# Chess Assignment

van den Bergh Laurin
*Institute of Informatics*
*University of Zurich*
Zurich, Switzerland
laurin.vandenbergh@uzh.ch — Matriculation number: 16-744-401

*Abstract*—**This document is a model and instructions for LATEX. This and the IEEEtran.cls file define the components of your paper [title, text, heads, etc.]. \*CRITICAL: Do Not Use Symbols, Special Characters, Footnotes, or Math in Paper Title or Abstract. 1 sentence summary, 1 sentence method, and 1 sentence results, and 1 sentence conclusion.**

*Index Terms*—**reinforcement learning, chess**

## I. INTRODUCTION

### A. - describe general setting

In this assignment, we explore three different reinforcement learning algorithms to learn how to play a simplified version of chess. These three algorithms are SARSA, Q-Learning and DQN[1].

### B. - touch on methods

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.

### C. - touch on results and conclusion

seeds and nonseeded experiments foreshadow results from qlearning

## II. METHODS

### A. Environment

This version of chess takes place on a 4 by 4 board and can be thought of as a late game 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. For DQN no sequence is needed, because chess fulfills Markov property and gives rise to a Markov Decision Process (MDP). Also no preprocessing of sequence is needed anymore because state encoding is already given. [1] [2]

---

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

Chess has/fulfills the Markov property, so applying reinforcement learning algorithms like Q-Learning and SARSA is theoretically justified.

(1) Describe the algorithms Q-learning and SARSA. Explain the differences and the possible advantages/disadvantages between the two methods. No coding is needed to reply to this question.

### B. SARSA and Q-Learning[2]

*1) Temporal-Difference Algorithms:* maybe shorten 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. why? else what?

*2) differences/what they have in common:* SARSA and Q-Learning address the temporal-credit-assignment problem [3], that is, trying to attribute future rewards to previous actions. These future rewards get discounted with the hyper-parameter $\gamma$ (see Section III-C). Both algorithms repeatedly choose and take actions in the environment according to some policy $\pi$, e.g. an $\epsilon$-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 $\pi$, 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 $\pi$.

Q-learning, on the other hand, is an off-policy algorithm, which means that it takes its actions $a$ according to its policy $\pi$, but during the update steps it assumes a greedy policy, i.e.

---

[2]Answer to task 1.

optimal play, for future actions $a'$. Q-Learning has the update rule:

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

*3) - advangages/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. $Q_*(s,a) = \max_\pi Q_\pi(s,a)$, and in a greedy stetting 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. [4] 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 [2], 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 Replay[3]

replay memory of fixed size, queue (get code in appendix) (2) [Group Only] Describe the experience replay technique. Cite relevant papers.

Experience replay is a technique proposed by Lin [5] 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 [1], [2].

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 then once, leading to increased data efficiency and faster convergence [1], [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 [1]. 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. [1] 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. [2] later improved upon the method and presented the DQN algorithm, as it is known today, where they address the problem of the Q-values $Q(s,a)$ being correlated to the target values $y = r + \gamma \max_{a'} Q(s', a')$ 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.

### E. Experiments

In order to address all tasks, we divided the tasks into several independent experiments. First, we conducted seeded runs[5] 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-1). 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 $\beta$ and $\gamma$[6] we trained 49 agents with different combinations for $\beta$ and $\gamma$ 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 (see Figures 4 and 5). These runs are seeded identically to the seeded runs mentioned above.

### F. Implementation and Hyper-parameters

*1) Implementation:* We implemented all algorithms from scratch according to Sutton et al. [4] (SARSA and Q-

---

[3]Answer to "group only" task 2.

[4]Answer to task 5: Describing the used method.
[5]Answers to task 3 and 5.
[6]Answer to task 4.

Learning[7]) and Mnih et al. [2] (DQN[8]). For the implementation see file `neural_net.py` 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 $\frac{2}{n_{\text{in}}+n_{\text{out}}}$, where $n_{\text{in}}$ and $n_{\text{out}}$ 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

*2) Hyper-Parameters:* For all experiments we used the default hyperparameters provided in the `Assignment.ipynb` file unless otherwise noted. 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.

