279
221                                     280
222 # update only action that was taken, i.e. 281
223     all rows apart from the one 282
224     corresponding to the action taken ( 283
225     action index) are 0 284
226 delta = delta*action_binary_mask 285
227                                     286
228 self.compute_gradients(delta, a1, h1, a2, 287
229 self.update_parameters(self.eta) 288
230                                     289
231
232 def backward_dqn(self, batch, gamma): 290
233     """ 291
234     backward for method "dqn" 292
235     """ 293
236                                     294
237 # ===== compute targets y and feature matrix 295
238     X ===== 296
239                                     297
240 # turn batch into individual tuples, numpy 298
241     arrays, or lists
242 states = batch.state
243 rewards = np.array(list(batch.reward)) 299
244 actions = np.array(list(batch.action)) 300
245 next_states = list(batch.next_state)
246 dones = np.array(list(batch.done))
247                                     301
248 # compute targets y and feature matrix X 302
249 y = np.zeros((self.O, len(dones))) 303
250 for j in np.arange(len(dones)): 304
251     if dones[j]: # if done, set y_j = r_j 305
252         y[actions[j], j] = rewards[j] 306
253     else: 307
254         # compute Q_prime 308
255         Q_target = self.forward(next_states 309
256             j], target=True)[-1] 310
257         y[actions[j], j] = rewards[j] + 311
258             gamma*np.max(Q_target) 312
259                                     313
260                                     314
261 # convert states to feature matrix X 315
262 X = np.hstack((states)) 316
263                                     317
264                                     318
265 # ===== compute TD error (aka delta) ===== 319
266                                     320
267 a1, h1, a2, Qvalues = self.forward(X) 320
268 delta = y - Qvalues # -> shape (O, 321
269     batch_size)
270
271 self.compute_gradients(delta, a1, h1, a2, 322
272 self.update_parameters(self.eta) 323
273                                     324
274
275 def compute_gradients(self, delta, a1, h1, a2, 324
276 ): 325
277 # ===== compute gradient of the loss with 326
278     respect to the weights =====
279                                     327
280 # common part of the gradient TODO: check 328
281     dimensions 329
282 self.dL_da2 = delta * self.grad_act_f_2(a2 330

```

```

# gradient of loss wrt W2
self.dL_dW2 = self.dL_da2 @ h1.T

# gradient of loss wrt W1
self.dL_dW1 = ( (self.W2.T @ self.dL_da2) *
    self.grad_act_f_1(a1) ) @ x.T

```

```

def update_parameters(self, eta):

    # gradient clipping

    # dL_dW1_norm = np.linalg.norm(self.dL_dW1)
    # if dL_dW1_norm >= self.gradient_clip:
    #     self.dL_dW1 = self.gradient_clip *
    #         self.dL_dW1 / dL_dW1_norm

    # dL_dW2_norm = np.linalg.norm(self.dL_dW2)
    # if dL_dW2_norm >= self.gradient_clip:
    #     self.dL_dW2 = self.gradient_clip *
    #         self.dL_dW2 / dL_dW2_norm

    # update W1 and W2
    self.W2 = self.W2 + eta * self.dL_dW2
    self.W1 = self.W1 + eta * self.dL_dW1

def train(self, env, N_episodes, eta, epsilon_0,
    beta, gamma, alpha=0.001, gradient_clip=1,
    batch_size=32, run_number=None):
    """
    alpha is used as weight for the exponential
    moving average displayed during training

    batch_size is only used for the DQN method.
    """

    # add training hyper parameters
    self.N_episodes = N_episodes
    self.eta = eta
    self.epsilon_0 = epsilon_0
    self.beta = beta
    self.gamma = gamma
    self.alpha = alpha
    self.gradient_clip = gradient_clip
    self.batch_size = batch_size

    training_start = time.time()

    try:

