

# Computational Intelligence Activity Log

Giulio Maselli, s306125

A.A. 2023/24



**Politecnico  
di Torino**

# 1 Lab 1: Set Covering A\*

This section provides a detailed look at the Python implementation of the set-covering problem using the A\* search algorithm.

```
1 import numpy as np
2 from queue import PriorityQueue
3 from random import random
4 from functools import reduce
5 from collections import namedtuple
6
7 PROBLEM_SIZE = 50
8 NUM_SETS = 100
9 SETS = tuple(np.array([random() < 0.2 for _ in range(PROBLEM_SIZE)
10 ]))
11
12     for _ in range(NUM_SETS))
13 State = namedtuple('State', ['taken', 'not_taken'])
14
15 # Stima del numero di elementi che devono essere ancora coperti dai
16   set nello stato attuale
17 def h(state):
18     return PROBLEM_SIZE - sum(
19         reduce(
20             np.logical_or,
21             [SETS[i] for i in state.taken],
22             np.array([False for _ in range(PROBLEM_SIZE)]),
23         )
24     )
25
26 # Numero di elementi coperti dai set nello stato attuale
27 def g(state):
28     covered_elements = sum(
29         reduce(
30             np.logical_or,
31             [SETS[i] for i in state.taken],
32             np.array([False for _ in range(PROBLEM_SIZE)]),
33         )
34     )
35     return PROBLEM_SIZE - covered_elements
36
37 # Verifica se lo stato attuale comprende tutti gli elementi
38 def goal_check(state):
39     return np.all(reduce(np.logical_or, [SETS[i] for i in state.
40 taken], np.array([False for _ in range(PROBLEM_SIZE)])))
41
42 assert goal_check(State(set(range(NUM_SETS)), set()), "Problem not
43   solvable"
44
45 frontier = PriorityQueue()
46 initial_state = State(set(), set(range(NUM_SETS)))
47 frontier.put((h(initial_state), initial_state))
48
49 counter = 0
50 _, current_state = frontier.get()
```

```

48 # A*
49 while not goal_check(current_state):
50     counter += 1
51     for action in current_state.not_taken:
52         new_state = State(
53             current_state.taken ^ {action},
54             current_state.not_taken ^ {action})
55         frontier.put((g(new_state) + h(new_state), new_state))
56     _, current_state = frontier.get()
57 print(f"Solved in {counter} steps ({len(current_state.taken)} tiles
58       )")
59 print(current_state)
60 # Solved in 6 steps (6 tiles)
61 # State(taken={1, 70, 50, 51, 53, 86}, not_taken={0, 2, 3, 4, 5, 6,
62       7, 8, 9, 10, 11, 12, 13, 14, 15...})

```

Here the A\* algorithm initializes the search with an empty '**taken**' set, and all sets in '**not\_taken**', then uses a priority queue to manage states, prioritizing them based on the sum of the cost so far ('**g**') and the estimated cost ('**h**')

## 2 Lab 2: Nim with ES

That's my implementation of a game-playing agent for the game of Nim, using Evolution Strategies (ES). Nim is a mathematical game of strategy in which two players take turns removing objects (matches) from distinct piles.

The backbone code already had three implemented strategies, *pure\_random*, *optimal* and *gabriele*, which I transformed into *evolvable\_based\_on\_ES\_1\_lambda*.

1. The 'select\_parent' function randomly chooses a 'parent' from the current population, which will then be used to generate new moves (offspring). This selection is done in such a way that moves with a better 'fitness' (in this context a measure of how close a move is to winning the game) are more likely to be chosen.

```
1 def select_parent(population, tournament_size=2):
2     return min(random.choices(population, k=tournament_size),
3                 key=lambda i: i.fitness)
```

2. The 'evolve\_one\_plus\_lambda' function is responsible for creating new offspring from the selected parents. This is done by applying the 'mutation' function to the parents to create a new set of moves.

```
1     def evolve_one_plus_lambda(parent, mutation_rate, state):
2         child = mutation(parent, state)
3         return child
```

3. The 'mutation' function introduces variation into the population. It takes a parent move and creates a new move (child) by randomly altering the parent's move. This could involve changing the row from which objects are removed or the number of objects taken. Mutation introduces new strategies into the population and is essential for exploring the space of possible strategies.

```
1     def mutation(parent, nim):
2         if parent is None:
3             # in assenza di genitore, genera una mossa casuale
4             row = random.choice([r for r, c in enumerate(nim.rows)
5                                 if c > 0])
6             num_objects = random.randint(1, nim.rows[row])
7             a = nim_sum(nim)
8             offspring = Move(row, num_objects, a)
9         else:
10            if nim.k is None:
11                elements = random.randrange(1, nim.rows[parent.row]
12                                           + 1)
13            else:
14                elements = min(nim.k, random.randrange(1, nim.rows
15                                                       [parent.row] + 1))
16            temp_rows = nim.rows.copy()
17            temp_rows[parent.row] -= elements
18            offspring = Move(parent.row, elements, nim_sum(
19                             temp_rows))
20        return offspring
```

