# Reinforcement learning: Project 4

Florence Hugh, August Jonasson, Amanda Tisell

March 21, 2025

## Introduction

In this project our goal is to train an agent to play the game Othello, which is a modern version of Reversi[1]. Othello is a turn-based strategy board game for two players played on an $8 \times 8$ grid. Players take turns placing one piece on an empty spot on the board with their assigned color. After a play is made, the pieces of the opponent's color that lie bounded by the current player's color are turned over. The game ends when the board is filled or when none of the players can make a move. The player who controls the most pieces at the end wins the game. Figure 1 shows an example of the first few moves in a game.
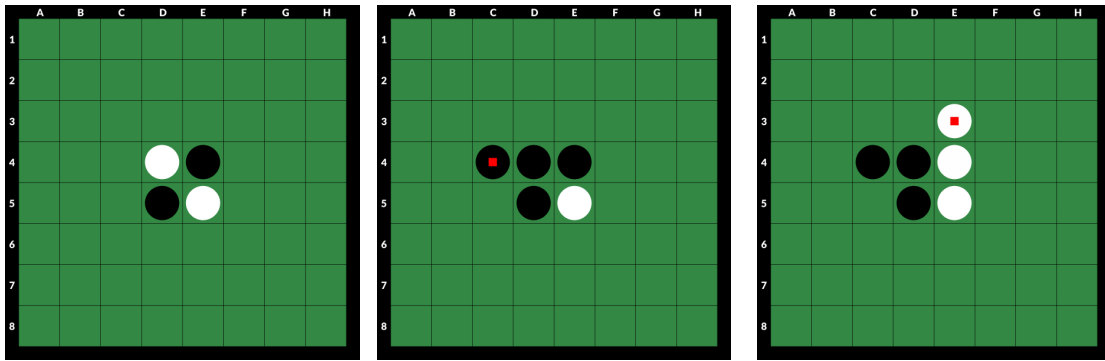


Figure 1: Starting state and first two moves of a game.

### Some challenges and how we address them

Since Othello is a two-player game, we need to consider the agent/opponent setting. A natural idea would be to have one set of values for black and a completely different set of values for white, thus yielding two separate policies. When dealing with a two-player setting such as this, however, it turns out that we can avoid having to train two agents simultaneously by using self-play in combination with a minimax algorithm, thus allowing for both the black and white pieces to play against the same set of values.

Othello, similarly to Chess and AlphaGo, is unsolved and has a huge state space (around $10^{28}$ possible board positions on an $8 \times 8$ grid). Handling this size would require a large amount of compute which we do not have. Therefore, we decide to cut down the size of the board to a $6 \times 6$ grid instead. This results in a state space of around $10^{12}$, which is still large enough to prevent accidentally brute-forcing the problem, but much more manageable within the scope of this project. A possibly even more important

---

[1]https://en.wikipedia.org/wiki/Reversi

advantage of this simplification is that $6 \times 6$ Othello is solved, with the white pieces having a clear advantage[2]. This means that it will now be much easier to assess the results of the training, since we now know what behavior to look for. Ideally, we would like to see the agent converge towards white winning every game.

As in most reinforcement learning systems, we have a setting of sparse rewards (only at the end of games) and the issue of balancing exploration/exploitation properly. We want to be able to learn meaningful patterns in a reasonable amount of time without having to go through the whole decision tree of possible games. Our attempt at addressing both of these problems (which is heavily inspired by the AlphaGo article Silver et al. 2016) is to use a Monte Carlo tree search algorithm in combination with a neural network-approximated policy that, simply put, favors deep over wide exploration such that rewards become more frequent, thus speeding up the learning.

## MDP framing

The environment is a $6 \times 6$ grid where the dynamics are completely deterministic. The number of unique board positions in $6 \times 6$ Othello has been approximated to $10^{12}$ ($3^{32} \cdot 2^4 \approx 10^{16}$ and then divided by some number that accounts for what proportion of these actually represent legal board positions). A state will be represented as a $6 \times 6$ matrix with entries being

$$s_{ij} = \begin{cases} -1 & \text{if square } (i,j) \text{ is occupied by a white piece,} \\ 0 & \text{if square } (i,j) \text{ is not occupied by any piece,} \\ 1 & \text{if square } (i,j) \text{ is occupied by a black piece.} \end{cases} \tag{1}$$

The action space of size 32 includes all squares on the board, excluding the four starting pieces. However, most of these actions will be illegal. According to a paper that claims to have weakly solved Othello, there are on average 10 legal actions per move on a $8 \times 8$ grid (Takizawa 2023). This leads us to believe that a $6 \times 6$ grid may have on average 5-8 legal actions from a given state. An action $a$ is represented as a tuple $(i,j)$, which is the position where the piece is placed on the board.

A reward is only given at the end of each game and will be defined as

$$R(s,a) = \begin{cases} 1 & \text{if black wins the game} \\ -1 & \text{if white wins the game} \\ 0 & \text{if the game is drawn} \end{cases} \tag{2}$$

---

[2]https://en.wikipedia.org/wiki/Computer_Othello

# Monte Carlo tree search (MCTS)

Monte Carlo tree search (Sutton, Barto, et al. 2018, ch. 8.11) is a type of decision-time planning algorithm based on Monte Carlo control applied to simulations starting from the root state, that is, a type of rollout algorithm.

Decision-time planning (Sutton, Barto, et al. 2018, ch. 8.8) is a type of planning that focuses on a particular state. That is, instead of simulating experience for the whole state space, it only plans ahead from the current state $s_t$. Using simulated experience to gradually improve a policy or value function or to select an action for the current state. The values and policy are specified to the current state and the action choices available there. In decision-time planning there is a tradeoff between better policies and the time that is needed to simulate enough trajectories to obtain good value estimates.

A rollout algorithm (Sutton, Barto, et al. 2018, ch. 8.10) is a decision-time planning algorithm based on Monte Carlo control applied to simulated trajectories that all begin in the current state. Each rollout produce Monte Carlo estimates of action-values for each current state and for a given policy, usually called the rollout policy. The aim of a rollout algorithm is to improve the rollout policy and not find the optimal policy like in regular Monte Carlo methods.

In MCTS, as in a rollout algorithm, each execution is an iterative process that simulates many trajectories, each trajectory starting from the current state $s_t$, running to a terminal state $s_T$. This iterative process consists of four steps; selection, expansion, rollout (simulation) and backpropagation.
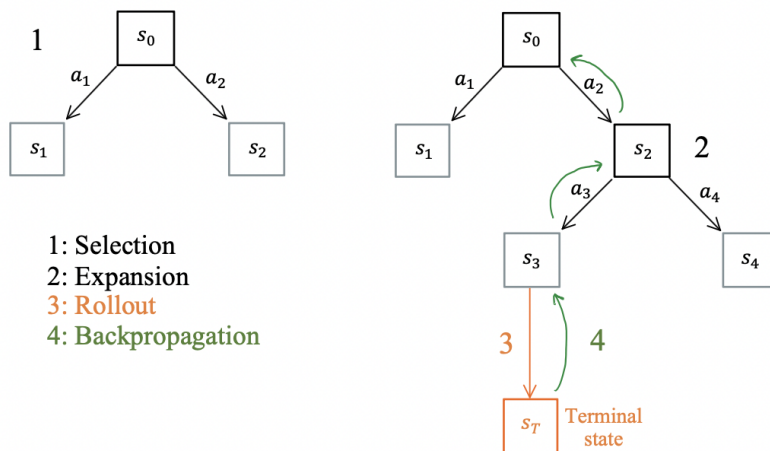


Figure 2: Monte Carlo tree search

These four steps continues to execute, starting each time at the tree's root node (initial state), until some criterion is fulfilled. This criterion can for instance be iterating for a fixed number of iterations (fixed number of terminal states reached with the tree policy).