        # initialize histories for important
        # metrics
        self.R_history = np.full([self.
            N_episodes, 1], np.nan)
        self.N_moves_history = np.full([self.
            N_episodes, 1], np.nan)
        self.dL_dW1_norm_history = np.full([self.
            N_episodes, 1], np.nan)
        self.dL_dW2_norm_history = np.full([self.
            N_episodes, 1], np.nan)

        # progress bar
        episodes = tqdm(np.arange(self.
            N_episodes), unit="episodes")
        ema_previous = 0

        n_steps = 0

```

```

331         for n in episodes:
332             epsilon_f = self.epsilon_0 / (1 +
333                 beta * n)  ## DECAYING EPSILON
334             Done = 0
335             i = 1
336             ## COUNTER FOR NUMBER OF ACTIONS
337             S, X, allowed_a = env.
338                 Initialise_game()  ##
339                 INITIALISE GAME
340             X = np.expand_dims(X, axis=1)
341             ## MAKE X A
342             TWO DIMENSIONAL ARRAY
343             X = np.copy(np.vstack((np.array
344                 ([[1]]), X)))  # add bias term
345             if self.method == "sarsa":
346                 # compute Q values for the given
347                 state
348                 a1, h1, a2, Qvalues = self.
349                     forward(X)  # -> shape (O,
350                     1)
351                 # choose an action A using
352                 epsilon-greedy policy
353                 A_binary_mask, A_ind =
354                     EpsilonGreedy_Policy(Qvalues,
355                         allowed_a, epsilon_f)  #
356                         -> shape (O, 1)
357             while Done==0:
358                 ##
359                 START THE EPISODE
360                 if (self.method == "qlearning"
361                     or (self.method == "dqn")):
362                     # compute Q values for the
363                     given state
364                     a1, h1, a2, Qvalues = self.
365                         forward(X)  # -> shape
366                         (O, 1)
367                     # choose an action A using
368                     epsilon-greedy policy
369                     A_binary_mask, A_ind =
370                         EpsilonGreedy_Policy(
371                             Qvalues, allowed_a,
372                             epsilon_f)  # -> shape
373                             (O, 1)
374                     # take action and observe reward
375                     R and state S_prime
376                     S_prime, X_prime,
377                         allowed_a_prime, R, Done =
378                         env.OneStep(A_ind)
379                     X_prime = np.expand_dims(X_prime,
380                         axis=1)
381                     X_prime = np.copy(np.vstack((np.
382                         array([[1]]), X_prime)))
383                     add bias term
384                     n_steps += 1
385                     if self.method == "dqn":
386                         # store the transition in
387                         memory
388                         self.replay_memory.push(X,
389                             A_ind, R, X_prime, Done)
390                     # sample a batch of
391                     transitions
392                     transactions = self.
393                         replay_memory.sample(
394                             self.batch_size)
395                     # turn list of transactions
396                     into transaction of
397                     lists
398                     batch = Transition(*zip(*
399                         transactions))
400                     # backward step and
401                     parameter update
402                     self.backward_dqn(batch,
403                         self.gamma)
404                     # update Q values indirectly by
405                     updating the weights and
406                     biases directly
407             if Done==1:  # THE EPISODE HAS
408                 ENDED, UPDATE...BE CAREFUL,
409                 THIS IS THE LAST STEP OF THE
410                 EPISODE
411                 if (self.method == "
412                     qlearning") or (self.
413                         method == "sarsa"):
414                     # compute gradients and
415                     update weights
416                     self.backward(R, X,
417                         Qvalues, None, a1,
418                         h1, a2, None,
419                         future_reward=False,
420                         action_binary_mask=
421                             A_binary_mask)
422                     # store history
423                     # todo: record max possible
424                     reward per episode
425                     self.R_history[n] = np.copy(
426                         R)  # reward per episode
427                     self.N_moves_history[n] = np.
428                         copy(i)  # nr moves per
429                         episode
430                     # store norm of gradients
431                     self.dL_dW1_norm_history[n]
432                         = np.linalg.norm(self.
433                             dL_dW1)
434                     self.dL_dW2_norm_history[n]
435                         = np.linalg.norm(self.
436                             dL_dW2)
437                     # compute exponential moving
438                     average (EMA) to
439                     display during training
440                     ema = alpha*R + (1-alpha)*
441                         ema_previous
442                     if n == 0:  # first episode
443                         ema = R
444                     ema_previous = ema
445                     if run_number is not None:
446                         episodes.set_description
447                             (f"Run = {run_number
448                                 }; EMA Reward = {ema
449                                     :.2f}")
450                     else:
451                         episodes.set_description
452                             (f"EMA Reward = {ema
453                                 :.2f}")

```



