

NEURAL NETWORKS AND GAMES: OTHELLO

Samuel Effendy
Vira Kasprova

Othello Reinforcement Learning

Reinforcement learning has been shown to be an effective method to solve complex problems in the domain of board games, as it trains the agent to maximize its expected reward by learning the value of each action through entering specific states. In this paper, we explore its application on Othello, a two-player deterministic, zero-sum, perfect information game.

Theoretical Background

Othello.

Othello is a perfect information game because every player is perfectly informed of every event in the game, including the initialization of the board. It is also zero-sum, since every gained advantage of a player is equal to the gained disadvantage of their opponent, e.g. every tile gained by a player is lost to another player. Next, it is deterministic, because there are no random elements in its rules. This game is a perfect test bed for reinforcement algorithms like Q-learning, since there is a clear distinction between reward and consequences as well as a numeric value for each action to train the agent.

The board is initially empty but for the four center squares, which feature discs on the two squares on the major diagonal ([4][4] and [5][5]) and the two squares on the opposing diagonal, respectively [4][5] and [5][4].

The player who moves first is referred to as “Black” and the player who moves second is referred to as “White”. A disc must be placed on an empty square in order for there to be a sequence of one or more discs belonging to the opponent, followed by one’s own disc, in at least one direction (horizontally, vertically, or diagonally) from the square played on. In such a series, the opponent’s discs are flipped and changed to one’s own color.

Othello’s state space has a maximum of 60 movements and a dimension of around 1028. Two players use 64 two-sided discs that are black on one side and white on the other to play Othello on an 8 by 8 board. The discs are placed on the board with the white side up by one player and the black side up by the other player.

Parameterization of Othello.

We chose to write our parameterization in Python due to its accessibility and library support for machine learning. To represent the board in our version, we will utilize an 8 by 8 character array. Instead of using black and white pieces, we will use Xs and Os. The lower left corner is at coordinate [0][0], while the upper right corner is at coordinate [7][7].

```
1  def __init__(self,p1,p2):
2      self.board = np.full((8,8),".")
3      #set four entries to be two X's and two 0's
4      self.board[3][4] = "X";
5      self.board[4][3] = "X";
6      self.board[3][3] = "0";
7      self.board[4][4] = "0";
8      self.p1 = p1
9      self.p2 = p2
10     self.isEnd = False
11     self.boardHash = None
12     self.playerSymbol = "X" #first player symbol
```

Python code

Since the rules for movement in Othello are complicated, we broke the movements down into three functions. The `validMove()` checks if the move is valid or not. The move function takes as input the board, the XY coordinate to place the piece, and the piece we are placing (X or O). If a piece can be flipped in any of the eight directions and the space on the board is unoccupied, the move is valid. We check all adjacent cells to our current location by calling `checkFlip()` function for left, right, down, up, down-left, down-right, up-left, and up-right.

```
1 def validMove(boardList, x, y, piece):
2     # Check that the coordinates are empty
3     if boardList[x][y] != '.':
4         return False
5     # Figure out the character of the opponent's piece
6     opponent = 'O'
7     if piece == 'O':
8         opponent = 'X'
9     # If we can flip in any direction, it is valid
10    # Check to the left
11    if checkFlip(boardList, x - 1, y, -1, 0, piece, opponent):
12        return True
13    # Check to the right
14    if checkFlip(boardList, x + 1, y, 1, 0, piece, opponent):
15        return True
16    # Check down
17    if checkFlip(boardList, x, y - 1, 0, -1, piece, opponent):
18        return True
19    # Check up
20    if checkFlip(boardList, x, y + 1, 0, 1, piece, opponent):
21        return True
22    # Check down-left
23    if checkFlip(boardList, x - 1, y - 1, -1, -1, piece, opponent):
24        return True
25    # Check down-right
26    if checkFlip(boardList, x + 1, y - 1, 1, -1, piece, opponent):
27        return True
28    # Check up-left
29    if checkFlip(boardList, x - 1, y + 1, -1, 1, piece, opponent):
30        return True
31    # Check up-right
32    if checkFlip(boardList, x + 1, y + 1, 1, 1, piece, opponent):
33        return True
34
35    return False # If we get here, we didn't find a valid flip direction
```

