

Chess Assignment

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

van den Bergh Laurin

Institute of Informatics

University of Zurich

Zurich, Switzerland

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

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

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

I. INTRODUCTION

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

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

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

II. METHODS

A. Environment

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

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

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

B. SARSA and Q-Learning²

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

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

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

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

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

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

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

²Answer to task 1.

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

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

C. Experience Replay³

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

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

D. Deep Q-Networks (DQN)⁴

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

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

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

E. Experiments

In order to address all tasks, we divided the tasks into several independent experiments. First, we conducted seeded runs⁵ for all three algorithms using seed 21 for reproducibility, which was chosen a-priori. These seeded runs serve as examples to compare the algorithm’s online performance qualitatively. The seeds are used such that the weights of all neural networks are instantiated identically for all algorithms and they subsequently serve as seeds for any random number used during training. This makes sure that all agents start with the same initial conditions and that the results are reproducible (see Section V-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 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 (see Figures 4 and 5). These runs are seeded identically to the seeded runs mentioned above.

F. Implementation and Hyper-parameters

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

³Answer to “group only” task 2.

⁴Answer to task 5: Describing the used method.

⁵Answers to task 3 and 5.

⁶Answer to task 4.

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

⁸Answer to task 5.

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

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

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

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

III. RESULTS

A. Seeded Runs⁹

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

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

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

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

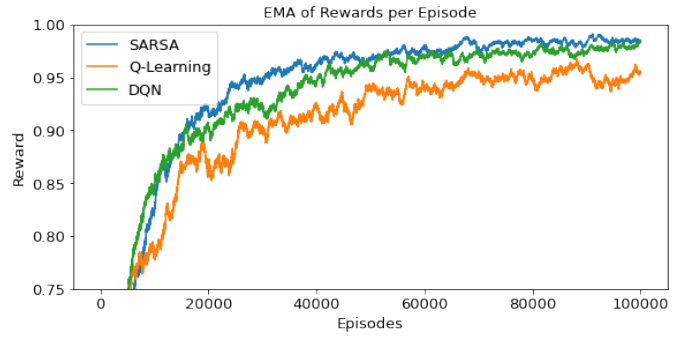


Fig. 1: Exponential moving average of the rewards achieved during training for 100000 episodes with identical hyper-parameters, weight initialization and model architecture.

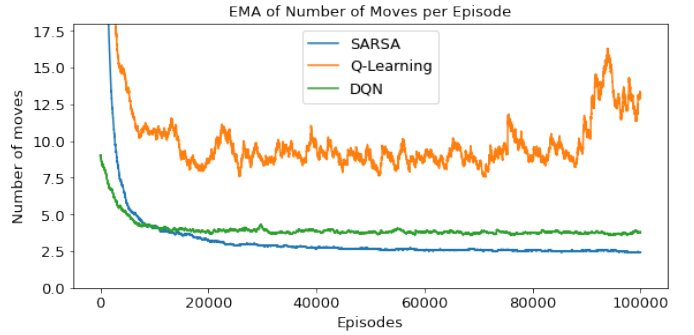


Fig. 2: Exponential moving average of the number of moves per episode achieved during training for 100000 episodes with identical hyper-parameters, weight initialization and model architecture.

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

B. Simulation Study (Non-seeded Runs)

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

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

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

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

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

C. Hyper-parameters¹¹

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

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

IV. CONCLUSION

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

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

All three algorithms were able to learn to play the simplified version of chess to a very high degree even without

