

CS7IS2 - Artificial Intelligence Assignment 3 Report

Kaaviya Paranji Ramkumar

April 2025

1 Introduction

In this assignment, I've explored the fascinating world of game-playing AI algorithms by implementing and comparing different approaches for two classic board games: Tic Tac Toe and Connect 4. The central goal was to understand how different AI strategies perform across games of varying complexity and to analyze the trade-offs between different approaches.

Algorithms implemented

- **Minimax** is a decision-making algorithm for turn-based games that recursively evaluates all possible future game states. It assumes both players play optimally - the maximizing player (the AI) tries to maximize its score, while the minimizing player (opponent) tries to minimize the AI's score. By exploring the game tree, Minimax can select the move that leads to the best possible outcome assuming optimal play.
- **Alpha-beta pruning** is an optimization of the Minimax algorithm that significantly reduces the number of nodes that need to be evaluated. It works by eliminating branches of the search tree that cannot possibly influence the final decision. This allows for deeper searches in the same amount of time without affecting the final choice of move.
- **Q-Learning** is a model-free reinforcement learning algorithm that learns an optimal policy by estimating the value of state-action pairs. The algorithm builds a Q-table mapping states and actions to expected rewards, learning through trial and error. Over time, it discovers which actions lead to higher rewards in each state, without needing to know the rules of the game in advance.

2 Implementation

2.1 Code organization and object-oriented design

The experimental framework for this assignment was carefully planned and implemented using object-oriented programming principles to ensure modularity, re-usability, and clean separation of concerns. The code structure follows software engineering best practices, making it easy to extend and maintain.

The code base for maze generation and solving is currently maintained at
<https://github.com/prkaaviya/TCD-Artificial-Intelligence-A3>

The project is organized in a hierarchical structure with clearly defined components like:

```
.
|-- agents/
|   |-- default_agent.py      # Default opponent (better than random)
|   |-- minimax_agent.py     # Minimax implementation
|   |-- qlearning_agent.py    # Q-learning implementation
|-- evaluation/
|   |-- evaluator.py          # Game evaluator functions to evaluate agents
|-- games/
```

```
| |-- connect4.py          # Connect 4 implementation
| |-- tic_tac_toe.py       # Tic Tac Toe implementation
|-- main.py               # Main entry point
|-- results/              # Directory for storing results
'-- utils/
    '-- evaluation.py
```

2.2 Game environment

2.2.1 About Tic Tac Toe

For Tic Tac Toe, a straightforward game environment was implemented using a 3×3 NumPy array to represent the board state. The implementation includes:

- **State representation:** The board is represented as a 3×3 matrix where:
 - 0 represents an empty cell
 - 1 represents player 1's mark (X)
 - 2 represents player 2's mark (O)
- **Action space:** Actions are defined as (row, column) tuples, indicating where a player places their mark. With a 3×3 grid, there are initially 9 possible actions, decreasing as the game progresses.
- **Check valid action:** The game tracks available actions and verifies that a selected action corresponds to an empty cell.
- **Detect terminal state:** After each move, the game checks if a terminal state has been reached by:
 - Examining all rows, columns, and both diagonals for three matching marks in a row
 - Checking if the board is full (this represents draw condition)

The simplicity of Tic Tac Toe made the implementation straightforward, with minimal complexity in state management and win condition checking.

2.2.2 About Connect 4

Connect 4 required a more complex implementation to handle its larger board and the gravity mechanic:

- **State representation:** The board is represented as a 6×7 matrix (6 rows, 7 columns) using the same numerical encoding for players:
 - 0 for empty cells
 - 1 for player 1's discs (Red)
 - 2 for player 2's discs (Yellow)
- **Action space:** Unlike Tic Tac Toe, actions in Connect 4 are simply column indices (0-6) where a player drops their disc. The game handles the "gravity" by determining the lowest empty row in the chosen column.
- **Column tracking:** To efficiently determine valid moves and handle the gravity mechanic, the game maintains an array tracking the current height of each column.
- **Detect terminal state:** After each move, the game checks for terminal states by:
 - Examining only the lines affected by the most recent move (horizontal, vertical, and both diagonals)
 - This targeted checking is much more efficient than examining the entire board
- **Optimization:** Rather than checking the entire board after each move, the implementation only examines lines that could potentially contain a winning sequence including the newly placed piece.

2.2.3 Challenges in state and action representation

1. Size of state space

- Tic Tac Toe: Small and manageable ($\sim 5,478$ possible states)
- Connect 4: Enormous (~ 4.5 trillion possible states)

2. Action representation

- Tic Tac Toe: Direct (row, column) coordinates
- Connect 4: Column selection only, with row determined by gravity

3. Win condition complexity

- Tic Tac Toe: Checking 8 possible lines (3 rows, 3 columns, 2 diagonals)
- Connect 4: Checking many more potential lines, optimized by only examining those affected by the latest move

These implementation differences directly impact the performance of the various AI algorithms, particularly when scaling from the simpler Tic Tac Toe to the more complex Connect 4 environment.