Python code

We use makeMove() function to make an actual move. We check all adjacent cells to our current location by calling checkFlip() function for left, right, down, up, down-left, down-right, up-left, and up-right. If the flip is possible we call flipPieces() function to flip the value of the opponent piece to the player piece and vice versa.

```

1 def makeMove(self, boardList, x, y, piece): #this function is copied over to state class
2     # Put the piece at x,y
3     boardList[x][y] = piece
4     # Figure out the character of the opponent's piece
5     opponent = 'O'
6     if piece == 'O':
7         opponent = 'X'
8
9     # Check to the left
10    if checkFlip(boardList, x - 1, y, -1, 0, piece, opponent):
11        flipPieces(boardList, x - 1, y, -1, 0, piece, opponent)
12    # Check to the right
13    if checkFlip(boardList, x + 1, y, 1, 0, piece, opponent):
14        flipPieces(boardList, x + 1, y, 1, 0, piece, opponent)
15    # Check down
16    if checkFlip(boardList, x, y-1, 0, -1, piece, opponent):
17        flipPieces(boardList, x, y-1, 0, -1, piece, opponent)
18    # Check up
19    if checkFlip(boardList, x, y + 1, 0, 1, piece, opponent):
20        flipPieces(boardList, x, y + 1, 0, 1, piece, opponent)
21    # Check down-left
22    if checkFlip(boardList, x-1, y - 1, -1, -1, piece, opponent):
23        flipPieces(boardList, x-1, y - 1, -1, -1, piece, opponent)
24    # Check down-right
25    if checkFlip(boardList, x + 1, y - 1, 1, -1, piece, opponent):
26        flipPieces(boardList, x + 1, y - 1, 1, -1, piece, opponent)
27    # Check up-left
28    if checkFlip(boardList, x - 1, y + 1, -1, 1, piece, opponent):
29        flipPieces(boardList, x - 1, y + 1, -1, 1, piece, opponent)
30    # Check up-right
31    if checkFlip(boardList, x + 1, y + 1, 1, 1, piece, opponent):
32        flipPieces(boardList, x + 1, y + 1, 1, 1, piece, opponent)

```

Python code

This function has two helper functions. One function, `checkFlip()`, checks if an existing disk on the board can be flipped in one of eight directions (left, right, up, down, up-left, up-right, down-left, down-right), with the condition that an opponent disk is found between our current and next piece. It returns true or false.

```
1 def checkFlip(boardList, x, y, deltaX, deltaY, piece, opponent):
2     if (x >= 0 and x < 8 and y >= 0 and y < 8):
3         if(boardList[x][y] == opponent):
4             x += deltaX
5             y += deltaY
6             while (x >= 0) and (x < 8) and (y >= 0) and (y < 8):
7                 if(boardList[x][y] == '.'):
8                     return False
9                 if(boardList[x][y] == piece):
10                    return True
11 return False
```

Python code

This is then used in an if statement to decide whether or not the second helper function `flipPieces()` should be executed, with `flipPieces()` doing the actual flipping. It simply changes the value of the opponent piece to the player piece and vice versa.

```
1 def flipPieces(boardList, x, y, deltaX, deltaY, myPiece, opponentPiece):
2     while boardList[x][y] == opponentPiece:
3         boardList[x][y] = myPiece
4         x += deltaX
5         y += deltaY
```

Python code

The function that calculates the score is called `score()`. The function takes as input the board and the character (X or O). It loops through the board and counts the number of a particular character (X or O).

```
1 def score(boardList, piece):
2     total = 0
3     for x in range(8):
4         for y in range(8):
5             if boardList[x][y] == piece:
6                 total += 1
7     return total
```

Python code

CHALLENGES

We had difficulties integrating Q-learning into our parameterized program. This is because our code was written to play Othello manually, without opponent AI. Furthermore, the structure of our code was divided into snippets containing separate functions, without grouping of similar functions under the agent and the state it inhabits. Upon consulting our advisor, we realized that we needed to have two classes and sort our functions beneath it.

Finding the next move of our game was challenging after modifying our code, since we weren't sure of which functions to implement inside a class, and which class. Eventually, we figured out that all three functions should be called in the state class, because we are checking for moves in each state that the agent enters.

```
1 def availablePositions(self):
2     #this is the function that returns all possible moves from a particular board
3     positions = []
4     for i in range(8):
5         for j in range(8):
6             if validMove(self.board,i,j,self.playerSymbol):
7                 positions.append((i,j))
8             print(validMove(self.board,i,j,self.playerSymbol), end = ' ')
9     for i in positions:
10        print(i, end = ' ')
11    return positions
```

Python code

Checking if the game was finished was trivial in our original implementation, due to the rule that an Othello game ends if there are no empty tiles left, regardless of the scores of each player. That means that there could be a tie, and this was taken into consideration in our program. However, as we currently have no final snapshot of the game finished through implementation of Q-learning, we cannot say for sure.