## 1. Selection

In the selection part, a node (state) is chosen according to a tree policy. The tree policy balances exploration and exploitation, where it's action selection can for instance be based on $\epsilon$-greedy or an

upper confidence bound (UCB) selection formula

$$a_t = \arg\max_a \left( Q(s,a) + c\sqrt{\frac{\log N(s)}{N(s,a)}} \right). \tag{3}$$

The $Q(s,a)$ is an average action-value (updated by the backpropagation, explained in part 4), $c$ is a constant that controls the greediness of the action selection, $N(s)$ is the number of visits in current state $s$ and $N(s,a)$ is the number of visits in the next state given action $a$.

The MCTS traverses down the decision tree and whenever a state is visited, the next action is selected among the possible actions from this current state according to equation 3. This process continues until a leaf node (unexplored and not terminal state) is reached in which an expansion is applied.

## 2. Expansion

The decision tree is expanded from the current state by connecting all possible actions to their respectively next state. From these new states one is selected based on the tree policy used in the selection step. However, all these newly added states will have a $N(s,a)$ count equal to zero. Since the MCTS priorities untried actions, one of these states will simply be picked randomly as a zero count will be seen as a $a_t$ value of $\pm\infty$. Whenever a zero counted state-action is chosen we proceed to the next step, the rollout.

## 3. Rollout

During the rollout, a simulation is applied from the current state. This simulates a whole episode starting from the current state until a terminal state is reached. The actions selected during the simulation are generated using a rollout policy and does not necessarily follow the chosen tree policy. The rollout policy is usually a simpler policy allowing for faster computation. It could for example follow a policy where actions are chosen uniformly. When the terminal state is reached, a rollout reward $R_r$ is received. The reward will be backpropagated and finally, all values generated by the rollout will be discarded.

## 4. Backpropagation

In the final phase, called the backpropagation, the reward observed from doing the rollout is backpropagated through the tree, updating all state-action values from where we started the rollout (the leaf) up to the root. The update rule is the in-place update of the average, given by

$$Q_{n+1}(s,a) \leftarrow Q_n(s,a) + \frac{1}{N(s,a)} \left[ R_r - Q_n(s,a) \right], \tag{4}$$

(Sutton, Barto, et al. 2018, eq. (2.5)) where $Q_n(s,a)$ is the old state-action-value, $N(s,a)$ is the number of visits to state-action pair $s,a$, $R_r$ is the reward yielded from the rollout, and $Q_{n+1}(s,a)$ is the updated value.

# Modifying MCTS

In order to speed up the tree traversal and gain better control of the exploration (such that rewards become more frequent) we now make some modifications to regular MCTS.

First, to address the agent/opponent setting of the system, we solve this problem by using a minimax[3] algorithm in combination with self-play. As such, we reframe the system as an agent learning how to play Othello against itself, rather than an agent playing Othello against an opponent. The minimax algorithm is a decision rule strategy used in game theory in order to find optimal moves in zero-sum games (when a player's gain is balanced out by the other player's loss). The algorithm assumes that one player tries to maximize its score, while the opponent tries to minimize it. This fits well into our model as our reward structure results in a zero-sum game.

We have incorporated this algorithm into the Monte Carlo tree structure by adding additional information about which color's move it is at each state. By doing this the MCTS has an indicator of when to minimize/maximize the next state-action value. The action selection part of the tree will handle the minimax situation by using formula

$$a_t = \underset{a}{\overset{(\arg\min)}{\arg\max}} \left( Q(s,a) \underset{(-)}{\overset{+}{}} c \frac{P(s,a)\sqrt{N(s)}}{N(s,a)} \right), \tag{5}$$

where the black (white) player will try to maximize (minimize) the return.

In addition, a policy $P(s,a)$ is added in order to make the agent more selective toward states that have already seen a lot of visits, thus pushing the actual exploration down along a few branches instead of higher up among many branches. This will lead to rewards being less sparse in the training, and as such, will speed up learning. The policy itself is retrieved through a neural network, which we will present in its own section.

Lastly, the rollout policy is replaced by function approximation of state-action values, as constantly letting the games play out to their end is too slow. Now, instead of playing them out, we simply input the leaf state into a function and get an approximation of the value back immediately. This function approximation is also handled by a neural network which we will present in its entirety in the next section. The backpropagation still remains as before. The only difference now is that instead of the actual reward, it is the approximated value that gets backpropagated, unless we have actually reached a terminal state - then the real reward is backpropagated.
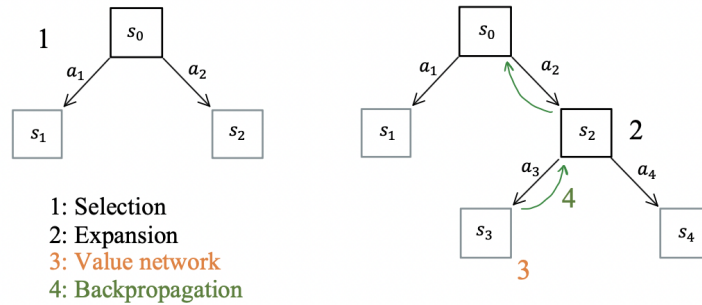


1: Selection
2: Expansion
3: Value network
4: Backpropagation

Figure 3: Modified Monte Carlo tree search

---

[3]https://en.wikipedia.org/wiki/Minimax#Minimax

# Value approximation network

The function approximation of the action-values is handled by a deep convolutional neural network. In this section we will briefly present the architecture of the network and also explain how, and on what data, the initial training was done. For specific hyperparameter settings (such as batch-sizes, learning rates, etc.) we refer to the code appendix (line 749 in the code).

The architecture of the value neural network is heavily inspired by resnet-18 (He et al. 2016). Figure 4 depicts the network in its entirety.
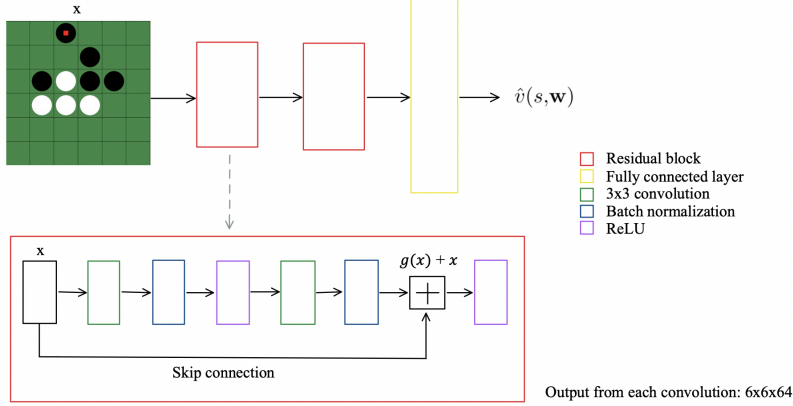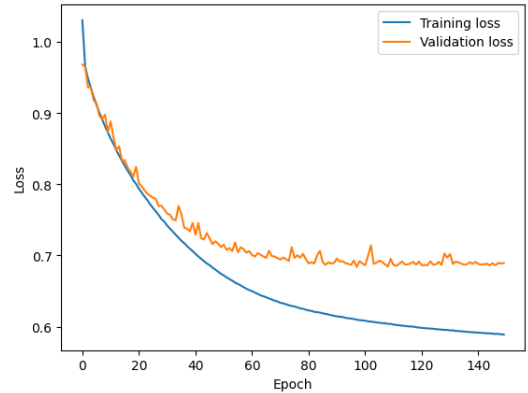


Figure 4: Value network architecture.

The input is a $6 \times 6$ matrix with elements according to (1) in the MDP framing, the targets are the Monte Carlo returns associated with the game from which that particular board position was retrieved, i.e. $-1$, $0$ or $1$, the output is the full-gradient action-value approximation, and finally, the loss function is the Mean Squared Value Error with $\mu(s)$ set to constant $1/|\text{data}|$ (Sutton, Barto, et al. 2018, p. 199).

