

Tackling Chess with Deep Reinforcement Learning

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. Q-Learning is an additional method beyond what was asked.

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 Methods*: 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.

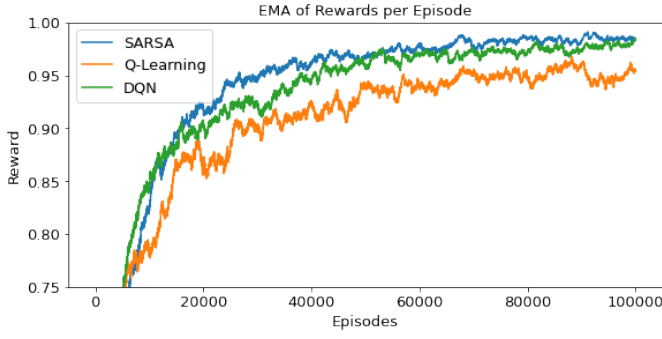


Fig. 1: Exponential moving average of the rewards achieved during training for 100000 episodes with identical hyper-parameters, weight initialization and model architecture.

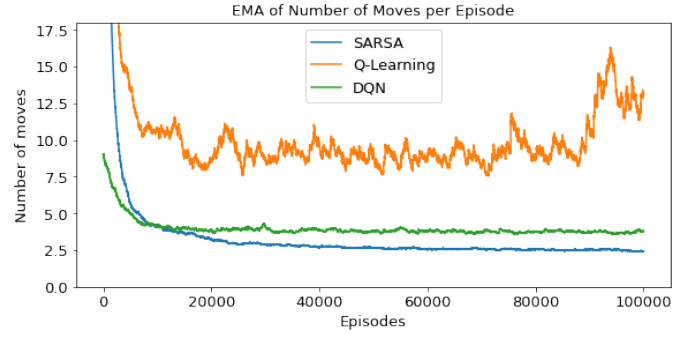


Fig. 2: Exponential moving average of the number of moves per episode achieved during training for 100000 episodes with identical hyper-parameters, weight initialization and model architecture.

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.

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 $\alpha = 0.001$.

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

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.

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.

The hyper-parameter β controls how fast the exploration probability ϵ will decay and therefore controls how the agent will tackle the exploration-exploitation problem. 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.

¹¹ Answer to task 4.

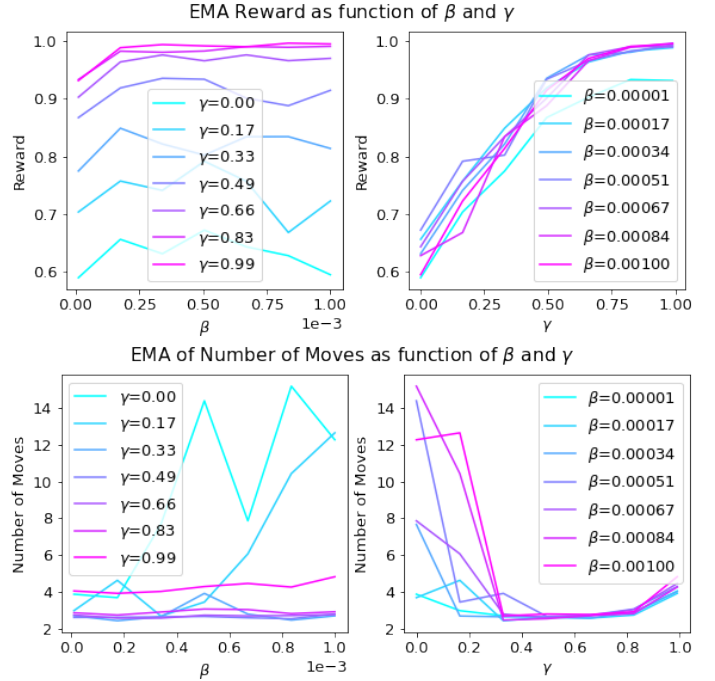


Fig. 3: Rewards and number of moves as functions of the speed of the decaying trend β and the discount factor γ after training a SARSA agent for 40000 episodes.