## III. RESULTS

### A. Seeded Runs[9]

(3) We provide you with a template code of the chess game. Fill the gaps in the program file chess implementing either Q Learning or SARSA; detailed instructions are available inside the file. We provide indicative parameter values. Produce two plots that show the reward per game and the number of moves per game vs training time. The plots will be noisy. Use an exponential moving average.

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 $\epsilon$ 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.

---

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

[8]Answer to task 5.

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



Fig. 1: Exponential moving average of the rewards achieved during training. All three algorithms, Q-Learning, SARSA and DQN were initialized with identical weights and trained with identical network architecture and hyper-parameters.



Fig. 2: caption

(5) Implement another Deep RL algorithm (it could be SARSA if you chose before Q Learning and vice versa). Compare the new results with your previous results.

As suggested by Mnih et al. [1], [2], DQN[10] 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. For the learning curves see Figures 4 and 5 in the Appendix.

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

---

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

some runs the rewards consistently dropped while the number of moves consistently increased. This shows that the measures taken by Mnih et al. [2] to combat the disadvantages of Q-Learning worked and increased the stability as well as the convergence speed.

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

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.

*C. Hyper-parameters*

(4) Change the discount factor $\gamma$ and the speed $\beta$ of the decaying trend of $\epsilon$ (for a definition of $\beta$ please see code). Analyse the results.

*1) gamma:* Figure 3(a) depicts the rewards and number of moves per episode as a function of $\beta$ and $\gamma$. We can see that the reward increases monotonically as $\gamma$ is increases, suggesting that a value of $\gamma \in [0.80, 1)$ should be chosen. This intuitively makes sense, as we have very sparse rewards and want the agent to "backpropagate" this reward through its sequence of actions.

*2) beta:* We can not see a clear relationship between $\beta$ and the rewards, apart from $\beta = 0$ being an inferior choice. In Figure 3(b) we can see that the that the number of steps taken by the agent decreases drastically when increasing $\gamma$ from very low levels, but again there is no clear pattern visible for $\beta$. In summary, for reasonalby chosen values of $\gamma$ the choice of $\beta$ seems to not have much of an influence for training periods up to 40000 episodes.

(6) [Group Only] Change the state representation or the administration of reward. Interpret your results.

(7) [Group Only] The values of the weights could explode during learning. This issue is called exploding gradients. Explain one possible way to fix this issue and implement it. For example, search and study RMSprop). Show in a plot how your solution fixed this issue.

## 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.

summarize results, takeaway

get correct bib format and complete entries



Fig. 3

## REFERENCES

[1] 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

[2] 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

[3] R. S. Sutton, "Temporal credit assignment in reinforcement learning," 1984.

[4] 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

[5] 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

[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.

## V. APPENDIX

*1) Reproducibility:* put the appendix after the bibliography?

```
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
```