4. The 'evolvable\_based\_on\_ES\_1\_lambda' function implements an evolutionary strategy (ES) specifically the (1+lambda) strategy, to evolve game moves for playing Nim. The process aims to discover an effective move by simulating evolution over a series of generations.

```
1  def evolvable_based_on_ES_1_lambda(self, state: Nim, alpha
    = 0.5):
2      # Popolazione iniziale di possibili mosse
3      population = [mutation(None, state) for _ in range(
        POPULATION_SIZE)]
4
5      for g in range(N_GENERATIONS):
6          offspring = []
7
8          for i in range(LAMBDA):
9              # Selezione del genitore (nel tuo caso,
              potrebbe essere tramite il torneo)
10             parent = select_parent(population)
11
12             # Generazione del figlio mediante mutazione
              del genitore
13             child = mutation(parent, state)
14
15             # Aggiornamento del genitore con alpha
16             parent.row = int(parent.row + alpha * (child.
                row - parent.row))
17             parent.num_objects = int(parent.num_objects +
                alpha * (child.num_objects - parent.num_objects))
18
19             # Aggiungi il figlio alla lista degli
              offspring
20             offspring.append(child)
21
22             # Unisci genitori e figli per formare la nuova
              popolazione
23             population += offspring
24
25             # Ordina la popolazione in base al fitness e
              seleziona i primi POPULATION_SIZE individui
26             population = sorted(population, key=lambda i: i.
                fitness, reverse=False)[:POPULATION_SIZE]
27             logging.debug(f"Actual best {population[0]}")
28
29             # Restituisci la migliore mossa
30             best_move = Nimply(population[0].row, population[0].
                num_objects)
31             return best_move
```

### 3 Lab 9: Genome and Evolutionary Strategy

This code outlines an implementation of an evolutionary algorithm designed to solve an optimization problem, with the ultimate goal of minimizing the number of fitness function calls. The optimization problem in question involves finding a binary genome sequence that achieves the highest fitness according to a given fitness function.

The key components are as usual the 'parent selection', 'crossover' and 'mutation'.

```
1 class EvolutionaryAlgorithm:
2     def __init__(self, problem_instance, genome_length,
3         population_size=50, generations=100, mutation_rate=0.01,
4         elitism_percentage=0.1):
5         self.problem = problem_instance
6         self.genome_length = genome_length
7         self.population_size = population_size
8         self.generations = generations
9         self.mutation_rate = mutation_rate
10        self.elitism_percentage = elitism_percentage
11        self.best_fitness_history = [] # To store the best fitness
12        in each generation
13
14    def initialize_population(self):
15        return [choices([0, 1], k=self.genome_length) for _ in
16            range(self.population_size)]
17
18    def crossover(self, parent1, parent2):
19        # xover
20        crossover_point = randint(1, self.genome_length - 1)
21        child = parent1[:crossover_point] + parent2[crossover_point
22            :]
23        return child
24
25    def mutate(self, genome):
26        # mutation
27        mutated_genome = [bit ^ (random() < self.mutation_rate) for
28            bit in genome]
29        return mutated_genome
30
31    def run_algorithm(self):
32        population = self.initialize_population()
33        best_fitness_per_generation = []
34
35        for generation in range(1, self.generations + 1):
36            # Evaluate fitness for each individual
37            fitness_values = [self.problem(ind) for ind in
38                population]
39
40            # Store the best fitness in each generation
41            best_fitness_per_generation.append(max(fitness_values))
42
43            # top individuals
44            sorted_indices = sorted(range(len(fitness_values)), key
```

```

39     =lambda k: fitness_values[k], reverse=True)
    selected_parents = [population[i] for i in
sorted_indices[:self.population_size // 2]]

40
41     # offspring creation
42     offspring = []
43     for i in range(0, len(selected_parents) - 1, 2):
44         parent1 = selected_parents[i]
45         parent2 = selected_parents[i + 1]
46         child = self.crossover(parent1, parent2)
47         child = self.mutate(child)
48         offspring.append(child)
49
50     # if the population size is odd, add the last ind
51     if len(selected_parents) % 2 == 1:
52         offspring.append(selected_parents[-1])
53
54     population = selected_parents + offspring
55
56     # best fitness in each generation
57     best_fitness = max(fitness_values)
58     print(f"Generation {generation}, Best Fitness: {
best_fitness:.2%}")
59
60     # best ind in the final gen
61     best_index = max(range(len(fitness_values)), key=lambda k:
fitness_values[k])
62     best_individual = population[best_index]
63
64     # Store the best fitness history for each generation
65     self.best_fitness_history = best_fitness_per_generation
66
67     return best_individual
68
69     def get_best_fitness_history(self):
70         return self.best_fitness_history

