Hosam Farouk

Q-Learning Agent **VS** Random Agent:



Wins: 50 | Losses: 0 | Draws: 0

Q-Learning Agent **VS** Defensive Agent:



Wins: 28 | Losses: 0 | Draws: 22

Q-Learning Agent **VS** Aggressive Agent:



Wins: 50 | Losses: 0 | Draws: 0

Function Implementations in 2-3 sentences:

### train():

- First initialised a variable episode to track how many episodes have been completed by the agent in the environment.
- Then loop for the number of iterations required and in the beginning of each iteration the environment is reset, and as long as the state is non-terminal it follows the epislon greedy strategy on whether to exploit or explore in the next step based on the num genereated by Math.random()
- After updating Q values based on the rewards and transitions, it extracts the optimal policy from the learned Q values

### extractPolicy():

- First initialized a Policy Object.

- Then it extracts the optimal policy from the Q-table by iterating through all states and determining the move with the highest Q value for each non-terminal state using the Collections.max method. The extractedPolicy.policy.put.. maps states to their corresponding optimal moves.

### getQValueForAllMoves():

- A helper method used to retrieve the Q values for all possible moves for a state which is saved in the variable moves using getPossibleMoves() method.

- It returns a HashMap where each move is a key, and its corresponding Q value is the value.

### updateQValue():

- A helper method that updates the Q value for a state-move pair using the Q learning formula, incorporating the immediate reward and the discounted max arg Q value of the next state. The updated Q-value is then stored in the Q table.