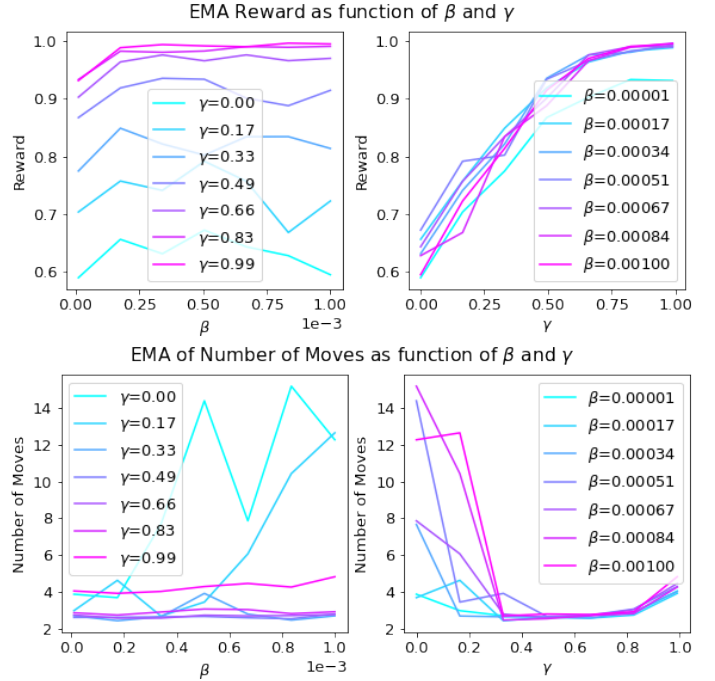


Fig. 3: Rewards and number of moves as functions of the speed of the decaying trend β and the discount factor γ after training a SARSA agent for 40000 episodes.

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

¹¹ Answer to task 4.

V. APPENDIX

1) *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 the repository as well.

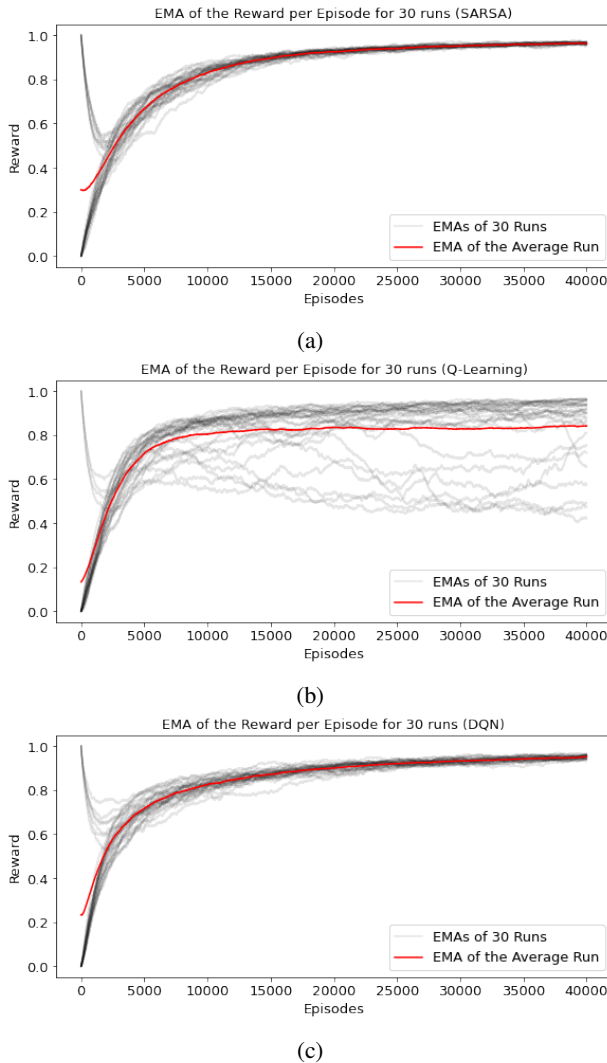


Fig. 4: Expected moving averages (EMA) of the rewards per episode for 30 runs for each algorithm (SARSA, Q-Learning, DQN). The red line depicts the EMA of the average across the 30 runs.