```

### Best Solution Overall:

- Population size: *1000*
- Generations: *100*
- Mutation rate: *0.01*
- Elitism percentage: *0.3*
- Fitness: *75.00%*
- Number of fitness calls: *75253*

## 4 Lab 10: Tic-Tac-Toe Q-Learning Project

This code outlines an implementation of a **Q-Learning algorithm** for playing Tic-Tac-Toe against a random opponent. The implementation involves training a Q-Learning agent (**QLearningPlayer**) to learn optimal strategies through repeated gameplay against a **RandomPlayer**, who selects moves randomly. The process aims to reinforce winning strategies and penalize losing ones, leveraging the Q-Learning framework to iteratively improve the agent's performance.

### 4.1 Core components

- **TicTacToe**: Represents the game environment with methods to make moves, check for wins or draws, and get available moves. The game board is a 3x3 numpy array, and players are represented by 1 (X) and -1 (O).
- **RandomPlayer**: An agent that selects moves randomly from the available spots on the board.
- **QLearningPlayer**: A reinforcement learning (RL) agent that uses Q-Learning to decide on the best moves. It maintains a Q-table mapping states to action values, updated based on the outcome of each game.

```
1 class QLearningPlayer:
2 def __init__(self, symbol, alpha=0.1, gamma=0.9, epsilon=0.2):
3     self.symbol = symbol # 1 per X, -1 per O
4     self.alpha = alpha # learning rate
5     self.gamma = gamma # discount factor
6     self.epsilon = epsilon # exploring prob (epsilon greedy)
7     self.q_table = {} # state -> action -> value
8
9 def get_state(self, game):
10     return str(game.board.reshape(9))
11
12 def choose_move(self, game):
13     state = self.get_state(game)
14     if np.random.rand() < self.epsilon:
15         # exploring
16         return random.choice(game.get_available_moves())
17     else:
18         # exploiting
19         self.q_table.setdefault(state, {})
20         if not self.q_table[state]:
21             # if none -> rand
22             return random.choice(game.get_available_moves())
23         best_move = max(self.q_table[state], key=self.q_table[
24             state].get)
25         return eval(best_move)
26
27 def update_q_values(self, prev_state, action, reward,
28     next_state):
29     self.q_table.setdefault(prev_state, {})
30     self.q_table.setdefault(next_state, {})
31     prev_q = self.q_table[prev_state].get(str(action), 0)
```



```

30     max_next_q = max(self.q_table[next_state].values(),
31                       default=0)
32     self.q_table[prev_state][str(action)] = prev_q + self.
33     alpha * (reward + self.gamma * max_next_q - prev_q)
34
35 def update_q_table(self, game_history):
36     # game_history = list (state, action, reward, next_state)
37     for i in range(len(game_history) - 1):
38         state, action, reward, next_state = game_history[i]
39         self.update_q_values(state, action, reward, next_state
40                             )
41
42     # update q-values for final_state
43     final_state, final_action, final_reward, _ = game_history
44     [-1]
45     self.update_q_values(final_state, final_action,
46                           final_reward, final_state)
47
48 def receive_reward(self, game):
49     # rewards values
50     if game.is_winner(self.symbol):
51         return 1 # win
52     elif game.is_draw():
53         return 0 # draw, maybe -0.5
54     else:
55         return -1 # loss
56
57 def save_q_table(self, filename='q_table.pkl'):
58     with open(filename, 'wb') as f:
59         pickle.dump(self.q_table, f)
60     print(f"Q-table saved in {filename}.")
61
62 def load_q_table(self, filename='q_table.pkl'):
63     with open(filename, 'rb') as f:
64         self.q_table = pickle.load(f)
65     print(f"Q-table loaded from {filename}.")