2.3 Implementations of different agents

2.3.1 Default agent

For the baseline opponent, a Default Agent was implemented that exhibits simple but strategic behavior. This agent follows a straightforward priority-based decision-making process:

- **Strategy:** The agent makes decisions based on three simple rules:
 1. If there exists a move that would result in an immediate win, take it
 2. If the opponent has a move that would result in their immediate win, block it
 3. Otherwise, select a random valid move
- **Implementation:** For each possible action, the agent simulates the result and checks if it leads to a terminal state. The implementation creates a copy of the game state to test moves without modifying the actual game.
- **Adaptability:** The same agent works for both Tic Tac Toe and Connect 4 without game-specific modifications, adapting to the different action spaces automatically.

This agent satisfies the requirement of being "better than random" while still being simple enough to test our more sophisticated algorithms against.

2.3.2 Minimax agent

The Minimax algorithm forms the foundation of traditional game-playing AI by recursively exploring the game tree to find optimal moves.

- **Core algorithm:** Minimax works by:
 - Recursively evaluating all possible future game states
 - Assuming that the opponent will play optimally
 - Alternating between maximizing the player's score and minimizing the opponent's score
 - Selecting the move that leads to the best possible outcome
- **Implementation details**

- Separate functions for maximizing and minimizing players
- Terminal state evaluation returning +1 for wins, -1 for losses, and 0 for draws
- Performance tracking for nodes visited and execution time
- **Depth limitation:** For Connect 4, depth-limited search was included to mitigate issues with scalability due to large state space:
 - Configurable maximum depth parameter to limit recursion
 - Evaluation function for non-terminal states when the depth limit is reached
- **Heuristic functions:** For depth-limited search in Connect 4, heuristic functions was designed to evaluate non-terminal states:
 - Evaluating all possible winning lines (horizontal, vertical, diagonal)
 - Assigning higher scores to positions with more of the player's pieces in a row
 - Considering both offensive potential (forming own lines) and defensive concerns (blocking opponent)

2.3.3 Minimax with Alpha-Beta pruning agent

Alpha-Beta pruning enhances the Minimax algorithm by eliminating branches that cannot influence the final decision.

- **Optimization principle:** Alpha-Beta works by:
 - Maintaining two values: alpha (best already found for maximizer) and beta (best already found for minimizer)
 - Pruning branches where it's proven that better options exist elsewhere in the tree
 - Producing identical decisions to regular Minimax, just more efficiently
- **Implementation details**
 - Modified versions of the max and min functions that track alpha and beta values
 - Early termination of branch exploration when a cutoff is detected
 - Same depth limitation and heuristic functions as the regular Minimax
- **Efficiency gains:** The implementation includes metrics to measure:
 - Number of nodes explored compared to regular Minimax
 - Execution time improvements
 - Decision quality validation (ensuring identical choices to regular Minimax)

2.3.4 Q-Learning agent

The Q-Learning agent takes a completely different approach, learning from experience rather than searching through future states.

- **Learning approach**
 - Builds a Q-table mapping state-action pairs to expected rewards
 - Updates Q-values using the formula: $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$
 - Where α is the learning rate, γ is the discount factor, r is the reward, s is the current state, a is the action, and s' is the resulting state
- **State representation**

- For Tic Tac Toe: Direct representation of the 3×3 board as tuples
- For Connect 4: Compact representation tracking pieces in each column to reduce state space
- **Exploration-exploitation balance**
 - Epsilon-greedy strategy with a decaying exploration rate
 - Initial exploration rate of 0.3, decaying by factor of 0.995 after each episode
 - Minimum exploration rate of 0.01 to ensure some exploration continues
- **Implementation features**
 - Separate training and exploitation modes
 - Persistent Q-table storage to save and load learned policies
 - Reward shaping to encourage winning (+1), avoid losing (-1), and slightly prefer draws over uncertain states (+0.2)

2.4 Evaluation Framework

To systematically compare the performance of the different agents, an evaluation system was implemented to completely automate the multiple playthroughs of the games between all agents:

- **Game evaluator structure**
 - Round-robin format where each agent plays against all others
 - Symmetric games (each agent plays as both first and second player)
 - Configurable number of games per evaluation
- **Performance metrics**
 - Win/loss/draw rates for each agent
 - Average game length (number of moves)
 - Execution time per move and per game
 - Nodes visited (for Minimax agents)
- **Visualization and analysis tools**
 - Bar charts comparing agent performance
 - CSV export for detailed statistical analysis
 - Detailed printout of each game results
- **Experimental controls**
 - Same initial conditions for all games
 - Independent agent instances for each game to prevent learning during evaluation
 - Special depth testing mode to analyze the effect of search depth on performance

3 Experimental Results and Analysis

In this section, I present the results of my experimental evaluation, examining the performance of each algorithm across both Tic Tac Toe and Connect 4 games. The analysis focuses on effectiveness (win rates), efficiency (nodes visited and execution time), and the relative strengths and weaknesses of each approach.