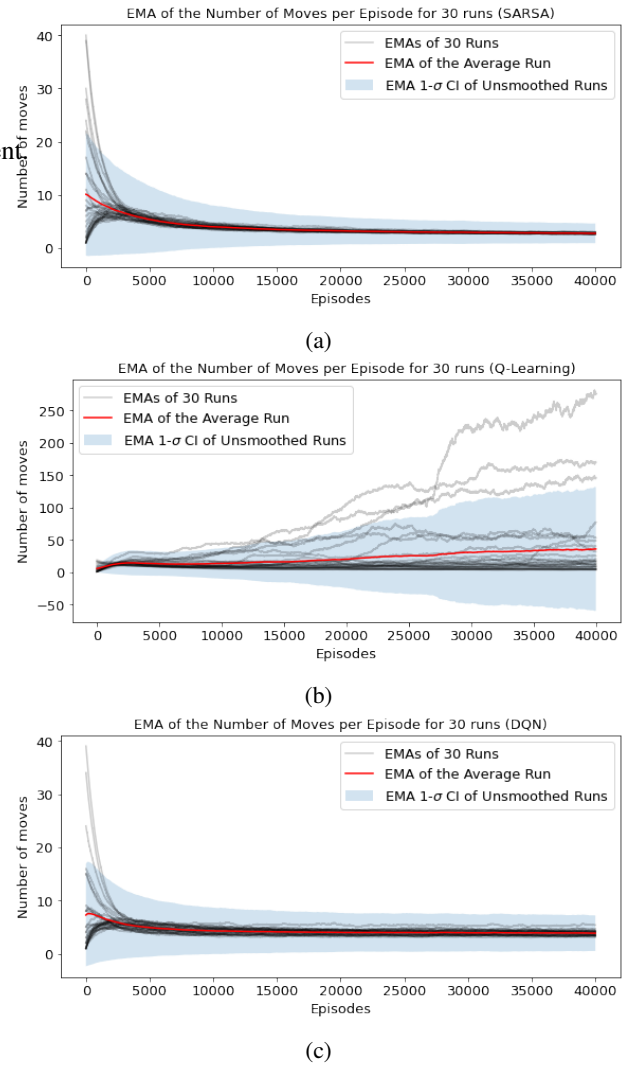


Fig. 5: Expected moving averages (EMA) of the number of moves per episode 30 runs for each algorithm (SARSA, Q-Learning, DQN). The red line depicts the EMA of the average across the 30 runs. The blue area depicts the EMA of the 1 standard deviation (σ) confidence interval (CI) of the 30 runs.

```

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
17
18 def EpsilonGreedy_Policy(Qvalues, allowed_a, epsilon
19 ):
```

```

19 """
20 returns: tuple
21     an action in form of a one-hot encoded
22     vector with the same shape/ensions as
23     Qvalues.
24     an action as decimal integer (0-based)
25 Assumes only a single state, i.e. online
26 learning and NOT (mini-)batch learning.
27 """
28 # get the Qvalues and the indices (relative of
29 all Qvalues) for the allowed actions
30 allowed_a_ind = np.where(allowed_a==1)[0]
31 Qvalues_allowed = Qvalues[allowed_a_ind]
32
33 # ----- epsilon greedy -----
34 # draw a random number and compare it to epsilon
35 rand_value = np.random.uniform(0, 1, 1)
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)
44
45 # get index of the action that was chosen (
46 relative to all actions, not only allowed)
47 ind_of_action_taken = allowed_a_ind[
48     action_taken_ind_of_allwed_only]
49
50 # ----- create usable output -----
51
52 # get the shape/ensions of the Qvalues
53 N_a, N_samples = np.shape(Qvalues) # N_samples
54 must be 1
55
56 # initialize all actions of binary mask to 0
57 A_binary_mask = np.zeros((N_a, N_samples))
58 # set the action that was chosen to 1
59 A_binary_mask[ind_of_action_taken, :] = 1
60
61 return A_binary_mask, ind_of_action_taken
62
63 # ===== activation functions and it's derivatives
64 =====
65
66 # relu and its derivative
67 def relu(x):
68     return np.maximum(0, x)
69
70 def heaviside(x):
71     return np.heaviside(x, 0)
72
73 # sigmoid and its derivative
74 def sigmoid(x):
75     return 1 / (1 + np.exp(-x))
76
77 def gradient_sigmoid(x):
78     return sigmoid(x) * (1 - sigmoid(x))
79
80 # tanh and its derivative
81 def tanh(x):