```

- Training:

```

1 def play_game(q_player, random_player, episodes=100000):
2     for episode in range(episodes):
3         game = TicTacToe()
4         game_history = [] # init game history
5
6         while True:
7             current_player = q_player if game.current_player ==
8             q_player.symbol else random_player
9             move = current_player.choose_move(game)
10            if move is None: # if none break loop
11                break
12
13            # save current_state if qlplayer
14            if isinstance(current_player, QLearningPlayer):
15                prev_state = current_player.get_state(game)
16
17            game.make_move(*move)

```

```

18         if isinstance(current_player, QLearningPlayer):
19             next_state = current_player.get_state(game)
20             reward = 0
21             game_history.append((prev_state, move, reward,
next_state))
22
23         # winner check
24         if game.is_winner(q_player.symbol) or game.is_winner(
random_player.symbol) or game.is_draw():
25             final_reward = q_player.receive_reward(game) #
final reward
26
27         # update reward for qlplayer
28         if isinstance(current_player, QLearningPlayer) and
game_history:
29             prev_state, move, _, next_state = game_history
[-1]
30             game_history[-1] = (prev_state, move,
final_reward, next_state)
31             break
32
33         # update qtable with complete history
34         if game_history:
35             q_player.update_q_table(game_history)
36
37         game.reset() # next game
38
39     print("Learning completed after {} episodes".format(epochs))

```

- Testing win rate on 1000 games: [00:01;00:00, 756.95it/s]  
Win rate: 77.60%  
Total games: **1000**, Wins: **776**, Draws: **59**, Losses: **165**

## 5 Quixo

### 5.1 Introduction:

The essence of this project is to develop a framework where different types of AI players, including a Deep Q-Learning (DQL) agent and a Minimax agent, can learn and compete in the game of Quixo.

### 5.2 Key Components

#### 5.2.1 Game Environment (Game.py)

Central to the project, this class manages the game state, enforces the rules of Quixo, and determines the winner. It allows moves to be made on the board, checks for win conditions, and ensures the game progresses logically from start to finish.

#### 5.2.2 Players

- **RandomPlayer:** Serves as a baseline opponent, making moves by randomly selecting from available options.

```
1 class RandomPlayer(Player):
2     def __init__(self) -> None:
3         super().__init__()
4
5     def make_move(self, game) -> tuple[tuple[int, int], Move]:
6         from_pos = (random.randint(0, 4), random.randint(0, 4))
7
8         move = random.choice([Move.TOP, Move.BOTTOM, Move.LEFT,
9                               Move.RIGHT])
10        return from_pos, move
```

- **Minimax Player:** Implements the Minimax algorithm with a specified depth of lookahead. It evaluates game states to make the most advantageous moves, considering the possible responses of the opponent. This player introduces a strategic challenge, forcing learning agents to develop more sophisticated strategies to win. The class includes methods for initializing the player, evaluating the game state, running the Minimax algorithm to choose the best move, and executing the selected move.

```
1 class MiniMaxPlayer(Player):
2     def __init__(self, depth: int, player_index) -> None:
3         super().__init__()
4         self.depth = depth
5         self.player_index = player_index
```

```

7     def evaluate_game_state(self, board):
8         player_score, player_4_in_row = count_aligned(board,
9             self.player_index, 4)
10        opponent_score, opponent_4_in_row = count_aligned(
11            board, 1 - self.player_index, 4)
12
13        # evaluation of important pieces (corners) and
14        aligned pieces
15
16        strategic_value = 0
17        corners = [(0, 0), (0, len(board)-1), (len(board)-1,
18            0), (len(board)-1, len(board)-1)]
19        for x, y in corners:
20            if board[x, y] == self.player_index:
21                strategic_value += 1
22            elif board[x, y] == 1 - self.player_index:
23                strategic_value -= 1
24
25        # critical situations
26        if opponent_4_in_row:
27            return -10 # avoid opponent's win
28        if player_4_in_row:
29            return 10 # win the game
30
31        return player_score - opponent_score + strategic_value
32
33    def minimax(self, board, depth, alpha, beta,
34        maximizing_player):
35
36        if depth == 0:
37            return self.evaluate_game_state(board)
38
39        if maximizing_player:
40            max_eval = -float('inf')
41            for move in get_possible_moves(board, self.
42                player_index):
43                # apply the move and compute new board state
44                new_board = apply_move(board, move, self.
45                    player_index)
46                eval = self.minimax(new_board, depth - 1,
47                    alpha, beta, False)
48                # print(f"Valutazione mossa {move}: {eval}")
49            # Debugging
50            max_eval = max(max_eval, eval)
51            alpha = max(alpha, eval)
52            if beta <= alpha:
53                break
54            return max_eval
55        else:
56            min_eval = float('inf')
57            for move in get_possible_moves(board, 1 - self.
58                player_index):
59                # apply the move and compute new board state
60                new_board = apply_move(board, move, 1 - self.
61                    player_index)
62                eval = self.minimax(new_board, depth - 1,