3.1 Setup and methodology

To ensure a comprehensive and fair comparison, I conducted evaluation playthrough with multiple games between all possible agent combinations. For each game environment:

- 100 games were played for each playthrough in Tic Tac Toe and Connect 4 (1,200 games total each game)
- Each algorithm played both as first player and second player
- All algorithms faced off against each other and the default agent
- Performance metrics (win rates, execution times, nodes visited) were recorded for each game

The Q-Learning agent was trained separately before evaluation, with 5,000 episodes for Tic Tac Toe and 50,000 episodes for Connect 4. This disparity in training episodes is because of much larger state space needed in Connect 4.

For Connect 4, I first confirmed the computational infeasibility of unlimited-depth Minimax by attempting to run it without depth limitation. As expected, the algorithm could not complete even a single move in a reasonable timeframe (30 minutes), demonstrating the necessity of depth-limited search for larger game spaces.

The commands used to execute these experiments were:

```
# For Connect 4 depth testing
python main.py --game connect4 --mode depth_test --num_games 3

# For Q-Learning training
python main.py --game tictactoe --mode train --train_episodes 5000
python main.py --game connect4 --mode train --train_episodes 50000

# For Tic Tac Toe games
python main.py --game tictactoe --mode evaluator --num_games 100 --load_qtable

# For Connect 4 games
python main.py --game connect4 --mode evaluator --num_games 1000 --load_qtable
```

3.2 Alpha-Beta pruning efficiency

The application of alpha-beta pruning during my experiments provided remarkable improvements which enhanced computational performance. Alpha-beta pruning leads to significant performance improvements across both game types although they manifest in varying degrees (Figure 1).

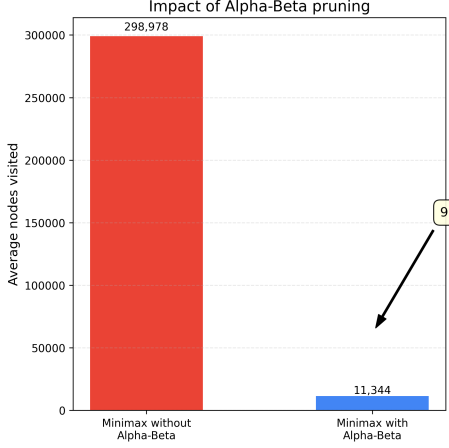
The Tic Tac Toe utilizes normal Minimax to analyze 298,978 nodes but alpha-beta pruning decreases this number to only 11,344 nodes which represents an impressive 96.2% reduction in nodes. The advantage in execution time directly mirrors this reduction resulting in 2.21 second game execution time reduction to just 0.08 seconds that represents a total 96.4% improvement.

Alpha-beta pruning reduced the node count for Connect 4 from 1,976 to 1,044 and cut execution time per game to 0.27 seconds from 0.54 seconds based on a depth limit of 3.

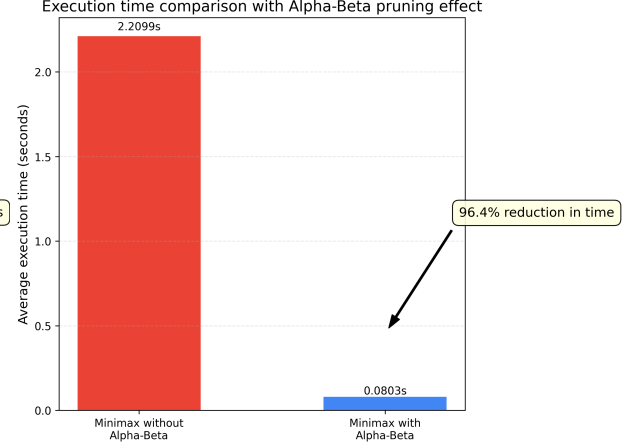
The shallow node exploration in Connect 4 detection accounts for the lower percentage of performance enhancement because of its size-bound constraints. Setting the depth limit to 3 restricts the search space severely thus creating fewer chances for optimizing the search. Even though Connect 4 improvements show fewer percentage changes the absolute gains from pruning strategies still stand substantial.