Fig. 4



Fig. 5

```
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
              vector with the same shapeensions as
              Qvalues.
22          an action as decimal integer (0−based)
23
24      Assumes only a single state, i.e. online
              learning and NOT (mini−)batch learning.
25      """
26      # get the Qvalues and the indices (relative of
              all Qvalues) for the allowed actions
```

```
27  allowed_a_ind = np.where(allowed_a==1)[0]
28  Qvalues_allowed = Qvalues[allowed_a_ind]
29
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
            is smaller than epsilon, draw a random
            action
        action_taken_ind_of_allwed_only = np.random.
              randint(0, len(allowed_a_ind))
37  else:  # greedy action
        action_taken_ind_of_allwed_only = np.argmax(
              Qvalues_allowed)
38
39  # get index of the action that was chosen (
            relative to all actions, not only allowed)
40  ind_of_action_taken = allowed_a_ind[
            action_taken_ind_of_allwed_only]
41
42
43
44  # ———————— create usable output ————————
45
```

```python
46
47      # get the shapeensions of the Qvalues
48      N_a, N_samples = np.shape(Qvalues)  # N_samples
            must be 1
49
50      # initialize all actions of binary mask to 0
51      A_binary_mask = np.zeros((N_a,N_samples))
52      # set the action that was chosen to 1
53      A_binary_mask[ind_of_action_taken,:] = 1
54
55      return A_binary_mask, ind_of_action_taken
56
57
58
59  # ===== activation functions and it's derivatives
        ======
60
61  # relu and its derivative
62  def relu(x):
63      return np.maximum(0,x)
64
65  def heaviside(x):
66      return np.heaviside(x,0)
67
68  # sigmoid and its derivative
69  def sigmoid(x):
70      return 1 / (1 + np.exp(−x))
71
72  def gradient_sigmoid(x):
73      return sigmoid(x) * (1 − sigmoid(x))
74
75  # tanh and its derivative
76  def tanh(x):
77      return np.tanh(x)
78
79  def gradient_tanh(x):
80      return 1 − np.tanh(x)**2
81
82  # identity and its derivative
83  def identity(x):
84      return x
85
86  def const(x):
87      return np.ones(x.shape)
88
89
90  def act_f_and_gradient(activation_function="relu"):
91      if activation_function == "relu":
92          return relu, heaviside
93      elif activation_function == "sigmoid":
94          return sigmoid, gradient_sigmoid
95      elif activation_function == "tanh":
96          return tanh, gradient_tanh
97      else:  # identity and constant 1
98          return identity, const
99
100
101
102 # ===== Replay Memory for Experience Replay (with
        DQN) =====
103
104 Transition = namedtuple('Transition', ("state", "
        action", "reward", "next_state", "done"))
105
106 class ReplayMemory(object):
107     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
                data
115         if len(self) < batch_size:
116             batch_size = len(self)
117         return random.sample(self.memory, batch_size
                )
118
119     def __len__(self):
120         return len(self.memory)
121
122
123
124
125 # ===== Neural Network ======
126
127 class NeuralNetwork(object):
128
129     def __init__(self, N_in, N_h, N_a,
                activation_function_1="relu",
                activation_function_2=None, method="
                qlearning", seed=None, capacity=100_000, C
                =100):
130         """
131         activation functions: "relu", "sigmoid", "
                tanh", None
132         methods: "qlearning", "sarsa", "dqn"
133         """
134         self.D = N_in  # input dimension (without
                bias)
135         self.K = N_h   # nr hidden neurons (without
                bias)
136         self.O = N_a   # nr output neurons (letter O
                , not digit 0)
137
138         # store method and seed
139         self.method = method
140         self.seed = seed
141
142         if self.method == "dqn":
143             self.capacity = capacity
144             self.replay_memory = ReplayMemory(
                    capacity)
145             self.C = C
146
147         # set activation function and gradient
                function
148         self.act_f_1_name = activation_function_1
149         self.act_f_2_name = activation_function_2
150         self.act_f_1, self.grad_act_f_1 =
                act_f_and_gradient(activation_function_1
                )
151         self.act_f_2, self.grad_act_f_2 =
                act_f_and_gradient(activation_function_2
                )
152
153
154         # initialize the weights and biases and set
                grobal seed
155         np.random.seed(self.seed)
156
157         # self.Wl = np.random.randn(self.K+1, self.D
                +1)/np.sqrt(self.D+1)  # standard normal
                distribution, shape: (K+1, D+1)
158         # glorot/xavier normal initialization
159         # self.Wl = np.random.randn(self.K+1, self.D
                +1)*np.sqrt(2/ (self.D+1 + self.K+1))  #
                standard normal distribution, shape: (K
                +1, D+1)
160         self.Wl = np.random.standard_normal((self.K
                +1, self.D+1))*np.sqrt(2/ (self.D+1 +
                self.K+1))  # standard normal
                distribution, shape: (K+1, D+1)