Figure 5: Initial training of the value network.



The data used for the initial training of the value network consists of unique board positions retrieved from one hundred thousand randomly simulated games, with each one labeled according to the reward that was given at the end of its respective game. In total, this yields a dataset of size $\approx 2 \cdot 10^6$, which is then divided into a training dataset and a validation dataset according to a 80-20 split. Figure 5 depicts the training curve for 150 epochs. From the validation curve, we see no indication of overfitting, and the training curve seems to converge. As such, we deduce that the training is healthy. Looking at the actual loss, however, which is squared, we notice that for returns in the range $[-1, 1]$ it seems relatively large. This might be an indication that our network is not deep enough to successfully capture the complex relationships between board positions and returns. Given more time and resources, adding more residual blocks is a potential improvement for the project. See code appendix for the full `tensorflow` implementation.
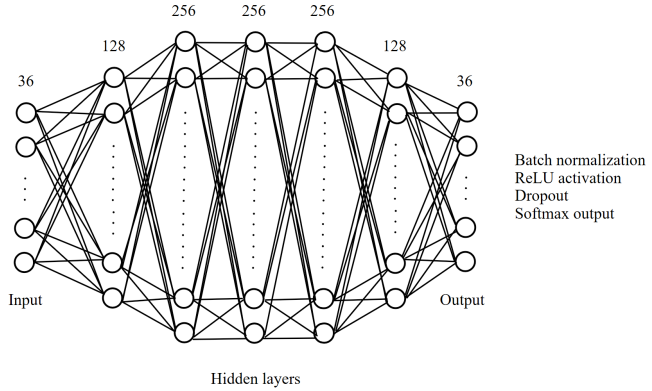
# Policy approximation network

The policy approximation is handled by a plain, fully connected, feedforward neural network with linear inputs up until and including the output layer, which uses a softmax function in order to model the probability distribution over the legal actions in the given state. The architecture is summarized in Figure 6. Again, refer to the code appendix for specific hyperparameter settings (line 802 in the code).
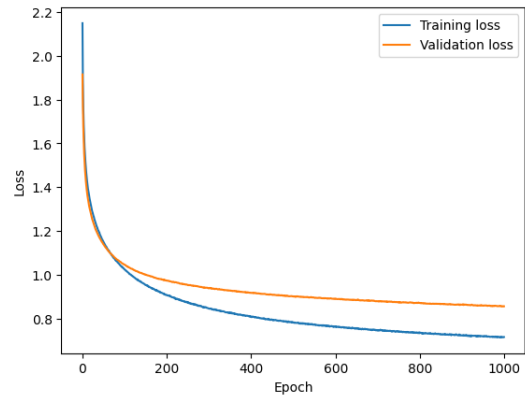
Figure 6: Initial training of the value network.

Since this network only requires the states as a sort of identifier to know which actions are legal, the board position matrices can be flattened to a 36-dimensional vector before being input into the network, thus avoiding the need for convolutions to preserve more complex relations. Most likely, the states can be mapped to an even lower dimension and still act as unique identifiers. This is a potential improvement for the network, in terms of lowering the necessary complexity. The targets here are the normalized number of times each action had been taken from the given state in that particular tree, thus forming an approximate probability distribution over the actions. Since the number of actions can vary depending on the state, the output dimension is kept at a constant 36. This means that the target will be of this dimension as well, but with zeros on all positions that are not legal actions from the given state.

This means that the predicted probabilities will not form a perfect distribution over the legal actions, but the hope is that the bulk of the probability mass will still fall here, as the remaining actions are fitted to zeros. Dealing with this is another potential improvement. Finally, the Kullback-Leibler divergence[4] was used as loss function (minimization of distance between probability distributions).

Figure 7: Initial training of the value network.

The dataset used for the initial training of the policy network is a set of one thousand games simulated using the regular *upper-confidence bound* formula (eq. (3)) but with the value network instead of the rollout policy, thus forming a sort of intermediary step before transitioning completely to the model in Figure 3. Ideally, this is where we would want to use prerecorded "expert" experience, but due to not being able find any such databases for $6 \times 6$ Othello, we settle for this solution, and assume that the MCTS with the value network alone is sufficient. The training can be seen in Figure 7 and looks healthy. Since the targets represent unique state-action visits, we are ensured that there are no identical matches between the training and validation set. This indicates that the model generalizes well to unseen data.

---

[4]https://en.wikipedia.org/wiki/Kullback%E2%80%93Leibler_divergence

# Training cycle

With the modified MCTS, the value network and the policy network, we are now ready to define a *training cycle*, which is shown in Figure 8. Using the initially trained neural networks, we begin with simulating games from the Monte Carlo tree with a predefined number of episodes before termination. These games are then stored in a replay buffer which upon termination of the tree is sent into training of the neural networks for a number of epochs. Finally the updated networks are sent back to the tree, where a new step of the cycle can begin.



Figure 8: The agent training cycle.

In the actual implementation of the model, we let one tree iterate until it had reached terminal states 25 times, for 30 cycles. Even a single cycle often took as long as 30 minutes to complete due to the tree exploring along several different branches before reaching terminal states. In total, with the initial training of the value and policy networks included, approximately 15 hours was spent on training the model (on a CPU).



Figure 9: Value network training from ten selected cycles.

Figures 9 and 10 depict the training curves of the neural networks within the 30 training cycles. Overall the training looks very healthy in both networks. Cycle 9 within the value network has a validation

Figure 10: Policy network training from ten selected cycles.

curve that indicates a poor fit, which would be interesting to examine further. Cycle 30 from the value network has a validation curve that lies entirely below the training curve. This is actually because the training has converged such that the training and validation datasets are perfectly correlated. The regularization in the training is what causes the validation curve to fall below. In the policy network training we don't see any anomalies, and we can now move on to the results of the agent.

# Results

Figure 11 shows the result after training the model on 30 training cycles with 25 episodes each. As seen in Figure 11, the white piece player becomes increasingly more dominant and, at the end, completely dominating. This is a promising result as $6 \times 6$ Othello is solved with the result that white has an advantage.



Figure 11: Reward distribution of each training cycle.



Figure 12: Terminal state distribution of cycle 5 (left) and cycle 30 (right).

Delving into some of the specific cycles now, we can see from Figure 12, for example, that even though white wins the first 6 episodes in the cycle, the fact that it wins in 6 different ways indicates that black is fighting back. Eventually, black even manages to get back with 6 wins in the same terminal state. Why white allows for this same terminal state to be reached so many times in a row could have something to do with the fact that white previously has gotten high rewards along this branch and therefore requires a couple of iterations before it realizes that better options exist.

Finally, though, which can be seen from Figure 12, white takes over completely and manages to force the game towards a single terminal state each game in the whole cycle. This happened in four other cycles towards the end as well, further strengthening the hypothesis that our training converged.

10

# Discussion

The final result indicates that our agent has trained well. However, to gain more insight into the agents behavior, a deeper analysis of the trained decision tree is required. This involves examining the structure of the tree to find meaningful patterns, such as frequently visited paths and areas where the agent shows uncertainty or bias, and map out specific trajectories and compare them to solutions of $6 \times 6$ Othello. Since the game is solved, it would actually be possible to see if an optimal policy has been reached.

We mention that initially we worked on a $4 \times 4$ grid where a large amount of time was spent on implementing the methods. This turned out to be too simple as we ended up brute-forcing the whole tree, resulting in an optimal solution during the first cycle. Then we jumped straight to the $8 \times 8$ grid in which we realized we did not have enough computational power. Since our code was written for the general case, we seemlessly ended up with the $6 \times 6$ grid eventually.