3.3 Algorithm performance analysis

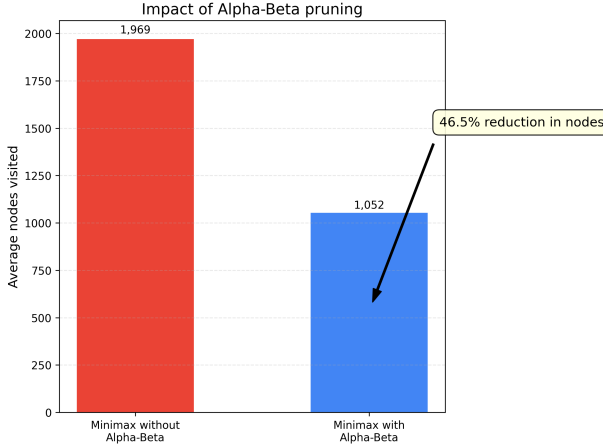
Table 1 presents the overall performance metrics for all agents in both games. Several interesting patterns emerge from these results.



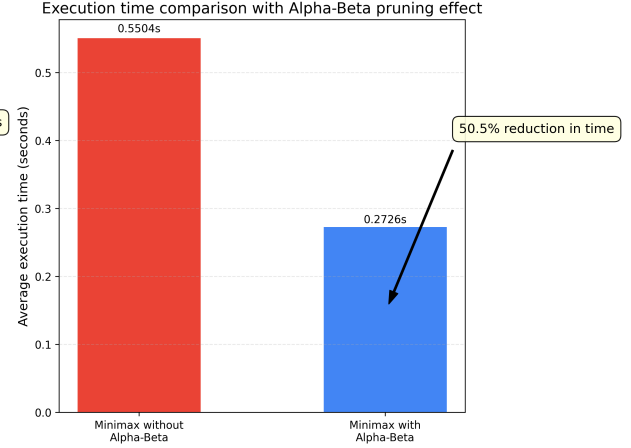
(a) Nodes visited in Tic Tac Toe



(b) Execution time in Tic Tac Toe



(c) Nodes visited in Connect 4



(d) Execution time in Connect 4

Figure 1: Impact of alpha-beta pruning on computational efficiency in both games

3.3.1 Tic Tac Toe performance

The Default Agent achieved 53.17% wins overall in Tic Tac Toe while both Minimax variants maintained similar levels at 50% and the Q-Learning Agent won at 43%.

Tic Tac Toe presents an evident advantage to the first player according to analysis of game results. All games played by both Minimax variants resulted either in a victory or draw when the agents played as the first player.

Minimax techniques maintain equivalent decision quality through identical move selection and virtually the same winning performance though they execute at highly differing computational speeds. Alpha-beta pruning maintains the decision quality when the algorithm performs calculations and removes unnecessary processing steps from the solution finding process.

The Q-Learning Agent demonstrates a winning percentage of 43% through rapid learning during 5,000 training episodes although it takes less than 0.0001 seconds to process each move. Minimax delivers optimal moves through complete searches yet Q-learning achieves effective strategies solely from experiences demonstrating reinforcement learning potential even when training amounts are restricted.

The Default Agent demonstrates strong performance in Tic Tac Toe due to the game's easy nature which allows simple strategies like winning directly and stopping opponent wins to prove successful.

(a) Overall performance metrics for Tic Tac Toe (1200 games in all)

Algorithm	Win rate	Draw rate	Avg. moves	Avg. time (s)
Default Agent	53.17%	3.67%	5.95	0.0002
Minimax without Alpha-Beta	50.17%	0.17%	5.12	2.2264
Minimax with Alpha-Beta	49.83%	0.17%	5.12	0.0809
Q-Learning Agent	43.00%	3.67%	6.10	<0.0001

(b) Overall performance metrics for Connect 4 (1200 games in all with depth limit = 3)

Algorithm	Win rate	Draw rate	Avg. moves	Avg. time (s)
Default Agent	61.50%	0.00%	13.74	0.0007
Minimax without Alpha-Beta	49.83%	0.00%	10.65	0.5504
Minimax with Alpha-Beta	49.83%	0.00%	10.72	0.2726
Q-Learning Agent	38.33%	0.00%	17.39	0.0001

Table 1: Performance comparison of all algorithms across both game environments

3.3.2 Connect 4 performance

The algorithms in Connect 4 generate wider performance changes among each other compared to Tic Tac Toe. The Default Agent surpasses the Tic Tac Toe performance by winning 61.50% of its matches but the Q-Learning Agent obtains a win rate of only 38.33%.

The growing difference between the agents stands out because Connect 4 introduces enhanced levels of complexity to the board game. The simple heuristic approach used by the Default Agent generates better performance because it can overcome opponents with shallow search limits and opponents with incomplete learning pipelines.

Minimax agents perform inadequately with a depth limit of 3 because they cannot maintain complete visibility to the end of the game as Minimax without depth limitations does in Tic Tac Toe. The heuristic evaluation function assists depth frontier decision-making but fails to replace the deficiencies brought by restricted lookahead capabilities.