161         # self.Wl = np.random.randn(self.K+1, self.D
                +1)  # standard normal distribution,
                shape: (K+1, D+1)
162
```

```python
163            # self.W2 = np.random.randn(self.O, self.K
                   +1)/np.sqrt(self.K+1)  # standard normal
                   distribution, shape: (O, K+1)
164            # glorot/xavier normal initialization
165            self.W2 = np.random.standard_normal((self.O,
                   self.K+1))*np.sqrt(2/ (self.K+1 + self.
                   O))  # standard normal distribution,
                   shape: (O, K+1)
166            # self.W2 = np.random.randn(self.O, self.K
                   +1)  # standard normal distribution,
                   shape: (O, K+1)
167
168            if self.method == "dqn":
169                self.W1_target = np.copy(self.W1)
170                self.W2_target = np.copy(self.W2)
171
172
173        def forward(self, x, target=False):
174            """
175            x has shape: (D+1, 1) (constant bias 1 must
                   be added beforehand added)
176            target: if True, use the weights of the
                   target network
177
178            returns:
179                last logits (i.e. Qvalues) of shape (O,
                   1)
180            """
181
182            if target == True:
183                W1 = np.copy(self.W1_target)
184                W2 = np.copy(self.W2_target)
185            else:
186                W1 = np.copy(self.W1)
187                W2 = np.copy(self.W2)
188
189            # forward pass/propagation
190            a1 = W1 @ x
191            h1 = self.act_f_1(a1)
192            h1[0,:] = 1  # set first row (bias to second
                   layer) to 1 (this ignores the weights
                   for the k+1th hidden neuron, because
                   this should not exist; this allows to
                   only use matrix multiplication and
                   simplify the gradients as we only need 2
                   instead of 4)
193            a2 = W2 @ h1
194            h2 = self.act_f_2(a2)
195            return a1, h1, a2, h2
196
197
198        def backward(self, R, x, Qvalues, Q_prime, a1,
                   h1, a2, gamma, future_reward,
                   action_binary_mask):
199            """
200            backward for methods "qlearning" and "sarsa"
201
202            x has shape (D+1, 1) (constant bias 1 must
                   be added beforehand)
203            set future_reward=True for future reward
                   with gamma>0, False for immediate reward
                   .
204            Q_prime must be chosen according to the
                   method on x_prime (on- or off-policy)
205            """
206
207            # backward pass/backpropagation
208            # compute the gradient of the square loss
                   with respect to the parameters
209
210            # ===== compute TD error (aka delta) =====
211
212            # make reward of shape (O, 1)
213            R_rep = np.tile(R, (self.O, 1))
214            if future_reward:  # future reward
215                delta = R_rep + gamma*Q_prime - Qvalues
                       # -> shape (O, 1)
216            else:  # immediate reward
217                delta = R_rep - Qvalues  # -> shape (O,
                       1)
218
219            # update only action that was taken, i.e.
                   all rows apart from the one
                   corresponding to the action taken (
                   action index) are 0
220            delta = delta*action_binary_mask
221
222
223            self.compute_gradients(delta, a1, h1, a2, x)
224            self.update_parameters(self.eta)
225
226
227        def backward_dqn(self, batch, gamma):
228            """
229            backward for method "dqn"
230            """
231
232            # ===== compute targets y and feature matrix
                   X =====
233
234            # turn batch into individual tuples, numpy
                   arrays, or lists
235            states = batch.state
236            rewards = np.array(list(batch.reward))
237            actions = np.array(list(batch.action))
238            next_states = list(batch.next_state)
239            dones = np.array(list(batch.done))
240
241            # compute targets y and feature matrix X
242            y = np.zeros((self.O, len(dones)))
243            for j in np.arange(len(dones)):
244                if dones[j]:  # if done, set y_j = r_j
245                    y[actions[j], j] = rewards[j]
246                else:
247                    # compute Q_prime
248                    Q_target = self.forward(next_states[
                           j], target=True)[-1]
249                    y[actions[j], j] = rewards[j] +
                           gamma*np.max(Q_target)
250
251            # convert states to feature matrix X
252            X = np.hstack((states))
253
254
255            # ===== compute TD error (aka delta) =====
256
257            a1, h1, a2, Qvalues = self.forward(X)
258            delta = y - Qvalues  # -> shape (O,
                   batch_size)
259
260            self.compute_gradients(delta, a1, h1, a2, X)
261            self.update_parameters(self.eta)
262
263
264        def compute_gradients(self, delta, a1, h1, a2, x