82     return np.tanh(x)
83
84 def gradient_tanh(x):
85     return 1 - np.tanh(x)**2
86
87
88 # identity and its derivative
89 def identity(x):
90     return x
91
92 def const(x):
93     return np.ones(x.shape)
94
95 def act_f_and_gradient(activation_function="relu"):
96     if activation_function == "relu":
97         return relu, heaviside
98     elif activation_function == "sigmoid":
99         return sigmoid, gradient_sigmoid
100     elif activation_function == "tanh":
101         return tanh, gradient_tanh
102     else: # identity and constant 1
103         return identity, const
104
105 # ===== Replay Memory for Experience Replay (with
106 DQN) =====
107
108 Transition = namedtuple('Transition', ("state", "
109 action", "reward", "next_state", "done"))
110
111 class ReplayMemory(object):
112     def __init__(self, capacity):
113         self.memory = deque(maxlen=capacity)
114
115     def push(self, *args):
116         self.memory.append(Transition(*args))
117
118     def sample(self, batch_size):
119         # if less data than batch size, return all
120         data
121         if len(self) < batch_size:
122             batch_size = len(self)
123         return random.sample(self.memory, batch_size)
124
125     def __len__(self):
126         return len(self.memory)
127
128 # ===== Neural Network =====
129
130 class NeuralNetwork(object):
131     def __init__(self, N_in, N_h, N_a,
132         activation_function_1="relu",
133         activation_function_2=None, method="
134 qlearning", seed=None, capacity=100_000, C
135 =100):
136         """
137         activation functions: "relu", "sigmoid", "
138         tanh", None
139         methods: "qlearning", "sarsa", "dqn"
140         """
141         self.D = N_in # input dimension (without
142         bias)
143         self.K = N_h # nr hidden neurons (without
144         bias)
145         self.O = N_a # nr output neurons (letter O
146         , not digit 0)
147
148         # store method and seed
149         self.method = method
150         self.seed = seed
151
152         if self.method == "dqn":