The Q-Learning Agent achieves fewer optimal policy solutions in Connect 4 rather than Tic Tac Toe mainly due to the game’s broader state space that requires extra training attempts to develop a complete policy. Despite receiving 50,000 training episodes which amounted to ten times the number provided to the Tic Tac Toe agent the Q-Learning agent recorded a lower win rate suggesting additional training sessions or better state analysis would help larger board games.

3.4 Scalability analysis

The performance gap between Tic Tac Toe algorithms and Connect 4 algorithms demonstrates essential scalability factors to consider. Minimax algorithms proved to be computationally impossible when executed without depth constraints since they failed to finish one move within thirty minutes because of their massive state space.

The application of depth-limited search in combination with heuristic evaluation function became necessary for playing Connect 4. A depth limit of 3 enabled Minimax execution yet required it to sacrifice its ability for perfect playing. Due to its depth limitation the algorithm keeps a look ahead of just three moves that may not let it detect winning strategies that need extended move sequences.

The Q-Learning agent encountered scalability problems when applied to the Connect 4 game. The complexity of the game required a compacted state representation to ensure learning success although the program performed 50,000 training episodes still failed to grasp complete tactical mastery.

Testing different depth limits for Connect 4 revealed:

- Depth 1: 44 nodes, 0.015s - Very fast but strategically weak
- Depth 2: 449 nodes, 0.121s - Reasonable compromise for fast play

- Depth 3: 1,111 nodes, 0.290s - Chosen for our experiments as a good balance
- Depth 4: 5,866 nodes, 1.430s - Better play but significantly slower
- Depth 5: 11,618 nodes, 2.861s - Strong play but prohibitively slow for evaluation
- Unlimited: Computationally infeasible (>30 minutes per move)

These scalability observations highlight a fundamental trade-off: as game complexity increases, algorithms must sacrifice either completeness (Minimax with depth limits) or optimality (Q-Learning with limited training) to remain computationally feasible.

4 Conclusion

I gained intriguing understanding about how strategic AI methods handle difficulties through this assignment of game-playing algorithms. The use of Alpha-beta pruning sharply decreased computational workload by more than 90 percent during Tic Tac Toe gameplay maintaining absolute decision making to demonstrate the value of intelligent approach choice over excess work persistence. The three algorithms demonstrated their unique capabilities during Connect 4 game play so that Minimax delivered optimal results even though it proved difficult to scale up while Q-learning provided adaptability despite its minimal execution time once trained and the basic Default Agent surprisingly outperformed its more advanced counterparts. The conclusions established a fundamental AI dilemma that exists throughout computer systems where performance speed must be weighed against the caliber of choices made in situations where resources are restricted and problems remain complex.

Appendices

A Default agent implementation code

```
1 class DefaultAgent:
2     """
3     A default opponent for game playing.
4
5     This agent is better than random but still simple:
6     1. If there's a winning move, take it
7     2. If opponent has a winning move, block it
8     3. Otherwise, make a random valid move
9     """
10
11 def __init__(self, player_id: int):
12     """
13     Initialize the agent.
14
15     Args:
16         player_id: The ID of the player (1 or 2)
17     """
18     self.player_id = player_id
19     self.opponent_id = 3 - player_id # if player_id is 1, opponent is 2 and vice versa
20     self.name = "Default Agent"
21
22 def select_action(self, game, training: bool = False) -> Union[Tuple[int, int], int]:
23     """
24     Select an action based on the current game state.
25
26     Args:
27         game: The game object (TicTacToe or Connect4)
28
29     Returns:
30         action: For TicTacToe: (row, col) tuple, for Connect4: column index
31     """
32     valid_actions = game.get_valid_actions()
33     if not valid_actions:
34         raise ValueError("No valid actions available")
35
36     # check if there's a winning move
37     for action in valid_actions:
38         # make a copy of the game to simulate moves
39         game_copy = self._copy_game(game)
40
41         # simulate taking the action
42         state, _, done, info = game_copy.step(action)
43
44         # if the game is over and agent won, take this action
45         if done and info["winner"] == self.player_id:
46             return action
47
48     # check if the opponent has a winning move and block it
49     # simulate opponent's turn first
50     for action in valid_actions:
51         game_copy = self._copy_game(game)
52
53         # change current player to opponent to simulate their move
54         game_copy.current_player = self.opponent_id
55
56         # simulate opponent taking the action
57         state, _, done, info = game_copy.step(action)
58
59         # if the game would be over and opponent would win, block this action
60         if done and info["winner"] == self.opponent_id:
61             return action
62
```

```

63         # otherwise make a random valid move
64         return random.choice(valid_actions)
65
66     def _copy_game(self, game):
67         """Create a deep copy of the game to simulate moves."""
68         if hasattr(game, 'rows') and hasattr(game, 'cols'):
69             game_copy = type(game)()
70             game_copy.board = game.board.copy()
71             game_copy.current_player = game.current_player
72             game_copy.done = game.done
73             game_copy.winner = game.winner
74             game_copy.column_heights = game.column_heights.copy()
75             return game_copy
76         else:
77             game_copy = type(game)()
78             game_copy.board = game.board.copy()
79             game_copy.current_player = game.current_player
80             game_copy.done = game.done
81             game_copy.winner = game.winner
82             game_copy.valid_actions = game.valid_actions.copy() \
83                 if hasattr(game, 'valid_actions') else []
84             return game_copy