                   ):
265            # ===== compute gradient of the loss with
                   respect to the weights =====
266
267            # common part of the gradient  TODO: check
                   dimensions
268            self.dL_da2 = delta * self.grad_act_f_2(a2)
269
270            # gradient of loss wrt W2
271            self.dL_dW2 = self.dL_da2 @ h1.T
272
273            # gradient of loss wrt W1
274
```

```python
275         self.dL_dW1 = ( (self.W2.T @ self.dL_da2) *
                self.grad_act_f_1(a1) ) @ x.T
276
277
278
279     def update_parameters(self, eta):
280
281         # gradient clipping
282
283         # dL_dW1_norm = np.linalg.norm(self.dL_dW1)
284         # if dL_dW1_norm >= self.gradient_clip:
285         #     self.dL_dW1 = self.gradient_clip *
                self.dL_dW1 / dL_dW1_norm
286
287         # dL_dW2_norm = np.linalg.norm(self.dL_dW2)
288         # if dL_dW2_norm >= self.gradient_clip:
289         #     self.dL_dW2 = self.gradient_clip *
                self.dL_dW2 / dL_dW2_norm
290
291         # update W1 and W2
292         self.W2 = self.W2 + eta * self.dL_dW2
293         self.W1 = self.W1 + eta * self.dL_dW1
294
295
296
297
298     def train(self, env, N_episodes, eta, epsilon_0,
            beta, gamma, alpha=0.001, gradient_clip=1,
            batch_size=32, run_number=None):
299         """
300         alpha is used as weight for the exponential
                moving average displayed during training
                .
301         batch_size is only used for the DQN method.
302         """
303
304         # add training hyper parameters
305         self.N_episodes = N_episodes
306         self.eta = eta
307         self.epsilon_0 = epsilon_0
308         self.beta = beta
309         self.gamma = gamma
310         self.alpha = alpha
311         self.gradient_clip = gradient_clip
312         self.batch_size = batch_size
313
314
315         training_start = time.time()
316
317         try:
318
319             # initialize histories for important
                    metrics
320             self.R_history = np.full([self.
                    N_episodes, 1], np.nan)
321             self.N_moves_history = np.full([self.
                    N_episodes, 1], np.nan)
322             self.dL_dW1_norm_history = np.full([self
                    .N_episodes, 1], np.nan)
323             self.dL_dW2_norm_history = np.full([self
                    .N_episodes, 1], np.nan)
324
325             # progress bar
326             episodes = tqdm(np.arange(self.
                    N_episodes), unit="episodes")
327             ema_previous = 0
328
329             n_steps = 0
330
331             for n in episodes:
332
333                 epsilon_f = self.epsilon_0 / (1 +
                        beta * n)    ## DECAYING EPSILON
334                 Done = 0

335                 ## SET DONE TO ZERO (BEGINNING
                    OF THE EPISODE)
                    i = 1

336                 ## COUNTER FOR NUMBER OF ACTIONS

337             S, X, allowed_a = env.
                    Initialise_game()       ##
                    INITIALISE GAME
338             X = np.expand_dims(X, axis=1)
                        ## MAKE X A
                    TWO DIMENSIONAL ARRAY
339             X = np.copy(np.vstack((np.array
                    ([[1]]), X)))  # add bias term

                    if self.method == "sarsa":
                        # compute Q values for the given
                            state
                        a1, h1, a2, Qvalues = self.
                            forward(X)   # -> shape (O,
                            1)

                        # choose an action A using
                            epsilon-greedy policy
                        A_binary_mask, A_ind =
                            EpsilonGreedy_Policy(Qvalues
                            , allowed_a, epsilon_f)  #
                            -> shape (O, 1)

                    while Done==0:
                                                ##
                        START THE EPISODE

                        if (self.method == "qlearning")
                            or (self.method == "dqn"):
                            # compute Q values for the
                                given state
                            a1, h1, a2, Qvalues = self.
                                forward(X)   # -> shape (
                                O, 1)

                            # choose an action A using
                                epsilon-greedy policy
                            A_binary_mask, A_ind =
                                EpsilonGreedy_Policy(
                                Qvalues, allowed_a,
                                epsilon_f)  # -> shape (
                                O, 1)

                        # take action and observe reward
                            R and state S_prime
                        S_prime, X_prime,
                            allowed_a_prime, R, Done =
                            env.OneStep(A_ind)
                        X_prime = np.expand_dims(X_prime
                            , axis=1)
                        X_prime = np.copy(np.vstack((np.
                            array([[1]]), X_prime)))  #
                            add bias term