```



```

143         self.capacity = capacity
144         self.replay_memory = ReplayMemory(
145             capacity)
146         self.C = C
147     # set activation function and gradient
148     # function
149     self.act_f_1_name = activation_function_1
150     self.act_f_2_name = activation_function_2
151     self.act_f_1, self.grad_act_f_1 =
152         act_f_and_gradient(activation_function_1)
153     self.act_f_2, self.grad_act_f_2 =
154         act_f_and_gradient(activation_function_2)
155     # initialize the weights and biases and set
156     # global seed
157     np.random.seed(self.seed)
158     # self.W1 = np.random.randn(self.K+1, self.D
159     # +1)/np.sqrt(self.D+1) # standard normal
160     # distribution, shape: (K+1, D+1)
161     # glorot/xavier normal initialization
162     # self.W1 = np.random.randn(self.K+1, self.D
163     # +1)*np.sqrt(2/ (self.D+1 + self.K+1)) #
164     # standard normal distribution, shape: (K
165     # +1, D+1)
166     self.W1 = np.random.standard_normal((self.K
167     +1, self.D+1))*np.sqrt(2/ (self.D+1 +
168     self.K+1)) # standard normal
169     # distribution, shape: (K+1, D+1)
170     # self.W1 = np.random.randn(self.K+1, self.D
171     # +1) # standard normal distribution,
172     # shape: (K+1, D+1)
173     # self.W2 = np.random.randn(self.O, self.K
174     # +1)/np.sqrt(self.K+1) # standard normal
175     # distribution, shape: (O, K+1)
176     # glorot/xavier normal initialization
177     self.W2 = np.random.standard_normal((self.O
178     +1, self.K+1))*np.sqrt(2/ (self.K+1 +
179     self.O)) # standard normal distribution,
180     # shape: (O, K+1)
181     # self.W2 = np.random.randn(self.O, self.K
182     # +1) # standard normal distribution,
183     # shape: (O, K+1)
184     if self.method == "dqn":
185         self.W1_target = np.copy(self.W1)
186         self.W2_target = np.copy(self.W2)
187     def forward(self, x, target=False):
188         """
189         x has shape: (D+1, 1) (constant bias 1 must
190         be added beforehand)
191         target: if True, use the weights of the
192         target network
193         returns:
194             last logits (i.e. Qvalues) of shape (O
195             +1)
196         """
197         if target == True:
198             W1 = np.copy(self.W1_target)
199             W2 = np.copy(self.W2_target)
200         else:
201             W1 = np.copy(self.W1)
202             W2 = np.copy(self.W2)
203     # forward pass/propagation
204
205     a1 = W1 @ x
206     h1 = self.act_f_1(a1)
207     h1[0,:] = 1 # set first row (bias to second
208     layer) to 1 (this ignores the weights
209     for the k+1th hidden neuron, because
210     this should not exist; this allows to
211     only use matrix multiplication and
212     simplify the gradients as we only need 2
213     instead of 4)
214     a2 = W2 @ h1
215     h2 = self.act_f_2(a2)
216     return a1, h1, a2, h2
217
218 def backward(self, R, x, Qvalues, Q_prime, a1,
219             h1, a2, gamma, future_reward,
220             action_binary_mask):
221     """
222     backward for methods "qlearning" and "sarsa"
223
224     x has shape (D+1, 1) (constant bias 1 must
225     be added beforehand)
226     set future_reward=True for future reward
227     with gamma>0, False for immediate reward
228     .
229     Q_prime must be chosen according to the
230     method on x_prime (on- or off-policy)
231     """
232     # backward pass/backpropagation
233     # compute the gradient of the square loss
234     # with respect to the parameters
235
236     # ===== compute TD error (aka delta) =====
237
238     # make reward of shape (O, 1)
239     R_rep = np.tile(R, (self.O, 1))
240     if future_reward: # future reward
241         delta = R_rep + gamma*Q_prime - Qvalues
242         # -> shape (O, 1)
243     else: # immediate reward
244         delta = R_rep - Qvalues # -> shape (O,
245         1)
246
247     # update only action that was taken, i.e.
248     # all rows apart from the one
249     # corresponding to the action taken (
250     # action index) are 0
251     delta = delta*action_binary_mask
252
253     self.compute_gradients(delta, a1, h1, a2, x)
254     self.update_parameters(self.eta)
255
256 def backward_dqn(self, batch, gamma):
257     """
258     backward for method "dqn"
259
260     # ===== compute targets y and feature matrix
261     # X =====
262
263     # turn batch into individual tuples, numpy
264     # arrays, or lists
265     states = batch.state
266     rewards = np.array(list(batch.reward))
267     actions = np.array(list(batch.action))
268     next_states = list(batch.next_state)
269     dones = np.array(list(batch.done))
270
271     # compute targets y and feature matrix X
272     y = np.zeros((self.O, len(dones)))
273     for j in np.arange(len(dones)):