```

405         break
406     else: # IF THE EPISODE IS NOT OVER...
407         if self.method == "qlearning":
408             # chose next action off policy
409             Q_prime = np.max(self.forward(X_prime)[-1])
410         elif self.method == "sarsa":
411             # chose next action on policy
412             a1_prime, h1_prime, a2_prime, Qvalues_prime = self.forward(X_prime)
413             # chose next action and save it
414             A_binary_mask_prime, A_ind_prime = EpsilonGreedy_Policy(Qvalues_prime, allowed_a_prime, epsilon_f)
415             # get Qvalue of next action
416             Q_prime = Qvalues_prime[A_ind_prime]
417         if (self.method == "qlearning") or (self.method == "sarsa"):
418             # backpropagation and weight update
419             self.backward(R, X, Qvalues, Q_prime, a1, h1, a2, self.gamma, future_reward=True, action_binary_mask=A_binary_mask)
420         # NEXT STATE AND CO. BECOME ACTUAL STATE...
421         if self.method == "sarsa":
422             A_binary_mask = np.copy(A_binary_mask_prime)
423             A_ind = np.copy(A_ind_prime)
424             a1 = np.copy(a1_prime)
425             h1 = np.copy(h1_prime)
426             a2 = np.copy(a2_prime)
427             Qvalues = np.copy(Qvalues_prime)
428             S = np.copy(S_prime)
429             X = np.copy(X_prime)
430             allowed_a = np.copy(allowed_a_prime)
431         i += 1 # UPDATE COUNTER FOR NUMBER OF ACTIONS
432         if (self.method == "dqn") and n_steps % self.C == 0:
433             # update target network every C steps
434             self.W1_target = np.copy(self.W1)
435             self.W2_target = np.copy(self.W2)
436             training_end = time.time()
437             self.training_time_in_seconds = training_end - training_start
438             return None
439         except KeyboardInterrupt as e:
440             # return nothing
441             training_end = time.time()
442             self.training_time_in_seconds = training_end - training_start
443             return None
444
445     def save(self, name_extension=None):
446         # create directory for the model
447         name = f"{self.method}_{self.act_f_1_name}_{self.act_f_2_name}"
448         if name_extension is not None:
449             name += f"_{name_extension}"
450
451         path = f"models/{name}"
452         if not os.path.isdir(path): os.mkdir(path)
453         print(f"saving to: {path}")
454
455         # save weights
456         np.save(f"{path}/W1.npy", self.W1)
457         np.save(f"{path}/W2.npy", self.W2)
458
459         # save training history
460         np.save(f"{path}/training_history_R.npy", self.R_history)
461         np.save(f"{path}/training_history_N_moves.npy", self.N_moves_history)
462         np.save(f"{path}/training_history_dL_dW1_norm.npy", self.dL_dW1_norm_history)
463         np.save(f"{path}/training_history_dL_dW2_norm.npy", self.dL_dW2_norm_history)
464
465         # save training parameters and other general info
466         params = {
467             "method": self.method,
468             "N_episodes": self.N_episodes,
469             "eta": self.eta,
470             "epsilon_0": self.epsilon_0,
471             "beta": self.beta,
472             "gamma": self.gamma,
473             "alpha": self.alpha,
474             # "gradient_clip": self.gradient_clip,
475             "seed": self.seed,
476             "D": self.D,
477             "K": self.K,
478             "O": self.O,
479             "training_time_in_seconds": self.training_time_in_seconds
480         }
481         if self.method == "dqn":
482             params["capacity"] = self.capacity
483             params["batch_size"] = self.batch_size
484             params["C"] = self.C

```