                        n_steps += 1

                        if self.method == "dqn":

                            # store the transition in
                                memory
                            self.replay_memory.push(X,
                                A_ind, R, X_prime, Done)

                            # sample a batch of
                                transitions
```

```python
372            transactions = self.
                 replay_memory.sample(
                 self.batch_size)
373            # turn list of transactions
                 into transaction of
                 lists
374            batch = Transition(*zip(*
                 transactions))
375
376            # backward step and
                 parameter update
377            self.backward_dqn(batch,
                 self.gamma)
378
379            # update Q values indirectly by
                 updating the weights and
                 biases directly
380
381            if Done==1:  # THE EPISODE HAS
                 ENDED, UPDATE...BE CAREFUL
                 THIS IS THE LAST STEP OF THE
                 EPISODE
382
383                if (self.method == "
                     qlearning") or (self.
                     method == "sarsa"):
384                    # compute gradients and
                         update weights
385                    self.backward(R, X,
                         Qvalues, None, a1,
                         h1, a2, None,
                         future_reward=False,
                         action_binary_mask=
                         A_binary_mask)
386
387                # store history
388                # todo: record max possible
                     reward per episode
389                self.R_history[n] = np.copy(
                     R)  # reward per episode
390                self.N_moves_history[n] = np
                     .copy(i)  # nr moves per
                     episode
391
392                # store norm of gradients
393                self.dL_dW1_norm_history[n]
                     = np.linalg.norm(self.
                     dL_dW1)
394                self.dL_dW2_norm_history[n]
                     = np.linalg.norm(self.
                     dL_dW2)
395
396                # compute exponential moving
                     average (EMA) to
                     display during training
397                ema = alpha*R + (1-alpha)*
                     ema_previous
398                if n == 0: # first episode
399                    ema = R
400                ema_previous = ema
401                if run_number is not None:
402                    episodes.set_description
                         (f"Run = {run_number
                         }; EMA Reward = {ema
                         :.2f}")
403                else:
404                    episodes.set_description
                         (f"EMA Reward = {ema
                         :.2f}")
405
406                break
407
408            else:  # IF THE EPISODE IS NOT
                 OVER...
```

```python
409                if self.method == "qlearning
                     ":
410                    # chose next action off-
                         policy
411                    Q_prime = np.max(self.
                         forward(X_prime)
412                         [-1])
413
414                elif self.method == "sarsa":
415                    # chose next action on-
                         policy
416
417                    a1_prime, h1_prime,
                         a2_prime,
418                         Qvalues_prime = self
                         .forward(X_prime)  #
                         -> shape (N_a, 1)
419
420                    # chose next action and
                         save it
421                    A_binary_mask_prime,
                         A_ind_prime =
                         EpsilonGreedy_Policy
                         (Qvalues_prime,
                         allowed_a_prime,
                         epsilon_f)
422
423                    # get Qvalue of next
                         action
424                    Q_prime = Qvalues_prime[
                         A_ind_prime]
425
426
427                if (self.method == "
                     qlearning") or (self.
                     method == "sarsa"):
428                    # backpropagation and
                         weight update
429                    self.backward(R, X,
                         Qvalues, Q_prime, a1
                         , h1, a2, self.gamma
                         , future_reward=True
                         , action_binary_mask
                         =A_binary_mask)
430
431                # NEXT STATE AND CO. BECOME
                     ACTUAL STATE...
432                if self.method == "sarsa":
433                    A_binary_mask = np.copy(
                         A_binary_mask_prime)
434                    A_ind = np.copy(
                         A_ind_prime)
435                    a1 = np.copy(a1_prime)
436                    h1 = np.copy(h1_prime)
437                    a2 = np.copy(a2_prime)
438                    Qvalues = np.copy(
                         Qvalues_prime)
439                S = np.copy(S_prime)
440                X = np.copy(X_prime)