```

```

alpha, beta, True)
    min_eval = min(min_eval, eval)
    beta = min(beta, eval)
    if beta <= alpha:
        break
    return min_eval

def make_move(self, game: 'Game') -> tuple[tuple[int, int], Move]:
    best_score = -float('inf')
    best_action = None # tuple (from_pos, move)

    for action in get_possible_moves(game.get_board(), self.player_index):
        # apply the move to obtain the new board
        new_board = apply_move(game.get_board(), action, self.player_index)
        score = self.minimax(new_board, self.depth, -float('inf'), float('inf'), True)
        if score > best_score:
            best_score = score
            best_action = action

        # print(f"MiniMaxPlayer best_action: from_pos={best_action[0]}, move={best_action[1]}")
        # print_board(new_board)
    return (best_action[0][1], best_action[0][0]), best_action[1] if best_action is not None else ((0, 0), Move.TOP)

'''
depth = 5 has a great win rate (between 75% and 95%), almost
20 seconds per game

depth = 5 vs. Random
100%
Percentuale vittorie: 94.0%
Percentuale pareggi: 0.0%
Percentuale sconfitte: 6.0%

depth = 3 vs. Random
100%
Percentuale vittorie: 86.0%
Percentuale pareggi: 0.0%
Percentuale sconfitte: 14.0%
'''

```

- **Q-Learning Player:** designed to implement a reinforcement learning agent using the Q-Learning algorithm. This agent is capable of learning optimal moves and strategies by interacting with the game environment and adapting based on the outcomes of its actions.

**Initialization (`__init__`):** Sets up the learning agent with initial parameters for learning rate (alpha), discount factor (gamma), exploration rate (epsilon), and initializes the Q-table.

**Q-Table Persistence (`load_q_table`, `save_q_table`):** These functions handle the loading and saving of the Q-table to disk, allowing the agent to retain its learned knowledge between game sessions and continue improving over time.

```

1 class QLearningPlayer(Player):
2     def __init__(self, player_index, alpha=0.2, gamma=0.9,
3         epsilon=0.2, epsilon_decay=1, epsilon_min=0, preload=True)
4         :
5         super().__init__()
6         self.alpha = alpha # Learning rate
7         self.gamma = gamma # Discount factor
8         self.epsilon = epsilon # Exploration probability
9         self.epsilon_decay = epsilon_decay # Epsilon decay
10        rate
11        self.epsilon_min = epsilon_min # Minimum epsilon
12        self.q_table = {} # Initialize Q-table
13        self.moves_history = [] # Track moves
14        self.player_index = player_index
15        self.win_count = 0
16        self.draw_count = 0
17        self.loss_count = 0
18
19        if preload:
20            self.q_table = self.load_q_table()
21
22    def load_q_table(self):
23        try:
24            with open("q_table.pickle", "rb") as f:
25                print("File Q-table caricato.")
26                return pickle.load(f)
27        except FileNotFoundError:
28            print("File Q-table non trovato, inizializzazione
29            di una nuova Q-table.")
30            return {}
31    def save_q_table(self):
32        try:
33            with open("q_table.pickle", "wb") as f:
34                pickle.dump(self.q_table, f)
35                print("Q-table salvata con successo.")
36        except Exception as e:
37            print(f"Errore nel salvataggio della Q-table: {e}")
38    )

```

**State Representation (state\_representation):** Converts the current game board into a standardized format that serves as a key for accessing the Q-table, ensuring that the agent can evaluate and learn from specific game states.

**Move Selection (make\_move):** Determines the agent's actions within the game by either exploring new moves randomly or exploiting known moves with high Q-values, based on the epsilon-greedy strategy. This function encapsulates the core of the agent's decision-making process.

```

1 def state_representation(self, board):
2     return str(board.reshape(-1))
3
4 def make_move(self, game) -> tuple[tuple[int, int], Move]:
5     board = game.get_board()
6     state = self.state_representation(board)
7
8     possible_moves = get_possible_moves(board, self.player_index)
9
10    if random.random() < self.epsilon or state not in self.q_table
11        :
12        selected_move = random.choice(possible_moves)
13    else:
14        max_q_value = float("-inf")
15        selected_action_key = None
16        for move in possible_moves:
17            action_key = self.move_to_key(move)
18            q_value = self.q_table.get(state, {}).get(action_key,
19                float("-inf"))
20            if q_value > max_q_value:
21                max_q_value = q_value
22                selected_action_key = action_key
23
24        # Se nessuna mossa ha Q superiore, scelta casuale tra le
25        mosse
26        if selected_action_key is None:
27            selected_move = random.choice(possible_moves)
28        else:
29            selected_move = self.key_to_move(selected_action_key)
30
31    # Simula mossa selezionata per valutare il suo effetto
32    simulated_board = apply_move(deepcopy(board), selected_move,
33        self.player_index)
34    simulated_state = self.state_representation(simulated_board)
35    reward = self.evaluate_move_effect(game, selected_move)
36
37    # Calcola il miglior Q per il prossimo stato simulato
38    next_max = max(self.q_table[simulated_state].values(), default
39        =0) if simulated_state in self.q_table else 0
40    self.update_q_table(state, selected_move, reward, next_max)
41    self.epsilon = max(self.epsilon * self.epsilon_decay, self.
42        epsilon_min)
43    #print(f"Epsilon: {self.epsilon}")
44
45    self.moves_history.append((state, selected_move,
46        simulated_state))