# Possible Improvements

## Neural Networks

The policy network needs further improvements seeing how the probability mass of the output did not match well with the legal actions deeper down the tree. Meaning that in some cases, as little as 50 % of the probability mass was distributed to the legal actions. One way to solve it could be by adding a masking part which zeros out the illegal actions and weighted up the legal actions so they sum to one. There is probably a better way to improve this but we would need more time to figure it out.

Another improvement could be to simply train one neural network instead of having two separate ones. This network would then have as output both the actions and the values. By only training one network, we may speed up the training cycle slightly. It was, however, the tree traversal that ate up most of the training time, due to the size of the state space. So, any gains in the traversal would be even more valuable.

## Parallel Processing

To train the tree, it needed to be traversed many times. Instead of just using one agent to traverse down the tree, we could have trained several agents simultaneously to enhance the process. There should also be a way for parallelizing the networks and the tree.

## Handle Symmetry

We know for sure that the first four positions in which black can put its pieces are completely symmetric. With this knowledge, we could have fixed the first black piece to one position and then only trained that part of the tree to optimize it performance on one branch. Then the result in this branch could somehow map the other starting positions to equivalent actions.

## Play against the agent

Ideally, it would be nice to be able to play with the agent. This way we could see if it actually wins all the games (given that it plays the white pieces) when playing against a human. We started building a graphical user interface (GUI) to play against the agent. However, because our agent had not explored enough state-action pairs, it could not make an action fast enough as it had to explore the tree from this unseen state. A way to solve this could be to either build a larger tree with the already trained weights from the networks (but this takes quite a time) or we could have added a time

horizon threshold such that the agent only had a set amount of time to think. As it was implemented now, it spent a long time thinking (planning ahead) on each move which made the games very long. We are unfortunately not entirely sure how well it performs against a human, as of now.

# References

He, Kaiming et al. (2016). "Deep residual learning for image recognition". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778.

Silver, David et al. (2016). "Mastering the game of Go with deep neural networks and tree search". In: *nature* 529.7587, pp. 484–489.

Sutton, Richard S, Andrew G Barto, et al. (2018). *Reinforcement learning: An introduction*. Vol. 1. 1. MIT press Cambridge.

Takizawa, Hiroki (2023). "Othello is solved". In: *arXiv preprint arXiv:2310.19387*.