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 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,” Ph.D. dissertation, 1984.
- [2] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. The MIT Press, 2018.
- [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*

V. APPENDIX

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 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 this 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 # ===== Epsilon-greedy Policy =====
15
16 def EpsilonGreedy_Policy(Qvalues, allowed_a, epsilon):
17     """
18     returns: tuple
19         an action in form of a one-hot encoded
20         vector with the same shape as Qvalues.
21         an action as decimal integer (0-based)
22
23     Assumes only a single state, i.e. online
24     learning and NOT (mini-)batch learning.
25     """
26     # get the Qvalues and the indices (relative of
27     # all Qvalues) for the allowed actions
28     allowed_a_ind = np.where(allowed_a==1)[0]
29     Qvalues_allowed = Qvalues[allowed_a_ind]
30
31     # ----- epsilon greedy -----
32
33     # draw a random number and compare it to epsilon
34     rand_value = np.random.uniform(0, 1, 1)
35
36     if rand_value < epsilon: # if the random number
37         is smaller than epsilon, draw a random
38         action
39         action_taken_ind_of_allwed_only = np.random.
40             randint(0, len(allowed_a_ind))
41     else: # greedy action
42         action_taken_ind_of_allwed_only = np.argmax(
43             Qvalues_allowed)
```

```
# get index of the action that was chosen (
# 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 shape 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):
```

```

108         self.memory = deque(maxlen=capacity)
109
110     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         """
136         activation functions: "relu", "sigmoid", "
137         tanh", None
138         methods: "qlearning", "sarsa", "dqn"
139         """
140         self.D = N_in # input dimension (without
141         bias)
142         self.K = N_h # nr hidden neurons (without
143         bias)
144         self.O = N_a # nr output neurons (letter
145         , not digit 0)
146
147         # store method and seed
148         self.method = method
149         self.seed = seed
150
151         if self.method == "dqn":
152             self.capacity = capacity
153             self.replay_memory = ReplayMemory(
154                 capacity)
155             self.C = C
156
157         # set activation function and gradient
158         function
159         self.act_f_1_name = activation_function_1
160         self.act_f_2_name = activation_function_2
161         self.act_f_1, self.grad_act_f_1 =
162         act_f_and_gradient(activation_function_1)
163         self.act_f_2, self.grad_act_f_2 =
164         act_f_and_gradient(activation_function_2)
165
166         # initialize the weights and biases and seed
167         global seed
168         np.random.seed(self.seed)
169
170         # self.W1 = np.random.randn(self.K+1, self.D
171         +1)/np.sqrt(self.D+1) # standard norm
172         distribution, shape: (K+1, D+1)
173         # glorot/xavier normal initialization
174         self.W1 = np.random.randn(self.K+1, self.D
175         +1)*np.sqrt(2/ (self.D+1 + self.K+1)) #
176         standard normal distribution, shape:
177         (K+1, D+1)
178
179         self.W2 = np.random.randn(self.O, self.K
180         +1)/np.sqrt(self.K+1) # standard normal
181         distribution, shape: (O, K+1)
182         # glorot/xavier normal initialization
183         self.W2 = np.random.randn(self.O, self.K
184         +1) # standard normal distribution,
185         shape: (O, K+1)
186
187         if self.method == "dqn":
188             self.W1_target = np.copy(self.W1)
189             self.W2_target = np.copy(self.W2)
190
191     def forward(self, x, target=False):
192
193         """
194         x has shape: (D+1, 1) (constant bias 1 must
195         be added beforehand added)
196         target: if True, use the weights of the
197         target network
198
199         returns:
200         last logits (i.e. Qvalues) of shape (O,
201         1)
202         """
203
204         if target == True:
205             W1 = np.copy(self.W1_target)
206             W2 = np.copy(self.W2_target)
207         else:
208             W1 = np.copy(self.W1)
209             W2 = np.copy(self.W2)
210
211         # forward pass/propagation
212         a1 = W1 @ x
213         h1 = self.act_f_1(a1)
214         h1[0,:] = 1 # set first row (bias to second
215         layer) to 1 (this ignores the weights
216         for the k+1th hidden neuron, because
217         this should not exist; this allows to
218         only use matrix multiplication and
219         simplify the gradients as we only need 2
220         instead of 4)
221         a2 = W2 @ h1
222         h2 = self.act_f_2(a2)
223         return a1, h1, a2, h2
224
225     def backward(self, R, x, Qvalues, Q_prime, a1,
226                 h1, a2, gamma, future_reward,
227                 action_binary_mask):
228
229         """
230         backward for methods "qlearning" and "sarsa"
231
232         x has shape (D+1, 1) (constant bias 1 must
233         be added beforehand)
234         set future_reward=True for future reward
235         with gamma>0, False for immediate reward
236         .
237         Q_prime must be chosen according to the
238         method on x_prime (on- or off-policy)
239         """

```