METHOD

Initially, we were thinking of using an open-source version of the AlphaZero algorithm, but we chose Q-learning to train our program. In particular, we used a Tic-Tac-Toe with Q-learning examples provided by our advisors as a baseline. After doing our research, Q-learning is efficient because it doesn't require a model of the environment (model-free), saving us time in implementing our program.

Furthermore, Q-learning finds an optimal rule to maximize the reward for the agent in every state, by adding the maximum reward from future states to the reward it achieves in the current state. It is most suitable for board games like Othello, since the game has deterministic states and actions, and we can override old information with new information instantly (e.g. white replaces all black tiles in between, so history of black tiles is no longer saved.)

```
X
Available moves are [(3, 2), (4, 5), (5, 4)]
(3, 2)
[['0' 'X' '.' '.' '.' '.' '.' '.']
 ['.' '0' '.' '.' '.' '.' '.' '.']
 ['.' 'X' '0' 'X' '.' '.' '.' '.']
 ['.' '.' 'X' 'X' 'X' '.' '.' '.']
 ['.' '.' '.' 'X' '0' '.' '.' '.']
 ['.' '.' '.' '.' '.' '.' '.' '.']
 ['.' '.' '.' '.' '.' '.' '.' '.']
 ['.' '.' '.' '.' '.' '.' '.' '.']]
0
Available moves are [(0, 2), (2, 0), (2, 4), (3, 1), (4, 2)]
(0, 2)
[['0' '0' '0' '.' '.' '.' '.' '.']
 ['.' '0' '.' '.' '.' '.' '.' '.']
 ['.' 'X' '0' 'X' '.' '.' '.' '.']
 ['.' '.' 'X' 'X' 'X' '.' '.' '.']
 ['.' '.' '.' 'X' '0' '.' '.' '.']
 ['.' '.' '.' '.' '.' '.' '.' '.']
 ['.' '.' '.' '.' '.' '.' '.' '.']
 ['.' '.' '.' '.' '.' '.' '.' '.']]
```

Game visualization and output

X and O are both played by the AI. Here, we see the AI choosing from available moves to either X or O for each turn.

Playing the game with two AI players using Q-learning requires us to automate checking for available positions, choosing the agent's actions, and updating and adding states. This is reflected in the function below.

```

1  def play(self, rounds=10000):
2      global player1wins
3      global player2wins
4      global ties
5      ties = 0
6      player1wins = 0
7      player2wins = 0
8      for i in range(rounds):
9          if i % 1000 == 0:
10             print("Round {}".format(i))
11             while not self.isEnd:
12                 # Player 1
13                 positions = self.availablePositions()
14                 if len(positions) != 0:
15                     p1_action = self.p1.chooseAction(positions, self.board, self.playerSymbol)
16                     # take action and update board state
17                     self.updateState(p1_action)
18                     board_hash = self.getHash()
19                     self.p1.addState(board_hash)
20                 #Player 2
21                 self.playerSymbol = "O"
22                 positions = self.availablePositions()
23                 if len(positions) != 0:
24                     # Player 2
25                     p2_action = self.p2.chooseAction(positions, self.board, self.playerSymbol)
26                     self.updateState(p2_action)
27                     board_hash = self.getHash()
28                     self.p2.addState(board_hash)
29
30             win = self.winner()
31             if win is not None:
32                 #print("Winner is", self.defineWinner())
33                 self.defineWinner()
34                 self.showBoard()
35                 # ended with p1 either win or draw
36                 self.giveReward()
37                 self.p1.reset()
38                 self.p2.reset()
39                 self.reset()
40             break

```

Python code

We wrote functions to test out simulations between random players, trained and random players, and between trained players. Then, we took the success rates of our simulations using the function below.

```

1  def percentage(numGames, player1wins, player2wins, ties):
2      #global player1wins
3      #global player2wins
4      #global ties
5      print("Player 1 won", player1wins, "times")
6      print("Player 2 won", player2wins, "times")
7      print("There were ties", ties, "times")
8      print("Number of games:", numGames, "\n")
9
10     percentOne = (player1wins / numGames) * 100
11     percentTwo = (player2wins / numGames) * 100
12     ties = (ties / numGames) * 100
13
14     print("Player 1 percentage:", percentOne)
15     print("Player 2 percentage:", percentTwo)
16     print("Ties percentage:", ties)