# Appendix A: Python code

```python
1  import copy
2  import numpy as np
3  import pickle
4  import random
5  import tensorflow as tf
6  from MCT_Othello_classes import *
7  from sklearn.model_selection import train_test_split
8  from matplotlib import pyplot as plt
9  import matplotlib.ticker as mticker
10 import keras
11 from keras import layers, regularizers, Model
12 from othello_rl_helper_fcts import *
13
14 # 0=blank, 1=tot_black, -1=white
15 class OthelloBoard():
16     dirx = [-1, 0, 1, -1, 1, -1, 0, 1]
17     diry = [-1, -1, -1, 0, 0, 1, 1, 1]
18
19     def __init__(self, n):
20         self.n = n
21         self.board = [[0 for _ in range(n)] for _ in range(n)]
22         self.to_play = 1 #keep track of players turn 1=tot_black -1=white
23         # self.pass_counter = 0
24         self.reset_game()
25
26     def reset_game(self):
27         n = self.n
28         # self.pass_counter = 0
29         self.to_play = 1
30         self.board = [[0 for _ in range(n)] for _ in range(n)]
31         board = self.board
32
33         # set up initial bricks
34         z = (n - 2) // 2
35         board[z][z] = -1
36         board[n - 1 - z][z] = 1
37         board[z][n - 1 - z] = 1
38         board[n - 1 - z][n - 1 - z] = -1
39
40         return board
41
42     def print_board(self):
43         n = self.n
44         board = self.board
45         m = len(str(n - 1))
46         for y in range(n):
47             row = ''
48             for x in range(n):
49                 row += str(board[y][x])
50                 row += ' ' * m
51             print(row + ' ' + str(y))
52         print("")
53         row = ''
54         for x in range(n):
55             row += str(x).zfill(m) + ' '
56         print(row + '\n')
57
58     def make_move(self, curr_state, action, to_play):
59         x = action[0]
60         y = action[1]
61         if self.check_valid_move(curr_state, x, y, to_play):
62             n = self.n
63             bricks_taken = 0 # total number of opponent pieces taken
```

```
64
65              curr_state[y][x] = to_play
66              for d in range(len(self.dirx)): # 8 directions
67                  bricks = 0
68                  for i in range(n):
69                      dx = x + self.dirx[d] * (i + 1)
70                      dy = y + self.diry[d] * (i + 1)
71                      if dx < 0 or dx > n - 1 or dy < 0 or dy > n - 1:
72                          bricks = 0; break
73                      elif curr_state[dy][dx] == to_play:
74                          break
75                      elif curr_state[dy][dx] == 0:
76                          bricks = 0; break
77                      else:
78                          bricks += 1
79                  for i in range(bricks):
80                      dx = x + self.dirx[d] * (i + 1)
81                      dy = y + self.diry[d] * (i + 1)
82                      curr_state[dy][dx] = to_play
83                  bricks_taken += bricks
84              return (curr_state, bricks_taken)
85          else:
86              return print("Not valid move, retry")
87
88      def check_valid_move(self, curr_state, x, y, to_play):
89          """
90          Function checks playable moves. First if the agent is within the board, then
     checks
91          if the spot is occupied by a tot_black or white brick and finally, if the
     player do not
92          take any of the opponents bricks, then it is not a legal move.
93          """
94          if x < 0 or x > self.n - 1 or y < 0 or y > self.n - 1:
95              return False
96          if curr_state[y][x] != 0:
97              return False
98          (_, totctr) = self._check_valid_move(copy.deepcopy(curr_state), x, y, to_play)
99          if totctr == 0:
100             return False
101         return True
102
103     def _check_valid_move(self, board, x, y, to_play):
104         """
105         Helper function to check_valid_move function to not overwrite the playing
     board if move is illegal
106         """
107         n = self.n
108         bricks_taken = 0 # total number of opponent pieces taken
109
110         board[y][x] = to_play
111         for d in range(len(self.dirx)): # 8 directions
112             bricks = 0
113             for i in range(n):
114                 dx = x + self.dirx[d] * (i + 1)
115                 dy = y + self.diry[d] * (i + 1)
116                 if dx < 0 or dx > n - 1 or dy < 0 or dy > n - 1:
117                     bricks = 0; break
118                 elif board[dy][dx] == to_play:
119                     break
120                 elif board[dy][dx] == 0:
121                     bricks = 0; break
122                 else:
123                     bricks += 1
124             for i in range(bricks):
125                 dx = x + self.dirx[d] * (i + 1)
```

```python
126                    dy = y + self.diry[d] * (i + 1)
127                    board[dy][dx] = to_play
128                bricks_taken += bricks
129            return (board, bricks_taken)
130
131    def move_generator(self, curr_state, to_play):
132        possibleMoves = []
133        for i in range(self.n):
134            for j in range(self.n):
135                if(self.check_valid_move(curr_state, i, j, to_play)):
136                    possibleMoves.append((i, j))
137        return possibleMoves
138
139    def find_winner(self, curr_state):
140        tot_black = 0
141        tot_whites = 0
142
143        for i in range(self.n):
144            for j in range(self.n):
145                if (curr_state[i][j] == -1):
146                    tot_whites += 1
147                elif (curr_state[i][j] == 1):
148                    tot_black += 1
149
150        if (tot_black == tot_whites):
151            return 0
152        elif (tot_black > tot_whites):
153            return 1
154        else:
155            return -1
156
157
158
159 class MCTSNode(OthelloBoard):
160    def __init__(self, n, state, to_play, parent=None, parent_action=None):
161        super().__init__(n)
162        self.terminal_visits = 0
163        self.state = state
164        self.parent = parent
165        self.parent_action = parent_action
166        self.child_nodes = []
167        self._nof_visits = 0
168        self.player_turn = to_play
169        self.q_value = 0
170        self.p_action = 0
171        self.avg_q_value = 0
172        self.pass_counter = 0
173        self._untried_actions = self.untried_actions()
174        self.all_visited = False
175
176
177    def untried_actions(self):
178        self._untried_actions = self.move_generator(self.state, self.player_turn)
179
180        if len(self._untried_actions) == 0 and self.pass_counter != 2:
181            self.pass_counter += 1
182            self.player_turn *= -1
183            self.untried_actions()
184
185        return self._untried_actions
186
187    def expand(self):
188        action = self._untried_actions.pop()
189        next_state, _ = self.make_move(copy.deepcopy(self.state), action, self.player_
    turn)
```

```python
190             next_player = self.player_turn*-1
191             child_node = MCTSNode(self.n, next_state, next_player, parent=self, parent_
        action=action)
192             self.child_nodes.append(child_node)
193             if len(self._untried_actions) == 0:
194                 self.all_visited = True
195             return child_node
196
197     def update_q(self, val):
198             self.q_value = val
199             # return self.q_value
200
201     # generalize this function such that it works for something
202     def uniform_policy(self):
203             """
204             Initializes a policy uniformly over all legal actions.
205             """
206             nof_actions = len(self.child_nodes)
207             return 1 / nof_actions
208
209     def backpropagate(self, q_NN):
210             # self.acum_q_value += q_NN
211             self._nof_visits += 1
212             self.avg_q_value += (q_NN - self.avg_q_value)/self._nof_visits # Q_{n+1}
213             if self.parent: # Check if list is empty
214                 self.parent.backpropagate(q_NN)
215
216     def best_child(self, c):
217             """
218             Minimax for training two agents
219             """
220             if self.player_turn == 1: # max
221                 UCB_values = [child.avg_q_value + c * child.p_action * np.sqrt(self._nof_
        visits) / child._nof_visits
222                             for child in self.child_nodes]
223             #  + c_paramnp.sqrt(np.log(self._nof_visits)/child._nof_visits)
224                 return self.child_nodes[np.argmax(UCB_values)]
225             elif self.player_turn == -1: # min
226                 UCB_values = [child.avg_q_value - c * child.p_action * np.sqrt(self._nof_
        visits) / child._nof_visits
227                             for child in self.child_nodes]
228             #  - c_param*np.sqrt(np.log(self._nof_visits)/child._nof_visits)
229                 return self.child_nodes[np.argmin(UCB_values)]
230
231 def unpack_positions_returns(games):
232     """
233     Takes in games in the form of a list of tuples where the first element of the
234     tuple is a list of board positions representing one full game, and the second
235     element in the tuple is the reward (-1, 0, or 1) from that game.
236
237     Returns a list of tuples where each tuple consists of a single board position
238     and the reward associated with the game where that board position came from.
239     This is a (state s_t, return G_t) tuple in RL terms.
240     """
241     data = []
242     for game in games:
243         for position in game[0]:
244             data.append((position, game[1]))
245     return data
246
247 def unique_board_positions(state_return_tuples):
248     """
249     Takes in list of tuples (state, return) where state is a board position and
        returns
250     a list of only the unique tuples.
```

```python
251
252        This function will most likely be used only once , on the purely random simulated
           data .
253        """
254        unique_dict = {}
255        for state , target in state_return_tuples :
256            # Convert the matrix to a hashable representation using .tobytes()
257            key = (state.tobytes(), target)
258
259            # Only add unique tuples