```

207 # backward pass/backpropagation 268
208 # compute the gradient of the square loss 269
    with respect to the parameters 270
209
210 # ===== compute TD error (aka delta) ===== 271
211
212 # make reward of shape (O, 1) 272
213 R_rep = np.tile(R, (self.O, 1)) 273
214 if future_reward: # future reward 274
215     delta = R_rep + gamma*Q_prime - Qvalues 275
    # -> shape (O, 1) 276
216 else: # immediate reward 277
217     delta = R_rep - Qvalues # -> shape (O 278
    1) 279
218
219 # update only action that was taken, i.e. 280
    all rows apart from the one 281
    corresponding to the action taken ( 282
    action index) are 0 283
220 delta = delta*action_binary_mask 284
221
222
223 self.compute_gradients(delta, a1, h1, a2, 285
224 self.update_parameters(self.eta) 286
225
226
227 def backward_dqn(self, batch, gamma): 287
228     """ 288
229     backward for method "dqn" 289
230     """ 290
231
232 # ===== compute targets y and feature matrix 291
    X ===== 292
233
234 # turn batch into individual tuples, numpy 293
    arrays, or lists 294
235 states = batch.state 295
236 rewards = np.array(list(batch.reward)) 296
237 actions = np.array(list(batch.action)) 297
238 next_states = list(batch.next_state) 298
239 dones = np.array(list(batch.done)) 299
240
241 # compute targets y and feature matrix X 300
242 y = np.zeros((self.O, len(dones))) 301
243 for j in np.arange(len(dones)): 302
244     if dones[j]: # if done, set y_j = r_j 303
245         y[actions[j], j] = rewards[j] 304
246     else: 305
247         # compute Q_prime 306
248         Q_target = self.forward(next_state 307
    j], target=True)[-1] 308
249         y[actions[j], j] = rewards[j] + 309
    gamma*np.max(Q_target) 310
250
251
252 # convert states to feature matrix X 311
253 X = np.hstack((states)) 312
254
255
256 # ===== compute TD error (aka delta) ===== 313
257
258 a1, h1, a2, Qvalues = self.forward(X) 314
259 delta = y - Qvalues # -> shape (O, 315
    batch_size) 316
260
261 self.compute_gradients(delta, a1, h1, a2, 317
262 self.update_parameters(self.eta) 318
263
264
265 def compute_gradients(self, delta, a1, h1, a2, 319
266 ): 320
267 # ===== compute gradient of the loss with 321
    respect to the weights ===== 322
268
269
270 # common part of the gradient TODO: check 323
    dimensions 324
271 self.dL_da2 = delta * self.grad_act_f_2(a2) 325
272
273 # gradient of loss wrt W2 326
274 self.dL_dW2 = self.dL_da2 @ h1.T 327
275
276 # gradient of loss wrt W1 328
277 self.dL_dW1 = ( (self.W2.T @ self.dL_da2) * 329
    self.grad_act_f_1(a1) ) @ x.T 330
278
279 def update_parameters(self, eta): 331
280
281 # gradient clipping 332
282
283 # dL_dW1_norm = np.linalg.norm(self.dL_dW1) 333
284 # if dL_dW1_norm >= self.gradient_clip: 334
285     self.dL_dW1 = self.gradient_clip * 335
    self.dL_dW1 / dL_dW1_norm 336
286
287 # dL_dW2_norm = np.linalg.norm(self.dL_dW2) 337
288 # if dL_dW2_norm >= self.gradient_clip: 338
289     self.dL_dW2 = self.gradient_clip * 339
    self.dL_dW2 / dL_dW2_norm 340
290
291 # update W1 and W2 341
292 self.W2 = self.W2 + eta * self.dL_dW2 342
293 self.W1 = self.W1 + eta * self.dL_dW1 343
294
295 def train(self, env, N_episodes, eta, epsilon_0, 344
    beta, gamma, alpha=0.001, gradient_clip=1, 345
    batch_size=32, run_number=None): 346
296
297     """ 347
298     alpha is used as weight for the exponential 348
    moving average displayed during training 349
299
300     batch_size is only used for the DQN method. 350
301     """ 351
302
303 # add training hyper parameters 352
304 self.N_episodes = N_episodes 353
305 self.eta = eta 354
306 self.epsilon_0 = epsilon_0 355
307 self.beta = beta 356
308 self.gamma = gamma 357
309 self.alpha = alpha 358
310 self.gradient_clip = gradient_clip 359
311 self.batch_size = batch_size 360
312
313
314 training_start = time.time() 361
315
316 try: 362
317
318     # initialize histories for important 363
    metrics 364
319     self.R_history = np.full([self. 365
    N_episodes, 1], np.nan) 366
320     self.N_moves_history = np.full([self. 367
    N_episodes, 1], np.nan) 368
321     self.dL_dW1_norm_history = np.full([self. 369
    N_episodes, 1], np.nan) 370
322     self.dL_dW2_norm_history = np.full([self. 371
    N_episodes, 1], np.nan) 372
323
324 # progress bar 373
325 episodes = tqdm(np.arange(self. 374
    N_episodes), unit="episodes") 375
326 ema_previous = 0 376

```