```

505         with open(f"{path}/training_parameters.json"
506                 , "w") as f:
507             json.dump(params , f)
508
509 def load_from(method , act_f_1 , act_f_2 ,
510             name_extension=None):
511
512     # read values and store in neural network
513     instance
514     name = f"{method}_{act_f_1}_{act_f_2}"
515     if name_extension is not None:
516         name += f"_{name_extension}"
517
518     path = f"models/{name}"
519     # print(f"loading from: {path}")
520
521     # initialize neural network
522     nn = NeuralNetwork(0,0,0, activation_function_1=
523         act_f_1 , activation_function_2=act_f_2 ,
524         method=method)
525
526     # network weights
527     nn.W1 = np.load(f"{path}/W1.npy")
528     nn.W2 = np.load(f"{path}/W2.npy")
529
530     # network training history
531     nn.R_history = np.load(f"{path}/
532         training_history_R.npy")
533     nn.N_moves_history = np.load(f"{path}/
534         training_history_N_moves.npy")
535     nn.dL_dW1_norm_history = np.load(f"{path}/
536         training_history_dL_dW1_norm.npy")
537     nn.dL_dW2_norm_history = np.load(f"{path}/
538         training_history_dL_dW2_norm.npy")
539
540     # network training parameters
541     with open(f"{path}/training_parameters.json", "r"
542             ) as f:
543         params = json.load(f)
544
545         # set parameters to the network instance
546         nn.method = params["method"]
547         nn.N_episodes = int(params["N_episodes"])
548         nn.eta = float(params["eta"])
549         nn.epsilon_0 = float(params["epsilon_0"])
550         nn.beta = float(params["beta"])
551         nn.gamma = float(params["gamma"])
552         nn.alpha = float(params["alpha"])
553         # nn.gradient_clip = float(params["
554             gradient_clip"])
555
556         try:
557             nn.seed = int(params["seed"])
558         except:
559             nn.seed = params["seed"]
560         nn.D = int(params["D"])
561         nn.K = int(params["K"])
562         nn.O = int(params["O"])
563         nn.training_time_in_seconds = float(params["
564             training_time_in_seconds"])
565
566         if nn.method == "dqn":
567             nn.capacity = int(params["capacity"])
568             nn.batch_size = int(params["batch_size"
569                 ])
570             nn.C = int(params["C"])
571
572     if nn.method == "dqn":
573         nn.W1_target = np.copy(nn.W1)
574         nn.W2_target = np.copy(nn.W2)
575
576     return nn

```

Listing 1: Object oriented implementation of the neural networks, which can be instantiated with specifications for the model architecture and a method: “sarsa”, “qlearning” or “dqn”. The training loop will adapt automatically.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def moving_average(a, n=3) :
5     steps = len(a)-n
6     ma = np.full(steps , np.nan)
7     for i in range(steps):
8         ma[i] = np.mean(a[i:i+n])
9     return ma, np.arange(steps)
10
11
12 def exponential_moving_average(array , alpha=0.001):
13     """
14     Calculate exponential moving average of an array
15     """
16     ema = np.full(len(array), np.nan)
17     ema[0] = array[0]
18     for i in range(1, len(array)):
19         ema[i] = alpha * array[i] + (1 - alpha) *
20             ema[i-1]
21     return ema
22
23
24 def save_avg_statistics(histories , method):
25     # unpack histories
26     R_histories = [history[0] for history in
27         histories]
28     N_moves_histories = [history[1] for history in
29         histories]
30     training_times = [history[2] for history in
31         histories]
32     layer1_gradient_norms_histories = [history[3]
33         for history in histories]
34     layer2_gradient_norms_histories = [history[4]
35         for history in histories]
36
37     # turn into numpy arrays
38     R_histories = np.hstack(R_histories)
39     N_moves_histories = np.hstack(N_moves_histories)
40     training_times = np.hstack(training_times)
41     layer1_gradient_norms_histories = np.hstack(
42         layer1_gradient_norms_histories)
43     layer2_gradient_norms_histories = np.hstack(
44         layer2_gradient_norms_histories)
45
46     # compute mean and standard deviation for each
47     row of the histories
48     R_mean = np.mean(R_histories , axis=1)
49     R_std = np.std(R_histories , axis=1)
50
51     N_moves_mean = np.mean(N_moves_histories , axis
52         =1)
53     N_moves_std = np.std(N_moves_histories , axis=1)
54
55     layer1_gradient_norms_mean = np.mean(
56         layer1_gradient_norms_histories , axis=1)
57     layer1_gradient_norms_std = np.std(
58         layer1_gradient_norms_histories , axis=1)
59
60     layer2_gradient_norms_mean = np.mean(
61         layer2_gradient_norms_histories , axis=1)
62     layer2_gradient_norms_std = np.std(
63         layer2_gradient_norms_histories , axis=1)
64
65     # save to file