260            if key not in unique_dict :
261                unique_dict[key] = (state, target)
262
263        # Saving the unique (board , reward) tuples as data list
264        return list(unique_dict.values())
265
266    def value_predictors_targets(state_return_tuples):
267        """
268        Takes in list of tuples (state , return) where state is board position and return
           is
269        observed reward from game .
270
271        And returns them in the form of X, y, ready to be used in a model .
272        """
273        # orders data into predictros and targets
274        X = np.array([x for x, _ in state_return_tuples])
275        y = np.array([y for _, y in state_return_tuples])
276
277        # expanding X to include #channels =1
278        X = np.expand_dims(X, axis=-1)
279
280        return X, y
281
282    # Unpickling and saving the games to a list
283    def read_files(path):
284        with open(path, "rb") as fp:
285            boards = pickle.load(fp)
286        return boards
287
288    def create_dataset(X, y, batch_size=128, shuffle_buffer_size=10000):
289        """
290        X and y should already be processed according to above preprocessing
291        functions before being passed to this function .
292
293        This function is simply to make training more efficient from memory .
294        """
295        # Create dataset from tensor slices
296        dataset = tf.data.Dataset.from_tensor_slices((X, y))
297
298        # Shuffle the dataset (reshuffles each epoch)
299        dataset = dataset.shuffle(buffer_size=shuffle_buffer_size, reshuffle_each_
           iteration=True)
300
301        # Batch the dataset
302        dataset = dataset.batch(batch_size)
303
304        # Prefetch for performance optimization
305        dataset = dataset.prefetch(tf.data.AUTOTUNE)
306
307        return dataset
308
309    def map_states_to_action_visits(sav):
310        state_to_action_visits = {}
311        for state, action, count in sav:
312            key = (state.tobytes())
```

```python
313            if key not in state_to_action_visits:
314                state_to_action_visits[key] = [(action, count)]
315            else:
316                state_to_action_visits[key].append((action, count))
317
318        return state_to_action_visits
319
320 def map_actions_to_integers(n):
321     actions = [(i,j) for i in range(n) for j in range(n)]
322     map_actions = dict(zip(actions, [i for i in range(n*n)]))
323     return map_actions
324
325
326 def policy_predictors_targets(sav, n):
327     """
328     Takes in a list of tuples (state, action, count) where state is board position,
329     action is the action taken in this state and the count is the number of visits
330     to the next state when taking the action.
331
332     It should be the N(s,a) count from the UCB algorithm.
333
334     The predictors are the states and the target will be the N(s,a) count
335     """
336     map_sav = map_states_to_action_visits(sav)
337     mapped_actions = map_actions_to_integers(n)
338     X = []
339     y = []
340
341     for state, value in map_sav.items():
342         x = np.frombuffer(state, dtype=int)
343         X.append(x)
344         y_i = np.zeros(n*n)
345         total_visits = 0
346         for action, visits in value:
347             total_visits += visits
348             y_i[mapped_actions[action]] = visits
349
350         y_i = y_i/total_visits
351         y.append(y_i)
352
353     X = np.array(X)
354     y = np.array(y)
355
356     return X, y
357
358 def episode(node):
359     episode_list = [np.array(node.state)]
360     if node.parent:
361         episode_list.extend(episode(node.parent))
362     return episode_list
363
364 def treetraversal(node, res):
365     """
366     Recursive helper function to state_action_visits.
367     """
368     if not node:
369         return
370
371     if node.parent != None:
372         tup = (np.array(node.parent.state), node.parent_action, node._nof_visits)
373         res.append(tup)
374
375     for child in node.child_nodes:
376         treetraversal(child, res)
377
```

```python
378
379  def state_action_visits(root):
380      """
381      Given MCTS, recursively goes through the tree and returns a list with tuples
382      containing state s, action a and N(s,a).
383      """
384      res = []
385      treetraversal(root, res)
386      return res
387
388
389  def results_distribution(episodes):
390      """
391      Given episodes, return a tuple with win, draw and loss.
392      """
393      win = 0
394      loss = 0
395      draw = 0
396
397      for i in range(len(episodes)):
398          if episodes[i][1] == -1:
399              loss += 1
400
401          elif episodes[i][1] == 0:
402              draw += 1
403
404          elif episodes[i][1] == 1:
405              win += 1
406
407      return (win, draw, loss)
408
409
410  def terminal_state_visits(episodes):
411      """
412      Given episodes, returns a dictionary with terminal state as key
413      and number of visits in that terminal state as value.
414      """
415      term_state_visits = {}
416
417      for i in range(len(episodes)):
418          term_state = (episodes[i][0][0].tobytes(),)   # Make terminal state hashable
419          term_state_visits[term_state] = term_state_visits.get(term_state, 0) + 1
420
421      return term_state_visits
422
423
424  def bar_plot_term_states(episodes):
425      """
426      Given episodes, plot the distribution over how often a terminal state is visited.
427      """
428      win, draw, loss = results_distribution(episodes)
429      tot_games = len(episodes)
430      win_proc = np.round(win/tot_games*100, decimals=4)
431      draw_proc = np.round(draw/tot_games*100, decimals=4)
432      loss_proc = np.round(loss/tot_games*100, decimals=4)
433
434      term_state_visits = terminal_state_visits(episodes)
435      x_vals = [i for i in range(len(term_state_visits))]
436      plt.bar(x_vals, term_state_visits.values())
437
438      # Add a text box with game results
439      textstr = f"Total games = {tot_games}\nWin: {win_proc}%\nDraw: {draw_proc}%\nLoss:
          {loss_proc}%"
440
441      # Position the text box
```

```python
442     props = dict(boxstyle='round', facecolor='white', alpha=0.8)
443     plt.text(0.68, 0.9, textstr, transform=plt.gca().transAxes, fontsize=9,
444             verticalalignment='top', bbox=props)
445
446     plt.xlabel("Terminal states")
447     plt.ylabel("Number of visits")
448     plt.grid()
449
450
451 def online_value_training(value_model, replay_buffer, epochs=20, batch_size=64):
452     """
453     Given a replay buffer in the form of list of tuples, (where each tuple consists of
454     a list of board positions from one game, and the reward from said game) and a
        value model,
455     trains this model using the replay buffer through given number of epochs.
456
457     Returns the trained value_model
458     """
459     state_return_tuples = unpack_positions_returns(replay_buffer)
460     X_buffer, y_buffer = value_predictors_targets(state_return_tuples)
461     X_train_onl, X_test_onl, y_train_onl, y_test_onl = train_test_split(X_buffer, y_
        buffer,
462                                                         test_size=0.2,
463                                                         random_state
        =80085)
464
465     train_dataset = create_dataset(X_train_onl, y_train_onl)
466     val_dataset = create_dataset(X_test_onl, y_test_onl)
467
468     history_onl = value_model.fit(
469         train_dataset,
470         batch_size=batch_size,
471         epochs=epochs,
472         validation_data=val_dataset
473     )
474
475     return value_model, history_onl
476
477 def online_policy_training(policy_model, tree, epochs=20, n=6, batch_size=64):
478     """
479     Given a replay buffer in the form of list of state-actions and their visits and a
        policy model,
480     trains this model using the replay buffer through given number of epochs.
481
482     Returns the trained policy_model
483     """
484     savs = state_action_visits(tree)
485     X_buffer, y_buffer = policy_predictors_targets(savs, n)
486     X_train_onl, X_test_onl, y_train_onl, y_test_onl = train_test_split(X_buffer, y_
        buffer,
487                                                         test_size=0.2,
488                                                         random_state
        =80085)
489
490     train_dataset = create_dataset(X_train_onl, y_train_onl)
491     val_dataset = create_dataset(X_test_onl, y_test_onl)
492
493     history_onl = policy_model.fit(
494         train_dataset,
495         batch_size=batch_size,
496         epochs=epochs,
497         validation_data=val_dataset
498     )
499
500     return policy_model, history_onl
```

```python
501
502  def simulate_random_games(n):
503      game_boards = []
504      game = OthelloBoard(n)
505
506      while True:
507          moves = game.move_generator()
508          if moves == []:
509              game.pass_counter += 1
510              game.change_turn()
511
512          else:
513              game.pass_counter = 0 # reset counter
514              action = random.choice(moves)
515              board = game.make_move(action)[0]
516              game_boards.append(np.array(board))
517              game.change_turn()
518
519          if game.pass_counter == 2:
520              reward = game.find_winner()
521              break
522
523      return (game_boards, reward)
524
525  def save_games(n, nr_simulations):
526      all_game_boards = []
527
528      for _ in range(nr_simulations):
529          boards, actions = simulate_random_games(n)
530          all_game_boards.append(boards)
531
532      with open(f"othello_sim_boards_{nr_simulations}", "wb") as fp:   #Pickling
533          pickle.dump(all_game_boards, fp)
534
535  def build_tree(n, value_model, policy_model, c, nr_simulations=1000):
536
537      game = OthelloBoard(n)
538      root = MCTSNode(n, game.board, game.to_play)
539      episodes = []
540      mapped_actions = map_actions_to_integers(n)
541
542      while True:
543          terminal_state, reward = expand_tree_iteratively(root, value_model, policy_
     model, mapped_actions, c)