```



```
Available moves are []
Winner is 0!
[['0' '0' '0' '0' '0' '0' '0' '0']
 ['0' '0' '0' '0' '0' '0' 'X' '0']
 ['0' '0' '0' '0' '0' 'X' '0' '0']
 ['0' '0' '0' 'X' 'X' '0' 'X' '0']
 ['0' '0' 'X' '0' 'X' 'X' 'X' '0']
 ['0' '0' '0' '0' 'X' '0' 'X' '0']
 ['0' '0' '0' '0' '0' 'X' '0' '0']
 ['0' '0' '0' '0' '0' '0' '0' '0']]
Result of random AI vs random AI
Player 1 won 36 times
Player 2 won 63 times
There were ties 1 times
Number of games: 100

Player 1 percentage: 36.0
Player 2 percentage: 63.0
Ties percentage: 1.0
```

```

1 def randomVRandom():
2
3     numGames = 10
4     st = State(p1,p2)
5     st.play(numGames)
6     print("Result of random AI vs random AI")
7     percentage(numGames, player1wins, player2wins, ties)
8     st.reset()

```

Python code

```

Available moves are []
Winner is X!
[['X' 'X' 'X' 'X' 'X' 'X' 'X' '0']
 ['X' 'X' 'X' 'X' 'X' 'X' 'X' '0']
 ['X' 'X' 'X' '0' 'X' '0' 'X' '0']
 ['X' 'X' '0' 'X' 'X' 'X' '0' '0']
 ['X' '0' 'X' '0' 'X' 'X' 'X' '0']
 ['X' '0' '0' 'X' '0' 'X' '0' 'X']
 ['X' '0' '0' '0' 'X' '0' '0' '0']
 ['0' '0' '0' '0' '0' '0' '0' '0']]
Result of random AI vs trained AI
Player 1 won 36 times
Player 2 won 64 times
There were ties 0 times
Number of games: 100

Player 1 percentage: 36.0
Player 2 percentage: 64.0
Ties percentage: 0.0

```

```

1 def randomVTrained():
2
3     numGames = 10
4     st = State(p1,p2)
5     st.play(numGames)
6
7
8     p1.savePolicy()
9     st.play(numGames)
10    print("Result of random AI vs trained AI")
11    percentage(numGames, player1wins, player2wins, ties)
12
13    st.reset()

```

Python code

```

Available moves are []
Winner is 0!
[['0' '0' '0' '0' '0' '0' '0' '0']
 ['0' '0' '0' '0' '0' '0' '0' '0']
 ['0' '0' '0' '0' '0' '0' '0' '0']
 ['0' '0' '0' '0' 'X' '0' '0' '0']
 ['0' '0' '0' 'X' '0' '0' '0' '0']
 ['0' '0' '0' '0' 'X' '0' 'X' 'X']
 ['0' '0' '0' '0' '0' '0' '0' 'X']
 ['X' '.' 'X' 'X' 'X' 'X' 'X' 'X']]
Result of trained AI vs trained AI
Player 1 won 33 times
Player 2 won 67 times
There were ties 0 times
Number of games: 100

Player 1 percentage: 33.0
Player 2 percentage: 67.0
Ties percentage: 0.0

```

```

1 def trainedVTrained():
2
3     numGames = 10
4     st = State(p1,p2)
5     st.play(numGames)
6
7     p1.savePolicy()
8     p2.savePolicy()
9     p1.loadPolicy("policy_p1")
10    p2.loadPolicy("policy_p2")
11
12
13    st.play(numGames)
14    print("Result of trained AI vs trained AI")
15    percentage(numGames, player1wins, player2wins, ties)
16    st.reset()

```

Python code

What are some further directions you could take the program next? How could you further improve your trained player?

We might attempt to decrease the amount of time required to train the algorithm. We completed step 1 by creating a functional prototype, regardless of its performance or the problem's complexity. Next, we might take measures to reduce memory use and evaluate alternative network configurations and technical options to improve precision. In addition, we might attempt to apply several methods and compare them based on their time complexity and precision.

Next, we may create more effective graphical user interfaces to facilitate the audience's data reading. Lastly, we might potentially use several ML techniques that were utilized by other research teams, such as the minimax algorithm.