```

B Minimax agent implementation code

```

1 import numpy as np
2 import time
3 from typing import Tuple, List, Dict, Any, Union, Optional
4
5 class MinimaxAgent:
6     """
7     Minimax agent with optional alpha-beta pruning.
8
9     This agent uses the minimax algorithm to select the best action.
10    It can be configured to use alpha-beta pruning for efficiency.
11    For Connect4, a depth limit and heuristic evaluation function are used.
12    """
13
14    def __init__(self, player_id: int, use_alpha_beta: bool = True,
15                 max_depth: Optional[int] = None):
16        """
17        Initialize the Minimax agent.
18
19        Args:
20            player_id: The ID of the player (1 or 2)
21            use_alpha_beta: Whether to use alpha-beta pruning
22            max_depth: Maximum depth to search (None for unlimited)
23        """
24        self.player_id = player_id
25        self.opponent_id = 3 - player_id
26        self.use_alpha_beta = use_alpha_beta
27        self.max_depth = max_depth
28        self.name = f"Minimax {'with' if use_alpha_beta else 'without'} Alpha-Beta"
29        if max_depth:
30            self.name += f" (Depth {max_depth})"
31
32        self.nodes_visited = 0
33        self.execution_time = 0
34
35    def select_action(self, game, training: bool = False) -> Union[Tuple[int, int], int]:
36        """
37        Select the best action using minimax.
38
39        Args:
40            game: The game object
41
42        Returns:

```

```

43         The best action
44         """
45         self.nodes_visited = 0
46         start_time = time.time()
47
48         valid_actions = game.get_valid_actions()
49         if not valid_actions:
50             raise ValueError("No valid actions available")
51
52         # in case agent runs out of time/depth, play a fallback move
53         best_action = valid_actions[0]
54
55         if self.use_alpha_beta:
56             best_value = float('-inf')
57             alpha = float('-inf')
58             beta = float('inf')
59
60             for action in valid_actions:
61                 game_copy = self._copy_game(game)
62                 _, _, done, info = game_copy.step(action)
63
64                 # if this move wins immediately, take it
65                 if done and info["winner"] == self.player_id:
66                     self.execution_time = time.time() - start_time
67                     return action
68
69                 # otherwise evaluate the move
70                 if done:
71                     value = 0
72                 else:
73                     value = self._min_value(game_copy, 1, alpha, beta)
74
75                 if value > best_value:
76                     best_value = value
77                     best_action = action
78
79                 alpha = max(alpha, best_value)
80         else:
81             best_value = float('-inf')
82
83             for action in valid_actions:
84                 game_copy = self._copy_game(game)
85                 _, _, done, info = game_copy.step(action)
86
87                 if done and info["winner"] == self.player_id:
88                     self.execution_time = time.time() - start_time
89                     return action
90
91                 if done:
92                     value = 0
93                 else:
94                     value = self._min_value_no_pruning(game_copy, 1)
95
96                 if value > best_value:
97                     best_value = value
98                     best_action = action
99
100         self.execution_time = time.time() - start_time
101         return best_action
102
103     def _max_value(self, game, depth: int, alpha: float, beta: float) -> float:
104         """Maximizing player in minimax with alpha-beta pruning."""
105         self.nodes_visited += 1
106
107         # check if we have reached the maximum depth or a terminal state
108         if (self.max_depth is not None and depth >= self.max_depth) or game.is_terminal():
109             return self._evaluate(game)
110