544          episodes.append((episode(terminal_state)[:-1], reward))
545          print('Episode done!')
546
547          if len(episodes) == nr_simulations:
548              break
549
550      return root, episodes
551
552  def expand_tree_iteratively(root, value_model, policy_model, mapped_actions, c):
553      stack = [(root, root)]  # Stack holds pairs of (root, root)
554
555      while stack:
556          current_root, current_node = stack.pop()  # Get the last node to process
557
558          if current_node._untried_actions:  # If the node has untried actions, expand
     it
559              next_node = current_node.expand()
560              # Instead of using recursion, use a random value for the q_val
561              q_val = np.array(value_model(np.expand_dims(np.array(next_node.state),
     axis=(0, -1))))[0][0]
562              next_node.update_q(q_val)
```

```
563             next_node.backpropagate(next_node.q_value)
564
565             # Push the root back into the stack to continue processing from the root
566             stack.append((current_root, current_root))
567
568         else:  # If no untried actions, move to the best child or terminal node
569             if current_node.pass_counter != 2:
570
571                 flatten_state = np.ndarray.flatten(np.array(current_node.state))
572                 policy_dist = np.array(policy_model(np.expand_dims(flatten_state, axis
    =0)))[0])
573                 for child in current_node.child_nodes:
574                     child.p_action = policy_dist[mapped_actions[child.parent_action]]
575
576                 best_child = current_node.best_child(c)
577                 stack.append((current_root, best_child))  # Continue with the best
    child
578
579             else:  # Terminal node/state
580                 reward = current_node.find_winner(current_node.state)
581                 current_node.update_q(reward)
582                 current_node.backpropagate(reward)
583                 current_node.terminal_visits += 1
584                 return current_node, reward
585
586 def unpack_plot_history(history_list):
587     for i, history in enumerate(history_list):
588         if history != []:
589             label = f"C{i}"
590             plt.plot(history.history['loss'], label=f'Training loss', color = label)
591             plt.plot(history.history['val_loss'], linestyle="dashed", label = '
    Validation loss', color = label)
592
593     plt.xlabel('Epoch')
594     plt.ylabel('Loss')
595     plt.legend(loc='upper right')
596     plt.show()
597
598 def cycle_results_histogram(episodes_list):
599     """
600     Takes in an ordered (from first to last training cycle) episodes list where the
    length of the
601     list denotes the number of training cycles that were needed to create it (each
    element in the
602     list should be a list of episodes).
603
604     In return, plots a histogram with sections colored according to the number of
    white wins, black
605     wins, and draws (grey).
606
607     As white should win at 6x6 othello, we expect to see the bars further to the right
     to be
608     increasingly dominated by the white sections.
609     """
610     results = []
611     for episodes in episodes_list:
612         results.append(results_distribution(episodes))
613
614     # Generate x-axis indices for the number of tuples in the list
615     x = range(len(results))
616     labels = list(range(1, len(results) + 1))
617
618     # Unpack each tuple into separate segments
619     segment1 = [t[0] for t in results]  # Black section
620     segment2 = [t[1] for t in results]  # Grey section
```

```python
621      segment3 = [t[2] for t in results]  # White section
622
623      plt.figure(figsize=(11, 5))
624
625      # Plot the first segment (black)
626      plt.bar(x, segment1, color='black', edgecolor='black')
627
628      # Plot the second segment (grey), stacked on top of the first
629      plt.bar(x, segment2, bottom=segment1, color='grey', edgecolor='black')
630
631      # Compute the bottom for the third segment (segment1 + segment2)
632      bottom3 = [a + b for a, b in zip(segment1, segment2)]
633
634      # Plot the third segment (white)
635      plt.bar(x, segment3, bottom=bottom3, color='white', edgecolor='black')
636
637      # # Add a text annotation above each bar showing the white percentage.
638      # for i, t in enumerate(results):
639      #     total = sum(t)
640      #     # Calculate the white percentage (t[2] is the white segment)
641      #     white_percentage = (t[2] / total) * 100 if total else 0
642      #     # Use .3g to format the number with a maximum of 3 significant digits
643      #     plt.text(i, total + 1, f'{white_percentage:.3g}%', ha='center', va='bottom',
         fontsize=5)
644
645      plt.xlabel('Training cycle')
646      plt.ylabel('Episodes')
647      plt.ylim(0, max([sum(t) for t in results]) * 1.1)  # Slightly higher than max for
         visual clarity
648      plt.xticks(x, labels)
649      plt.show()
650
651  def terminal_visits_histogram_colored(episodes):
652      """
653      Takes in a list of episodes and returns a histogram with as many bars as
654      unique terminal states that were visited in those episodes, with each bar
655      colored according to who wins (or draw) in that position. White bars for
656      white wins, black bars for black wins, and grey for draws.
657      """
658
659      nr_games = len(episodes)
660      win, draw, loss = results_distribution(episodes)
661      w_rate, d_rate, l_rate = np.round(np.array([win, draw, loss])*(100/nr_games),
         decimals=0)
662
663      terminal_states = {}
664      for episode in episodes:
665          key = (episode[0][0].tobytes(), episode[1])
666          if key not in terminal_states:
667              terminal_states[key] = 1
668          else:
669              terminal_states[key] += 1
670
671      color_mapping = {
672          -1: 'white',  # white for -1
673          0: 'grey',   # grey for 0
674          1: 'black'   # black for 1
675      }
676
677      # Extract rewards, visits, and colors
678      identifiers = []
679      heights = []
680      colors = []
681      for (identifier, y_value), height in terminal_states.items():
682          if isinstance(identifier, bytes):
```

```python
683             identifier = identifier.decode('latin-1')
684         identifiers.append(identifier)
685         heights.append(height)
686         colors.append(color_mapping[y_value])
687
688     x = range(len(terminal_states))
689
690     plt.figure(figsize=(7.2, 5))
691     plt.bar(identifiers, heights, color=colors, edgecolor='black')
692
693     # Add a text box with game results
694     textstr = f"Total games = {nr_games}\nBlack wins: {w_rate}%\nDraws: {d_rate}%\
    nWhite wins: {l_rate}%"
695
696     # Position the text box
697     props = dict(boxstyle='round', facecolor='white', alpha=0.8)
698     plt.text(0.025, 0.975, textstr, transform=plt.gca().transAxes, fontsize=9,
699             verticalalignment='top', bbox=props)
700
701
702     plt.xticks(x, [str(i + 1) for i in x])
703
704     plt.gca().yaxis.set_major_locator(mticker.MaxNLocator(integer=True))
705
706     plt.xlabel('Unique terminal states')
707     plt.ylabel('Visits')
708
709     plt.show()
710
711 def remove_empty_lists(list):
712     # remove all empty lists from the list
713     list = [x for x in list if x]
714     return list
715
716 def unpack_plot_train_val_curves(history_list, cycle_labels):
717     """
718     All empty lists have to removed from the input list here in order for it to work.
719     See above helper function.
720
721     Cycle_labels should be a list of of labels for the cycles, i.e. [3,6,9,...] if
722     every third cycle is passed in the history list
723     """
724
725     fig = plt.figure(figsize=(12, 6))
726
727     for i, history in enumerate(history_list):
728         plt.subplot(2,5,i+1)
729         color = f"C{i}"
730
731         plt.plot(history.history['loss'], label=f'Training loss', color = color)
732         plt.plot(history.history['val_loss'], linestyle="dashed", label = 'Validation
    loss', color = color)
733         plt.title(f'Cycle {cycle_labels[i]}')
734         plt.tick_params(axis='x', which='both', bottom=False,
735             top=False, labelbottom=False)
736         plt.tick_params(axis='y', which='both', right=False,
737             left=False, labelleft=False)
738
739     textstr = '\u2500' * 5 + '  Training loss \n--------  Validation loss'
740     props = dict(boxstyle='round', facecolor='white', alpha=0.8)
741     plt.text(-0.125, -0.10, textstr, transform=plt.gca().transAxes, fontsize=12,
742             verticalalignment='top', bbox=props)
743
744     fig.text(0.5, 0.04, 'Epoch', va='center', ha='center', fontsize=14)
```

```
745     fig.text(0.09, 0.5, 'Loss', va='center', ha='center', rotation='vertical',
        fontsize=14)
746     plt.show()
747
748 ###########################################################################
749 # ----------------- Beginning of VALUE NN architecture -----------------
750 ###########################################################################
751
752 def residual_block(x, channels=64, kernel_size=(3,3), weight_decay=0.001):
753     shortcut = x  # No projection needed if dimensions already match.
754
755     x = layers.Conv2D(channels, kernel_size=kernel_size,
756                       padding='same', use_bias=False,
757                       kernel_regularizer=regularizers.l2(weight_decay))(x)
758     x = layers.BatchNormalization()(x)
759     x = layers.ReLU()(x)
760
761     x = layers.Conv2D(channels, kernel_size=kernel_size,
762                       padding='same', use_bias=False,
763                       kernel_regularizer=regularizers.l2(weight_decay))(x)
764     x = layers.BatchNormalization()(x)
765
766     # Direct addition is fine here.
767     x = layers.Add()([x, shortcut])
768     x = layers.ReLU()(x)
769
770     return x
771
772 # Build the model
773 inputs = layers.Input(shape=(6, 6, 1)) # 6x6 Othello
774
775 # Initial convolution block (producing 16 channels)
776 x = layers.Conv2D(64, kernel_size=(3,3),
777                   padding='same', use_bias=False,
778                   kernel_regularizer=regularizers.l2(0.001))(inputs)