441                allowed_a = np.copy(
                     allowed_a_prime)
442
443                i += 1  # UPDATE COUNTER FOR
                     NUMBER OF ACTIONS
444
445        if (self.method == "dqn") and (
               n_steps % self.C == 0):
446            # update target network
                 every C steps
447            self.W1_target = np.copy(
                 self.W1)
448            self.W2_target = np.copy(
```

```python
                        self.W2)
449
450
451             training_end = time.time()
452             self.training_time_in_seconds = \
                    training_end - training_start
453
454             return None
455
456
457         except KeyboardInterrupt as e:
458             # return nothing
459             training_end = time.time()
460             self.training_time_in_seconds = \
                    training_end - training_start
461
462             return None
463
464
465     def save(self, name_extension=None):
466         # create directory for the model
467         name = f"{self.method}_{self.act_f_1_name}_{self.act_f_2_name}"
468         if name_extension is not None:
469             name += f"_{name_extension}"
470
471         path = f"models/{name}"
472         if not os.path.isdir(path): os.mkdir(path)
473         print(f"saving to: {path}")
474
475         # save weights
476         np.save(f"{path}/W1.npy", self.W1)
477         np.save(f"{path}/W2.npy", self.W2)
478
479         # save training history
480         np.save(f"{path}/training_history_R.npy",
                    self.R_history)
481         np.save(f"{path}/training_history_N_moves.
                    npy", self.N_moves_history)
482         np.save(f"{path}/
                    training_history_dL_dW1_norm.npy", self.
                    dL_dW1_norm_history)
483         np.save(f"{path}/
                    training_history_dL_dW2_norm.npy", self.
                    dL_dW2_norm_history)
484
485         # save training parameters and other general
                    info
486         params = {
487             "method": self.method,
488             "N_episodes": self.N_episodes,
489             "eta": self.eta,
490             "epsilon_0": self.epsilon_0,
491             "beta": self.beta,
492             "gamma": self.gamma,
493             "alpha": self.alpha,
494             # "gradient_clip": self.gradient_clip,
495             "seed": self.seed,
496             "D": self.D,
497             "K": self.K,
498             "O": self.O,
499             "training_time_in_seconds": self.
                    training_time_in_seconds
500         }
501         if self.method == "dqn":
502             params["capacity"] = self.capacity
503             params["batch_size"] = self.batch_size
504             params["C"] = self.C
505         with open(f"{path}/training_parameters.json"
                    , "w") as f:
506             json.dump(params, f)
507
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
                    instance
513     name = f"{method}_{act_f_1}_{act_f_2}"
514     if name_extension is not None:
515         name += f"_{name_extension}"
516
517     path = f"models/{name}"
518     # print(f"loading from: {path}")
519
520     # initialize neural network
521     nn = NeuralNetwork(0,0,0, activation_function_1=
                    act_f_1, activation_function_2=act_f_2,
                    method=method)
522
523     # network weights
524     nn.W1 = np.load(f"{path}/W1.npy")
525     nn.W2 = np.load(f"{path}/W2.npy")
526
527     # network training history
528     nn.R_history = np.load(f"{path}/
                    training_history_R.npy")
529     nn.N_moves_history = np.load(f"{path}/
                    training_history_N_moves.npy")
530     nn.dL_dW1_norm_history = np.load(f"{path}/
                    training_history_dL_dW1_norm.npy")
531     nn.dL_dW2_norm_history = np.load(f"{path}/
                    training_history_dL_dW2_norm.npy")
532
533     # network training parameters
534     with open(f"{path}/training_parameters.json", "r
                    ") as f:
535         params = json.load(f)
536
537         # set parameters to the network instance
538         nn.method = params["method"]
539         nn.N_episodes = int(params["N_episodes"])
540         nn.eta = float(params["eta"])
541         nn.epsilon_0 = float(params["epsilon_0"])
542         nn.beta = float(params["beta"])
543         nn.gamma = float(params["gamma"])
544         nn.alpha = float(params["alpha"])
545         # nn.gradient_clip = float(params["
                    gradient_clip"])
546         try:
547             nn.seed = int(params["seed"])
548         except:
549             nn.seed = params["seed"]
550         nn.D = int(params["D"])
551         nn.K = int(params["K"])
552         nn.O = int(params["O"])
553         nn.training_time_in_seconds = float(params["
                    training_time_in_seconds"])
554
555         if nn.method == "dqn":
556             nn.capacity = int(params["capacity"])
557             nn.batch_size = int(params["batch_size"
                    ])
558             nn.C = int(params["C"])
559
560
561     if nn.method == "dqn":
562         nn.W1_target = np.copy(nn.W1)
563         nn.W2_target = np.copy(nn.W2)
564
    return nn
```

Listing 1: caption