```

328                                     367
329 n_steps = 0                                     368
330
331 for n in episodes:                               369
332     epsilon_f = self.epsilon_0 / (1 + 370
333         beta * n)  ## DECAYING EPSILON 371
334     Done = 0                                     372
335                                     373
336                                     374
337                                     375
338                                     376
339                                     377
340                                     378
341                                     379
342                                     380
343                                     381
344                                     382
345                                     383
346                                     384
347                                     385
348
349 while Done==0:
350                                     ##
351                                     START THE EPISODE
352                                     386
353                                     387
354                                     388
355                                     389
356                                     390
357                                     391
358                                     392
359                                     393
360                                     394
361                                     395
362                                     396
363                                     397
364                                     398
365                                     399
366                                     400
367                                     401
368                                     402
369                                     403
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999

```



```

404         episodes.set_description(
405             (f"EMA Reward = {ema_reward:.2f}")
406         )
407         break
408     else: # IF THE EPISODE IS NOT OVER...
409         if self.method == "qlearning":
410             # chose next action off policy
411             Q_prime = np.max(self.forward(X_prime)[-1])
412             elif self.method == "sarsa":
413                 # chose next action on policy
414                 a1_prime, h1_prime, a2_prime,
415                 Qvalues_prime = self.forward(X_prime)
416                 -> shape (N_a, 1)
417                 # chose next action and save it
418                 A_binary_mask_prime, A_ind_prime =
419                 EpsilonGreedyPolicy(Qvalues_prime,
420                                     allowed_a_prime,
421                                     epsilon_f)
422                 # get Qvalue of next action
423                 Q_prime = Qvalues_prime[A_ind_prime]
424                 if (self.method == "qlearning") or (self.method == "sarsa"):
425                     # backpropagation and weight update
426                     self.backward(R, X, Qvalues, Q_prime, a1,
427                                 h1, a2, self.gamma, future_reward=True,
428                                 action_binary_mask=A_binary_mask)
429                 # NEXT STATE AND CO. BECOME ACTUAL STATE...
430                 if self.method == "sarsa":
431                     A_binary_mask = np.copy(A_binary_mask_prime)
432                     A_ind = np.copy(A_ind_prime)
433                     a1 = np.copy(a1_prime)
434                     h1 = np.copy(h1_prime)
435                     a2 = np.copy(a2_prime)
436                     Qvalues = np.copy(Qvalues_prime)
437                     S = np.copy(S_prime)
438                     X = np.copy(X_prime)
439                     allowed_a = np.copy(allowed_a_prime)
440                     i += 1 # UPDATE COUNTER FOR NUMBER OF ACTIONS
441
442         if (self.method == "dqn") and (n_steps % self.C == 0):
443             # update target network every C steps
444             self.W1_target = np.copy(self.W1)
445             self.W2_target = np.copy(self.W2)
446
447         training_end = time.time()
448         self.training_time_in_seconds = training_end - training_start
449
450         return None
451
452     except KeyboardInterrupt as e:
453         # return nothing
454         training_end = time.time()
455         self.training_time_in_seconds = training_end - training_start
456
457         return None
458
459     def save(self, name_extension=None):
460         # create directory for the model
461         name = f"{self.method}_{self.act_f1_name}_{self.act_f2_name}"
462         if name_extension is not None:
463             name += f"_{name_extension}"
464
465         path = f"models/{name}"
466         if not os.path.isdir(path): os.mkdir(path)
467         print(f"saving to: {path}")
468
469         # save weights
470         np.save(f"{path}/W1.npy", self.W1)
471         np.save(f"{path}/W2.npy", self.W2)
472
473         # save training history
474         np.save(f"{path}/training_history_R.npy", self.R_history)
475         np.save(f"{path}/training_history_N_moves.npy", self.N_moves_history)
476         np.save(f"{path}/training_history_dL_dW1_norm.npy", self.dL_dW1_norm_history)
477         np.save(f"{path}/training_history_dL_dW2_norm.npy", self.dL_dW2_norm_history)
478
479         # save training parameters and other general info
480         params = {
481             "method": self.method,
482             "N_episodes": self.N_episodes,
483             "eta": self.eta,
484             "epsilon_0": self.epsilon_0,
485             "beta": self.beta,
486             "gamma": self.gamma,
487             "alpha": self.alpha,
488             # "gradient_clip": self.gradient_clip,
489             "seed": self.seed,
490             "D": self.D,
491             "K": self.K,
492             "O": self.O,
493             "training_time_in_seconds": self.training_time_in_seconds
494         }
495         if self.method == "dqn":
496             params["capacity"] = self.capacity

```

```

503         params["batch_size"] = self.batch_size
504         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)

```

```

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.