```

```

111     value = float('-inf')
112     valid_actions = game.get_valid_actions()
113
114     for action in valid_actions:
115         game_copy = self._copy_game(game)
116         _, _, done, info = game_copy.step(action)
117
118         if done:
119             if info["winner"] == self.player_id:
120                 child_value = 1.0
121             elif info["winner"] == 0:
122                 child_value = 0
123             else:
124                 child_value = -1.0
125         else:
126             child_value = self._min_value(game_copy, depth + 1, alpha, beta)
127
128         value = max(value, child_value)
129
130         if value >= beta:
131             return value # return if beta cutoff
132
133         alpha = max(alpha, value)
134
135     return value
136
137 def _min_value(self, game, depth: int, alpha: float, beta: float) -> float:
138     """Minimizing player in minimax with alpha-beta pruning."""
139     self.nodes_visited += 1
140
141     # check if we have reached the maximum depth or a terminal state
142     if (self.max_depth is not None and depth >= self.max_depth) or game.is_terminal():
143         return self._evaluate(game)
144
145     value = float('inf')
146     valid_actions = game.get_valid_actions()
147
148     for action in valid_actions:
149         game_copy = self._copy_game(game)
150         _, _, done, info = game_copy.step(action)
151
152         if done:
153             if info["winner"] == self.opponent_id:
154                 child_value = -1.0
155             elif info["winner"] == 0:
156                 child_value = 0
157             else:
158                 child_value = 1.0
159         else:
160             child_value = self._max_value(game_copy, depth + 1, alpha, beta)
161
162         value = min(value, child_value)
163
164         if value <= alpha:
165             return value # return if alpha cutoff
166
167         beta = min(beta, value)
168
169     return value
170
171 def _max_value_no_pruning(self, game, depth: int) -> float:
172     """Maximizing player in minimax without alpha-beta pruning."""
173     self.nodes_visited += 1
174
175     # check if we have reached the maximum depth or a terminal state
176     if (self.max_depth is not None and depth >= self.max_depth) or game.is_terminal():
177         return self._evaluate(game)
178

```

```

179     value = float('-inf')
180     valid_actions = game.get_valid_actions()
181
182     for action in valid_actions:
183         game_copy = self._copy_game(game)
184         _, _, done, info = game_copy.step(action)
185
186         if done:
187             if info["winner"] == self.player_id:
188                 child_value = 1.0 # Win
189             elif info["winner"] == 0:
190                 child_value = 0 # Draw
191             else:
192                 child_value = -1.0 # Loss
193         else:
194             child_value = self._min_value_no_pruning(game_copy, depth + 1)
195
196         value = max(value, child_value)
197
198     return value
199
200 def _min_value_no_pruning(self, game, depth: int) -> float:
201     """Minimizing player in minimax without alpha-beta pruning."""
202     self.nodes_visited += 1
203
204     # check if we have reached the maximum depth or a terminal state
205     if (self.max_depth is not None and depth >= self.max_depth) or game.is_terminal():
206         return self._evaluate(game)
207
208     value = float('inf')
209     valid_actions = game.get_valid_actions()
210
211     for action in valid_actions:
212         game_copy = self._copy_game(game)
213         _, _, done, info = game_copy.step(action)
214
215         if done:
216             if info["winner"] == self.opponent_id:
217                 child_value = -1.0 # Loss
218             elif info["winner"] == 0:
219                 child_value = 0 # Draw
220             else:
221                 child_value = 1.0 # Win
222         else:
223             child_value = self._max_value_no_pruning(game_copy, depth + 1)
224
225         value = min(value, child_value)
226
227     return value
228
229 def _evaluate(self, game) -> float:
230     """
231     Evaluate the current game state.
232
233     For terminal states:
234     - Win: +1.0
235     - Loss: -1.0
236     - Draw: 0.0
237
238     For non-terminal states (when using depth-limited search):
239     - Heuristic evaluation based on potential winning lines
240     """
241     # if the game is over, return the actual outcome
242     if game.is_terminal():
243         winner = game.get_winner()
244         if winner == self.player_id:
245             return 1.0 # the agent won
246         elif winner == self.opponent_id:

```

```

247         return -1.0 # the agent lost
248     else:
249         return 0.0 # the game is a draw
250
251     # for depth-limited search, use a custom heuristic evaluation
252     if hasattr(game, 'rows') and hasattr(game, 'cols'):
253         # heuristic for Connect4
254         return self._evaluate_connect4(game)
255     else:
256         # heuristic for TicTacToe
257         return self._evaluate_tictactoe(game)
258
259 def _evaluate_connect4(self, game) -> float:
260     """
261     Heuristic evaluation for Connect4.
262     Counts potential winning lines for both players.
263     """
264     board = game.board
265     score = 0
266
267     # check horizontal, vertical, and both diagonals for wins
268     # give higher scores to positions with more of the player's pieces in a row
269
270     # for horizontal
271     for row in range(game.rows):
272         for col in range(game.cols - 3):
273             window = board[row, col:col+4]
274             score += self._evaluate_window(window)
275
276     # for vertical
277     for col in range(game.cols):
278         for row in range(game.rows - 3):
279             window = board[row:row+4, col]
280             score += self._evaluate_window(window)
281
282     # for diagonal (in positive slope)
283     for row in range(game.rows - 3):
284         for col in range(game.cols - 3):
285             window = [board[row+i, col+i] for i in range(4)]
286             score += self._evaluate_window(window)
287
288     # for diagonal (in negative slope)
289     for row in range(3, game.rows):
290         for col in range(game.cols - 3):
291             window = [board[row-i, col+i] for i in range(4)]
292             score += self._evaluate_window(window)
293
294     return score
295
296 def _evaluate_window(self, window) -> float:
297     """Evaluate a window of 4 positions."""
298     player_count = np.sum(window == self.player_id)
299     opponent_count = np.sum(window == self.opponent_id)
300     empty_count = np.sum(window == 0)
301
302     # if there's a mix of both players' pieces, this window isn't winnable
303     if player_count > 0 and opponent_count > 0:
304         return 0
305
306     # score based on how many of player's pieces are in the window
307     if player_count > 0:
308         if player_count == 3 and empty_count == 1:
309             return 0.8 # near to winning
310         elif player_count == 2 and empty_count == 2:
311             return 0.3
312         elif player_count == 1 and empty_count == 3:
313             return 0.1
314