```

```

41 #print(f"{board}")
42 #print(f"Selected Move: {selected_move[0][0], selected_move
    [0][1]}")
43 #print(f"Possible Moves {possible_moves}")
44 selected_move = ((selected_move[0][1], selected_move[0][0]),
    selected_move[1])
45
46 return selected_move
47
48
49 def move_to_key(self, move):
50     # Converte la mossa in una chiave utilizzabile nella Q-
    table
51     return (move[0], move[1].value)
52
53 def key_to_move(self, key):
54     # Converte una chiave della Q-table in una mossa
55     return (key[0], Move(key[1]))
56

```

**Learning from Experience (update\_q\_table, learn):** Updates the Q-values in the Q-table based on the agent's experiences, applying the Q-Learning formula to adjust values towards optimal strategies over time.

**Evaluating Actions (evaluate\_move\_effect):** Assesses the immediate impact of actions to assign rewards, guiding the learning process by providing feedback on the effectiveness of moves in achieving favorable outcomes or preventing unfavorable ones.

```

1 def update_q_table(self, state, action, reward, next_max):
2     if state not in self.q_table:
3         self.q_table[state] = {}
4
5     action_key = self.move_to_key(action)
6     if action_key not in self.q_table[state]:
7         self.q_table[state][action_key] = 0 # Inizializza se
        non presente
8
9     # Aggiorna Q usando next_max, ovvero il massimo valore Q
    per il prossimo stato
10    self.q_table[state][action_key] += self.alpha * (reward +
        self.gamma * next_max - self.q_table[state][action_key])
11
12    '''
13    def adjust_epsilon_dynamically(self):
14        # Riduce epsilon piu' lentamente se l'agente sta vincendo
        o numero di pareggi alto
15        if self.win_count > self.loss_count or self.draw_count > (
            self.win_count + self.loss_count):
16            self.epsilon *= (self.epsilon_decay ** 0.5) #
            Riduzione piu' lenta
17        elif self.loss_count >= self.win_count:
18            # Aumenta epsilon leggermente se ci sono molte
            sconfitte, esploraz. piu' ampia
19            self.epsilon = min(self.epsilon / (self.epsilon_decay
                ** 0.5), 1.0)
20

```



```

21     self.epsilon = max(self.epsilon, self.epsilon_min)
22     '''
23
24 def evaluate_move_effect(self, game, move):
25     current_board = game.get_board()
26     ql_max_elements_before, ql_otk_before = count_aligned(
27         current_board, self.player_index)
28     opp_max_elements_before, opp_otk_before = count_aligned(
29         current_board, 1 - self.player_index)
30
31     # Simula la mossa
32     new_board = apply_move(deepcopy(current_board), move, self
33                             .player_index)
34     ql_max_elements_after, ql_otk_after = count_aligned(
35         new_board, self.player_index)
36     opp_max_elements_after, opp_otk_after = count_aligned(
37         new_board, 1 - self.player_index)
38
39     # reward basato su differenza e prevenzione di mosse
40     # critiche
41     reward = 0
42
43     # Miglioramento della posizione
44     if ql_max_elements_after > ql_max_elements_before:
45         reward += (ql_max_elements_after -
46                   ql_max_elements_before) * 2
47
48     # Blocco dell'avversario
49     if opp_otk_before and not opp_otk_after:
50         reward += 2
51
52     # Vittoria imminente
53     if ql_otk_after and not ql_otk_before:
54         reward += 3
55
56     # Vittoria
57     if ql_max_elements_after == 5:
58         reward += 6
59
60     # Avversario vicino alla vittoria
61     if not opp_otk_before and opp_otk_after:
62         reward -= 2
63
64     # Sconfitta
65     if opp_max_elements_after > 4:
66         reward -= 5
67
68     # mosse neutre
69     if ql_max_elements_after <= ql_max_elements_before and not
70       ql_otk_after:
71         reward -= 1
72
73     return reward
74
75 def learn(self, game, move, reward):
76     board = game.get_board()
77     state = self.state_representation(board)