```

```

52     np.save(f"statistics/{method}_R_mean.npy",
53             R_mean)
54     np.save(f"statistics/{method}_R_std.npy", R_std)
55     np.save(f"statistics/{method}_N_moves_mean.npy",
56             N_moves_mean)
57     np.save(f"statistics/{method}_N_moves_std.npy",
58             N_moves_std)
59     np.save(f"statistics/{method}_training_times.npy",
60             training_times)
61     np.save(f"statistics/{method}_layer1_gradient_norms_mean.npy",
62             layer1_gradient_norms_mean)
63     np.save(f"statistics/{method}_layer1_gradient_norms_std.npy",
64             layer1_gradient_norms_std)
65     np.save(f"statistics/{method}_layer2_gradient_norms_mean.npy",
66             layer2_gradient_norms_mean)
67     np.save(f"statistics/{method}_layer2_gradient_norms_std.npy",
68             layer2_gradient_norms_std)
69
70 def load_avg_statistics(method):
71     R_mean = np.load(f"statistics/{method}_R_mean.npy")
72     R_std = np.load(f"statistics/{method}_R_std.npy")
73
74     N_moves_mean = np.load(f"statistics/{method}_N_moves_mean.npy")
75     N_moves_std = np.load(f"statistics/{method}_N_moves_std.npy")
76
77     training_times = np.load(f"statistics/{method}_training_times.npy")
78
79     layer1_gradient_norms_mean = np.load(f"statistics/{method}_layer1_gradient_norms_mean.npy")
80     layer1_gradient_norms_std = np.load(f"statistics/{method}_layer1_gradient_norms_std.npy")
81
82     layer2_gradient_norms_mean = np.load(f"statistics/{method}_layer2_gradient_norms_mean.npy")
83     layer2_gradient_norms_std = np.load(f"statistics/{method}_layer2_gradient_norms_std.npy")
84
85     return R_mean, R_std, N_moves_mean, N_moves_std, training_times, layer1_gradient_norms_mean, layer1_gradient_norms_std, layer2_gradient_norms_mean, layer2_gradient_norms_std
86
87 def printable_name(method):
88     if method == "sarsa":
89         return "SARSA"
90     elif method == "qlearning":
91         return "Q-Learning"
92     elif method == "dqn":
93         return "DQN"
94     else:
95         return None

```

Listing 2: Helper functions used throughout the implementation of the neural network and the notebooks, where the experiments were conducted.