```

```

244         if dones[j]: # if done, set y_j = r_j 305
245             y[actions[j], j] = rewards[j] 306
246         else: 307
247             # compute Q_prime 308
248             Q_target = self.forward(next_states 309
249                                     j], target=True)[-1] 310
249             y[actions[j], j] = rewards[j] + 311
250                 gamma*np.max(Q_target) 312
251 313
252 # convert states to feature matrix X 315
253 X = np.hstack((states)) 316
254 317
255 318
256 # ===== compute TD error (aka delta) ===== 319
257
258 a1, h1, a2, Qvalues = self.forward(X) 320
259 delta = y - Qvalues # -> shape (O, 321
260         batch_size) 322
261
262 self.compute_gradients(delta, a1, h1, a2, X) 323
263 self.update_parameters(self.eta) 324
264
265 def compute_gradients(self, delta, a1, h1, a2, X): 325
266     # ===== compute gradient of the loss with 326
267     # respect to the weights ===== 327
268     # common part of the gradient TODO: check 328
269     # dimensions 329
270     self.dL_da2 = delta * self.grad_act_f_2(a2) 330
271     # gradient of loss wrt W2 331
272     self.dL_dW2 = self.dL_da2 @ h1.T 332
273     # gradient of loss wrt W1 333
274     self.dL_dW1 = ( (self.W2.T @ self.dL_da2) * 334
275                     self.grad_act_f_1(a1) ) @ x.T 335
276
277
278
279 def update_parameters(self, eta): 336
280     # gradient clipping 337
281
282     # dL_dW1_norm = np.linalg.norm(self.dL_dW1) 338
283     # if dL_dW1_norm >= self.gradient_clip: 339
284     #     self.dL_dW1 = self.gradient_clip * 340
285     #         self.dL_dW1 / dL_dW1_norm 341
286
287     # dL_dW2_norm = np.linalg.norm(self.dL_dW2) 342
288     # if dL_dW2_norm >= self.gradient_clip: 343
289     #     self.dL_dW2 = self.gradient_clip * 344
290     #         self.dL_dW2 / dL_dW2_norm 345
291
292 # update W1 and W2 346
293 self.W2 = self.W2 + eta * self.dL_dW2 347
294 self.W1 = self.W1 + eta * self.dL_dW1 348
295
296
297
298 def train(self, env, N_episodes, eta, epsilon_0, 349
299         beta, gamma, alpha=0.001, gradient_clip=1, 350
300         batch_size=32, run_number=None): 351
301     """ 352
302     alpha is used as weight for the exponential 353
303     moving average displayed during training 354
304     """ 355
305     batch_size is only used for the DQN method. 356
306     """ 357
307     # add training hyper parameters 358

```

```

self.N_episodes = N_episodes
self.eta = eta
self.epsilon_0 = epsilon_0
self.beta = beta
self.gamma = gamma
self.alpha = alpha
self.gradient_clip = gradient_clip
self.batch_size = batch_size

training_start = time.time()

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

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

    n_steps = 0

    for n in episodes:
        epsilon_f = self.epsilon_0 / (1 +
            beta * n) ## DECAYING EPSILON
        Done = 0

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

        ## COUNTER FOR NUMBER OF ACTIONS

        S, X, allowed_a = env.
            Initialise_game() ##
        INITIALISE GAME
        X = np.expand_dims(X, axis=1)
        ## MAKE X A
        TWO DIMENSIONAL ARRAY
        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"):