```

```

70     next_state = deepcopy(board)
71     apply_move(next_state, move, self.player_index) # Simula
        la mossa sullo stato
72
73     next_state_rep = self.state_representation(next_state)
74     next_max = max(self.q_table[next_state_rep].values(),
        default=0) if next_state_rep in self.q_table else 0
75     self.update_q_table(state, move, reward, next_max)
76

```

**Game Finalization (finalize\_game):** Processes the result of a game to update learning parameters and the Q-table based on the outcome, reinforcing successful strategies and penalizing unsuccessful ones. This step ensures that the agent learns from the entire game, not just individual moves.

```

1  def finalize_game(self, result):
2  if result == "win":
3      self.win_count += 1
4      reward = 10 # vittoria
5  elif result == "draw":
6      reward = 1 # pareggio
7      self.draw_count += 1
8  else: # "lose"
9      reward = -10 # sconfitta
10     self.loss_count += 1
11 if (self.win_count + self.draw_count + self.loss_count) % 10
    == 0:
12     self.win_count, self.draw_count, self.loss_count = 0, 0, 0
13 # Aggiorna la Q-table retroattivamente sulle mosse effettuate
    durante la partita
14 for state, action, next_state in reversed(self.moves_history):
15     current_q_value = self.q_table.get(state, {}).get(self.
        move_to_key(action), 0)
16     next_max = 0
17     if next_state in self.q_table:
18         next_max = max(self.q_table[next_state].values())
19     self.q_table.setdefault(state, {})[self.move_to_key(action
        )] = current_q_value + self.alpha * (reward + self.gamma *
        next_max - current_q_value)
20 # Resetta la history delle mosse per la prossima partita
21 self.moves_history.clear()
22

```

- **Deep-QLearning Player:** This agent leverages deep learning to approximate the optimal action-value function, enabling it to make strategic decisions based on the current state of the game.

**Model Construction (`_build_model`):** Builds the neural network model that approximates the Q-function, determining the value of taking a particular action in a given state. This network is the core of the DQN, guiding the agent's decision-making process.

**Initialization (`__init__`):** Sets up the agent with necessary parameters, including the size of the state and action spaces, and initializes the replay memory for experience replay, which is critical for stabilizing and improving learning outcomes.

```

1 class DQNAgent(Player):
2     def __init__(self, player_index, state_size, action_size):
3         super(DQNAgent, self).__init__()
4         self.state_size = state_size
5         self.action_size = action_size
6         self.memory = deque(maxlen=2000) # Buffer di
7             riproduzione per memorizzare esperienze passate
8         self.gamma = 0.95 # discount rate
9         self.epsilon = 0.1 # exploration rate
10        self.epsilon_min = 0.01
11        self.epsilon_decay = 1
12        self.learning_rate = 0.001 # Tasso di apprendimento
13        per l'ottimizzatore
14        self.player_index = player_index
15        self.model = self._build_model()
16
17    def _build_model(self):
18        # Costruisce il modello della rete neurale con
19        struttura lineare
20        model = nn.Sequential(
21            nn.Linear(self.state_size, 24),
22            nn.ReLU(),
23            nn.Linear(24, 24),
24            nn.ReLU(),
25            nn.Linear(24, self.action_size)
26        )
27        self.optimizer = optim.Adam(model.parameters(), lr=
28        self.learning_rate)
29        return model

```

**Experience Storage (remember):** Records experiences (state, action, reward, next state, done) in the replay memory, enabling the agent to learn from a diverse set of experiences by revisiting past states and outcomes.

**Replay (replay):** Samples a batch of experiences from the replay memory to update the neural network. This process involves calculating target Q-values and optimizing the network's parameters to reduce the discrepancy between predicted and target Q-values, facilitating learning from past actions.

```

1 def remember(self, state, action, reward, next_state, done):
2     # Memorizza l'esperienza passata nel buffer di
    riproduzione
3     self.memory.append((state, action, reward, next_state,
4         done))
5
6 def replay(self, batch_size, gamma):
7     # Effettua l'aggiornamento del modello utilizzando il
    buffer di riproduzione
8     minibatch = random.sample(self.memory, batch_size)
9     for state, action, reward, next_state, done in minibatch:
10        # Convert state and next_state to PyTorch tensors
11        state_tensor = torch.FloatTensor(state).unsqueeze(0)
12        # Adds a batch dimension
13        next_state_tensor = torch.FloatTensor(next_state).
14        unsqueeze(0) # Adds a batch dimension
15
16        target = reward
17        if not done:
18            # Use next_state_tensor instead of next_state
19            target = (reward + gamma * torch.max(self.model(
20                next_state_tensor)).item())
21
22        # Use state_tensor instead of state
23        target_f = self.model(state_tensor)
24        target_f[0][action] = target
25
26        self.optimizer.zero_grad()
27        # Calculate loss using the tensor, not the NumPy array
28        loss = nn.MSELoss()(target_f, self.model(state_tensor))
29
30    )
31
32    loss.backward()
33    self.optimizer.step()

```

**Action Selection (act):** Decides on actions based on the current state, employing an epsilon-greedy strategy that balances exploration (choosing random actions) and exploitation (choosing the best-known action according to the neural network).