779 x = layers.BatchNormalization()(x)
780 x = layers.ReLU()(x)
781
782 # Add a number of residual blocks with constant 64 channels
783 x = residual_block(x, channels=64, kernel_size=(3,3), weight_decay=0.001)
784 x = residual_block(x, channels=64, kernel_size=(3,3), weight_decay=0.001)
785
786 # Flatten the features and output a single scalar value (for a value network)
787 x = layers.Flatten()(x)
788 outputs = layers.Dense(1, name="value_output")(x)
789
790 # Create the model
791 value_model = Model(inputs=inputs, outputs=outputs)
792 value_model.summary()
793
794 # Optimizer for the value model
795 value_model.compile(
796     optimizer=keras.optimizers.SGD(learning_rate=1e-4, momentum=0.9),
797     loss = keras.losses.MeanSquaredError()
798     # metrics = [keras.metrics.RootMeanSquaredError]
799 )
800
801 ###########################################################################
802 #  --------------- Beginning of Policy NN architecture -----------------
803 ###########################################################################
804
805 # Policy model architecture
806 policy_model = keras.Sequential([
807
808     layers.Dense(128, input_shape=(36,)),
```

```python
809        layers.BatchNormalization(),
810        layers.Activation("relu"),
811        layers.Dropout(0.025),
812
813        layers.Dense(256),
814        layers.BatchNormalization(),
815        layers.Activation("relu"),
816        layers.Dropout(0.025),
817
818        layers.Dense(256),
819        layers.BatchNormalization(),
820        layers.Activation('relu'),
821        layers.Dropout(0.025),
822
823        layers.Dense(256),
824        layers.BatchNormalization(),
825        layers.Activation("relu"),
826        layers.Dropout(0.025),
827
828        layers.Dense(128),
829        layers.BatchNormalization(),
830        layers.Activation('relu'),
831        layers.Dropout(0.025),
832
833        layers.Dense(36),
834        layers.Activation('softmax')
835 ])
836
837 policy_model.summary()
838
839 policy_model.compile(optimizer=keras.optimizers.SGD(learning_rate=1e-3, momentum=0.9),
840                      loss=keras.losses.KLDivergence())
841
842 # %%
843 # Loading the model weights
844 # Load the model weights
845 # value_model.load_weights('zeroth_value_nn.weights.h5')
846 # policy_model.load_weights('zeroth_policy_nn.weights.h5')
847
848
849 ########################################################################
850 # ---------------------- Online training section ----------------------
851 ########################################################################
852
853 """
854 Given a replay buffer, we want to be able to continously feed new game
855 information into the value network in the form of mini batches. This
856 section aims to prepare for thats.
857
858 Replay buffer will be in the form of a list of tuples, where the first
859 elements in each tuple is a game consisting of a sequence of board positions
860 and the second element is the observed reward from that game.
861
862 Prerequisites for the following function is to already have pre-trained
863 value and policy networks.
864 """
865
866 # %%
867 def online_simulation(n, c, value_model, policy_model, nr_cycles, nr_episodes_per_tree
         = 160, history_val = [], history_pol = []):
868     ''' hehiha '''
869
870     trees = []
871     episodes = []
872
```

```python
873     value_model_history = [history_val]
874     policy_model_history = [history_pol]
875
876     for _ in range(nr_cycles):
877
878         ##### Build tree with updated model #####
879         tree, replay_buffer = build_tree(n, value_model, policy_model, c, nr_episodes_
        per_tree)
880
881         trees.append(tree)
882         episodes.append(replay_buffer)
883
884         ##### Update value network #####
885         value_model, value_history_temp = online_value_training(value_model, replay_
        buffer,
886                                                                  epochs=20, batch_size
        =64)
887         policy_model, policy_history_temp = online_policy_training(policy_model, tree,
888                                                                    epochs=50, batch_size
        =128)
889         value_model_history.append(value_history_temp)
890         policy_model_history.append(policy_history_temp)
891
892         print('Cycle done!')
893
894
895     return trees, episodes, value_model_history, policy_model_history
896
897
898 # %%
899 trees, episodes, val_hist, pol_hist = online_simulation(n=6, c=0.5,
900                                                         value_model=value_model,
901                                                         policy_model=policy_model,
902                                                         nr_cycles=10,
903                                                         nr_episodes_per_tree=25)
904
905 # %%
906 # saving the training from the online training
907 second_10cycles25episodes_c05_data = [trees, episodes, val_hist, pol_hist]
908 with open('second_10cycles25episodes_c05_data', 'wb') as handle:
909     pickle.dump(second_10cycles25episodes_c05_data,
910                 handle, protocol=pickle.HIGHEST_PROTOCOL)
911
912 value_model.save_weights('value_second_10cycles25episodes_c05.weights.h5')
913 policy_model.save_weights('policy_second_10cycles25episodes_c05.weights.h5')
914
915 #######################################################################
916 # ------------------ Initial training of VALUE NN --------------------
917 #######################################################################
918
919 # Initial value network data - from completely random games
920 # Unique board positions from first 100k random games
921 path = './othello_random_simulations/othello_sim_boards_100000_6x6'
922 games = read_files(path)
923
924 # %%
925 state_return_tuples = unpack_positions_returns(games)
926 state_return_tuples = unique_board_positions(state_return_tuples)
927
928 X, y = value_predictors_targets(state_return_tuples)
929
930 # splitting data into training and validation
931 X_train_val, X_test_val, y_train_val, y_test_val = train_test_split(X, y, test_size
    =0.2,
```

```
932                                                              random_state
      =80085)
933
934 train_dataset_value_model = create_dataset(X_train_val, y_train_val, batch_size=256,
935                                       shuffle_buffer_size=10000)
936 val_dataset_value_model = create_dataset(X_test_val, y_test_val, batch_size=256,
937                                     shuffle_buffer_size=10000)
938
939 # Training the value model
940 history_value = value_model.fit(
941     train_dataset_value_model,
942     epochs = 50,
943     validation_data=val_dataset_value_model
944     # callbacks=[keras.callbacks.EarlyStopping(monitor='val_loss', patience=10)]
945 )
946
947 # with open('initial_value_training_history', 'wb') as handle:
948 #     pickle.dump(history_value, handle, protocol=pickle.HIGHEST_PROTOCOL)
949
950 # # Saving the weights from the first good network
951 # value_model.save_weights('zeroth_value_nn.weights.h5')
952
953 # %%
954 ############################################################################
955 #   ------------------ Initial training of Policy NN   --------------------
956 ############################################################################
957
958 sav = read_files("./othello_random_simulations/othello_sim_sa_visits")
959 X, y = policy_predictors_targets(sav, 6)
960
961 X_train_pol, X_test_pol, y_train_pol, y_test_pol = train_test_split(X, y, test_size
      =0.2,
962                                             random_state=80085)
963
964 train_dataset_policy_model = create_dataset(X_train_pol, y_train_pol, batch_size=128)
965 val_dataset_policy_model = create_dataset(X_test_pol, y_test_pol, batch_size=128)
966
967 history_policy = policy_model.fit(train_dataset_policy_model,
968                       validation_data=val_dataset_policy_model,
969                       epochs=500)
970
971 # with open('initial_policy_training_history_2', 'wb') as handle:
972 #     pickle.dump(history_policy, handle, protocol=pickle.HIGHEST_PROTOCOL)
973
974 # saving weights from first good network
975 # policy_model.save_weights('zeroth_policy_nn_2.weights.h5')
976
977
978 ############################################################################
979 #   -------------- Concatenating history from different runs -------------
980 ############################################################################
981
982 # concatenating and plotting history from first and second iteration of VALUE network
983 initial_value_training_history_1 = read_files('initial_value_training_history')
984 initial_value_training_history_2 = read_files('initial_value_training_history_2')
985
986 complete_value_history_train_loss = initial_value_training_history_1.history['loss'] +
        initial_value_training_history_2.history['loss']
987 complete_value_history_val_loss = initial_value_training_history_1.history['val_loss']
        + initial_value_training_history_2.history['val_loss']
988
989 plt.plot(complete_value_history_train_loss, label='Training loss')
990 plt.plot(complete_value_history_val_loss, label = 'Validation loss')
991 plt.xlabel('Epoch')
992 plt.ylabel('Loss')
```

```python
993  plt.title('value network')
994  plt.legend(loc='upper right')
995  plt.show()
996
997  # %%
998  # concatenating and plotting history from first and second iteration of POLICY network
999  initial_policy_training_history_1 = read_files('initial_policy_training_history')
1000 initial_policy_training_history_2 = read_files('initial_policy_training_history_2')
1001
1002 complete_policy_history_train_loss = initial_policy_training_history_1.history['loss']
         + initial_policy_training_history_2.history['loss']
1003 complete_policy_history_val_loss = initial_policy_training_history_1.history['val_loss
     '] + initial_policy_training_history_2.history['val_loss']
1004
1005 plt.plot(complete_policy_history_train_loss, label='Training loss')
1006 plt.plot(complete_policy_history_val_loss, label = 'Validation loss')
1007 plt.xlabel('Epoch')
1008 plt.ylabel('Loss')
1009 plt.title('policy network')
1010 plt.legend(loc='upper right')
1011 plt.show()
```