```



```

352         # compute Q values for the
353         given state 389
354         a1, h1, a2, Qvalues = self.
355         forward(X) # -> shape 390
356         O, 1)
357
358         # choose an action A using 391
359         epsilon-greedy policy 392
360         A_binary_mask, A_ind = 393
361         EpsilonGreedy_Policy(
362             Qvalues, allowed_a,
363             epsilon_f) # -> shape 394
364         O, 1)
365
366         # take action and observe reward 395
367         R and state S_prime 396
368         S_prime, X_prime,
369         allowed_a_prime, R, Done = 397
370         env.OneStep(A_ind)
371         X_prime = np.expand_dims(X_prime, 398
372             , axis=1) 399
373         X_prime = np.copy(np.vstack((np. 400
374             array([[1]]), X_prime))) 401
375         add bias term 402
376
377         n_steps += 1
378
379         if self.method == "dqn": 403
380             404
381             # store the transition in
382             memory
383             self.replay_memory.push(X, 405
384                 A_ind, R, X_prime, Done) 406
385             407
386             # sample a batch of
387             transitions 408
388             transactions = self. 409
389             replay_memory.sample( 410
390                 self.batch_size)
391             # turn list of transactions 411
392             into transaction of
393             lists 412
394             batch = Transition(*zip(*
395                 transactions))
396             413
397             # backward step and
398             parameter update 414
399             self.backward_dqn(batch, 415
400                 self.gamma) 416
401             417
402             # update Q values indirectly by
403             updating the weights and
404             biases directly
405
406             if Done==1: # THE EPISODE HAS 418
407             ENDED, UPDATE...BE CAREFUL 419
408             THIS IS THE LAST STEP OF THE 420
409             EPISODE
410
411             if (self.method == "
412                 qlearning") or (self.
413                 method == "sarsa"):
414                 # compute gradients and
415                 update weights 421
416                 self.backward(R, X, 422
417                     Qvalues, None, a1,
418                     h1, a2, None, 423
419                     future_reward=False,
420                     action_binary_mask= 424
421                     A_binary_mask) 425
422
423             # store history
424             # todo: record max possible
425
426         reward per episode
427         self.R_history[n] = np.copy(
428             R) # reward per episode
429         self.N_moves_history[n] = np
430             .copy(i) # nr moves per
431             episode
432
433         # store norm of gradients
434         self.dL_dW1_norm_history[n]
435             = np.linalg.norm(self.
436                 dL_dW1)
437         self.dL_dW2_norm_history[n]
438             = np.linalg.norm(self.
439                 dL_dW2)
440
441         # compute exponential moving
442         average (EMA) to
443         display during training
444         ema = alpha*R + (1-alpha)*
445             ema_previous
446         if n == 0: # first episode
447             ema = R
448         ema_previous = ema
449         if run_number is not None:
450             episodes.set_description
451                 (f"Run = {run_number
452                     }; EMA Reward = {ema
453                         :.2f}")
454         else:
455             episodes.set_description
456                 (f"EMA Reward = {ema
457                     :.2f}")
458
459         break
460
461     else: # IF THE EPISODE IS NOT
462         OVER...
463
464         if self.method == "qlearning
465             ":
466             # chose next action off-
467             policy
468             Q_prime = np.max(self.
469                 forward(X_prime)
470                 [-1])
471
472         elif self.method == "sarsa":
473             # chose next action on-
474             policy
475
476             a1_prime, h1_prime,
477             a2_prime,
478             Qvalues_prime = self
479                 .forward(X_prime) #
480                 -> shape (N_a, 1)
481
482             # chose next action and
483             save it
484             A_binary_mask_prime,
485             A_ind_prime =
486                 EpsilonGreedy_Policy
487                 (Qvalues_prime,
488                 allowed_a_prime,
489                 epsilon_f)
490
491             # get Qvalue of next
492             action
493             Q_prime = Qvalues_prime[
494                 A_ind_prime]
495
496         if (self.method == "
497             qlearning") or (self.
498             method == "sarsa"):