**Move Translation (index\_to\_move):** Translates action indices determined by the DQN into specific moves within the game, ensuring that actions align with the game's rules and dynamics.

```

1 def act(self, state):
2     if np.random.rand() <= self.epsilon:
3         return random.randrange(self.action_size)
4     state = torch.FloatTensor(state).unsqueeze(0)
5     act_values = self.model(state)
6     return np.argmax(act_values.detach().numpy()[0]) #
    returns action
7
8 def index_to_move(self, action_index, game):
9     # Converte un indice di azione in una mossa specifica,
10    assicurandosi che la mossa sia valida.
11    possible_moves = get_possible_moves(game.get_board(), self
12        .player_index)

```

```

11     if action_index >= 0 and action_index < len(possible_moves
12 ):
13         return possible_moves[action_index]
14     else:
15         return random.choice(possible_moves)

```

## Model Saving and Loading Progress Tracking

```

1 def save_model(self):
2     torch.save(self.model.state_dict(), "dql_model.pth")
3
4 def save_progress(self):
5     progress = {
6         "epsilon": self.epsilon,
7         "memory": self.memory,
8         # Aggiungi altri attributi se necessario
9     }
10    with open("dql_progress.pickle", 'wb') as f:
11        pickle.dump(progress, f)
12
13 def load_model(self):
14    self.model.load_state_dict(torch.load("dql_model.pth"))
15    self.model.eval() # Imposta il modello in modalita' di
16                      # valutazione
17
18 def load_progress(self):
19    with open("dql_progress.pickle", 'rb') as f:
20        progress = pickle.load(f)
21        self.epsilon = progress["epsilon"]
22        self.memory = progress["memory"]

```

**Move Making (make\_move):** Integrates the DQN's action selection and the game's rules to choose and execute moves during gameplay. This function represents the agent's interaction with the game environment, applying its learned strategies to compete.

**Reward Calculation (calculate\_reward):** Assesses the outcomes of actions to assign rewards, guiding the agent's learning by providing feedback on the effectiveness of its decisions relative to the game's outcome.

```

1 def make_move(self, game):
2     state = normalize_state_simple(game.get_board(), game.
3     get_current_player())
4     action_index = self.act(state)
5     chosen_move = self.index_to_move(action_index, game)
6     (x, y), direction = chosen_move
7     swapped_position = (y, x)
8     swapped_move = (swapped_position, direction)
9     return swapped_move
10
11 def calculate_reward(self, game):
12     winner = game.check_winner()
13     if winner == self.player_index:
14         return 1

```

```

14     elif winner == -2:
15         return 0
16     elif winner != -1:
17         return -1
18     else:
19         return 0
20

```

**Step Simulation (step):** A hypothetical function that would simulate the effect of an action on the game state, returning the new state, reward, and whether the game has ended. This functionality is essential for evaluating potential moves and their outcomes within the learning process.

```

1  def step(board, action, player_index):
2
3      # Esegue un passo nel gioco data un'azione e restituisce
4      # il nuovo stato, la ricompensa e se c'e' un vincitore
5
6      new_board = apply_move(deepcopy(board), action,
7                              player_index)
8      winner = check_winner(new_board)
9
10     if winner == player_index:
11         reward = 2
12         done = True
13     elif winner == -2:
14         reward = 0
15         done = True
16     elif winner != -1:
17         reward = -1
18         done = True
19     else:
20         reward = 0
21         done = False
22
23     return new_board, reward, done

```

### 5.3 Results

Despite the struggle and initially unconvincing results, after all night long training sessions, seasoned with parameter tuning, the best agent to use for training agents with RL is Minimax with Depth 1 or 3, due to the very long processing times.

```
1      -Testing DQL versus RandomPlayer
2
3      Test completato dopo 100 partite
4      Vittorie: 78 (78.00%)
5      Pareggi: 0 (0.00%)
6      Sconfitte: 22 (22.00%)
7
8
9      -Training DQL versus QLearningPlayer
10
11     Test completato dopo 100 partite
12     Vittorie: 74 (74.00%)
13     Pareggi: 0 (0.00%)
14     Sconfitte: 26 (26.00%)
15
```