```

```

427         # backpropagation and weight update
428         self.backward(R, X, Qvalues, Q_prime, a1, h1, a2, self.gamma, future_reward=True, action_binary_mask=A_binary_mask)
429
430     # NEXT STATE AND CO. BECOME ACTUAL STATE...
431     if self.method == "sarsa":
432         A_binary_mask = np.copy(A_binary_mask_prime)
433         A_ind = np.copy(A_ind_prime)
434         a1 = np.copy(a1_prime)
435         h1 = np.copy(h1_prime)
436         a2 = np.copy(a2_prime)
437         Qvalues = np.copy(Qvalues_prime)
438         S = np.copy(S_prime)
439         X = np.copy(X_prime)
440         allowed_a = np.copy(allowed_a_prime)
441
442     i += 1 # UPDATE COUNTER FOR NUMBER OF ACTIONS
443
444     if (self.method == "dqn") and n_steps % self.C == 0:
445         # update target network every C steps
446         self.W1_target = np.copy(self.W1)
447         self.W2_target = np.copy(self.W2)
448
449     training_end = time.time()
450     self.training_time_in_seconds = training_end - training_start
451
452     return None
453
454 except KeyboardInterrupt as e:
455     # return nothing
456     training_end = time.time()
457     self.training_time_in_seconds = training_end - training_start
458
459     return None
460
461 def save(self, name_extension=None):
462     # create directory for the model
463     name = f"{self.method}_{self.act_f_1_name}_{self.act_f_2_name}"
464     if name_extension is not None:
465         name += f"_{name_extension}"
466
467     path = f"models/{name}"
468     if not os.path.isdir(path): os.mkdir(path)
469     print(f"saving to: {path}")
470
471     # save weights
472     np.save(f"{path}/W1.npy", self.W1)
473     np.save(f"{path}/W2.npy", self.W2)
474
475     # save training history
476     np.save(f"{path}/training_history_R.npy", self.R_history)
477
478     np.save(f"{path}/training_history_N_moves.npy", self.N_moves_history)
479     np.save(f"{path}/training_history_dL_dW1_norm.npy", self.dL_dW1_norm_history)
480     np.save(f"{path}/training_history_dL_dW2_norm.npy", self.dL_dW2_norm_history)
481
482     # save training parameters and other general info
483     params = {
484         "method": self.method,
485         "N_episodes": self.N_episodes,
486         "eta": self.eta,
487         "epsilon_0": self.epsilon_0,
488         "beta": self.beta,
489         "gamma": self.gamma,
490         "alpha": self.alpha,
491         # "gradient_clip": self.gradient_clip,
492         "seed": self.seed,
493         "D": self.D,
494         "K": self.K,
495         "O": self.O,
496         "training_time_in_seconds": self.training_time_in_seconds
497     }
498     if self.method == "dqn":
499         params["capacity"] = self.capacity
500         params["batch_size"] = self.batch_size
501         params["C"] = self.C
502     with open(f"{path}/training_parameters.json", "w") as f:
503         json.dump(params, f)
504
505 def load_from(method, act_f_1, act_f_2, name_extension=None):
506     # read values and store in neural network instance
507     name = f"{method}_{act_f_1}_{act_f_2}"
508     if name_extension is not None:
509         name += f"_{name_extension}"
510
511     path = f"models/{name}"
512     # print(f"loading from: {path}")
513
514     # initialize neural network
515     nn = NeuralNetwork(0,0,0, activation_function_1=act_f_1, activation_function_2=act_f_2, method=method)
516
517     # network weights
518     nn.W1 = np.load(f"{path}/W1.npy")
519     nn.W2 = np.load(f"{path}/W2.npy")
520
521     # network training history
522     nn.R_history = np.load(f"{path}/training_history_R.npy")
523     nn.N_moves_history = np.load(f"{path}/training_history_N_moves.npy")
524     nn.dL_dW1_norm_history = np.load(f"{path}/training_history_dL_dW1_norm.npy")
525     nn.dL_dW2_norm_history = np.load(f"{path}/training_history_dL_dW2_norm.npy")
526
527     # network training parameters
528     with open(f"{path}/training_parameters.json", "r") as f:
529         params = json.load(f)
530
531     # set parameters to the network instance
532     nn.method = params["method"]

```

```

538     nn.N_episodes = int(params["N_episodes"])
539     nn.eta = float(params["eta"])
540     nn.epsilon_0 = float(params["epsilon_0"])
541     nn.beta = float(params["beta"])
542     nn.gamma = float(params["gamma"])
543     nn.alpha = float(params["alpha"])
544     # nn.gradient_clip = float(params["
        gradient_clip"])
545     try:
546         nn.seed = int(params["seed"])
547     except:
548         nn.seed = params["seed"]
549     nn.D = int(params["D"])
550     nn.K = int(params["K"])
551     nn.O = int(params["O"])
552     nn.training_time_in_seconds = float(params["
        training_time_in_seconds"])
553
554     if nn.method == "dqn":
555         nn.capacity = int(params["capacity"])
556         nn.batch_size = int(params["batch_size"
            ])
557         nn.C = int(params["C"])
558
559
560     if nn.method == "dqn":
561         nn.W1_target = np.copy(nn.W1)
562         nn.W2_target = np.copy(nn.W2)
563
564     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.