```

```

315     # score based on how many of opponent's pieces are in the window
316     if opponent_count > 0:
317         if opponent_count == 3 and empty_count == 1:
318             return -0.8 # near to losing
319         elif opponent_count == 2 and empty_count == 2:
320             return -0.3
321         elif opponent_count == 1 and empty_count == 3:
322             return -0.1
323
324     return 0
325
326 def _evaluate_tictactoe(self, game) -> float:
327     """
328     Custom heuristic evaluation for TicTacToe.
329     Simpler than Connect4 since the state space is smaller.
330     """
331     board = game.board
332     score = 0
333
334     # check rows, columns, and diagonals for potential wins below
335
336     for row in range(3):
337         score += self._evaluate_line(board[row, :])
338
339     for col in range(3):
340         score += self._evaluate_line(board[:, col])
341
342     score += self._evaluate_line(np.array([board[0, 0], board[1, 1], board[2, 2]]))
343     score += self._evaluate_line(np.array([board[0, 2], board[1, 1], board[2, 0]]))
344
345     return score
346
347 def _evaluate_line(self, line) -> float:
348     """Evaluate a line of 3 positions for TicTacToe."""
349     player_count = np.sum(line == self.player_id)
350     opponent_count = np.sum(line == self.opponent_id)
351     empty_count = np.sum(line == 0)
352
353     # if both players have pieces in this line, it can't be won
354     if player_count > 0 and opponent_count > 0:
355         return 0
356
357     # score based on how many of player's pieces are in the line
358     if player_count > 0:
359         if player_count == 2 and empty_count == 1:
360             return 0.6 # near to winning
361         elif player_count == 1 and empty_count == 2:
362             return 0.2
363
364     # score based on how many of opponent's pieces are in the line
365     if opponent_count > 0:
366         if opponent_count == 2 and empty_count == 1:
367             return -0.6 # near to losing
368         elif opponent_count == 1 and empty_count == 2:
369             return -0.2
370
371     return 0
372
373 def _copy_game(self, game):
374     """Create a deep copy of the game to simulate moves."""
375     if hasattr(game, 'rows') and hasattr(game, 'cols'):
376         game_copy = type(game)()
377         game_copy.board = game.board.copy()
378         game_copy.current_player = game.current_player
379         game_copy.done = game.done
380         game_copy.winner = game.winner
381         game_copy.column_heights = game.column_heights.copy()
382     return game_copy

```



```

383         else:
384             game_copy = type(game)()
385             game_copy.board = game.board.copy()
386             game_copy.current_player = game.current_player
387             game_copy.done = game.done
388             game_copy.winner = game.winner
389             game_copy.valid_actions = game.valid_actions.copy() \
390                 if hasattr(game, 'valid_actions') else []
391             return game_copy

```

C Qlearning agent implementation code

```

1  import numpy as np
2  import random
3  import pickle
4  from typing import Tuple, Dict, List, Union, Any
5  import os
6
7  class QLearningAgent:
8      """
9      Q-learning agent for playing games.
10
11      This agent uses tabular Q-learning to learn a policy.
12      State representation is simplified to make learning tractable.
13      """
14
15     def __init__(self, player_id: int, learning_rate: float = 0.1, discount_factor: float =
16         0.9,
17         exploration_rate: float = 0.3, exploration_decay: float = 0.995,
18         min_exploration: float = 0.01,
19         load_qtable: bool = False, save_path: str = None):
20
21         """
22         Initialize the Q-learning agent.
23
24         Args:
25             player_id: The ID of the player (1 or 2)
26             learning_rate: Alpha - learning rate
27             discount_factor: Gamma - discount factor for future rewards
28             exploration_rate: Epsilon - exploration rate
29             exploration_decay: Factor to decay exploration rate after each episode
30             min_exploration: Minimum exploration rate
31             load_qtable: Whether to load an existing Q-table
32             save_path: Path to save/load Q-table
33         """
34         self.player_id = player_id
35         self.opponent_id = 3 - player_id
36         self.learning_rate = learning_rate
37         self.discount_factor = discount_factor
38         self.exploration_rate = exploration_rate
39         self.exploration_decay = exploration_decay
40         self.min_exploration = min_exploration
41         self.q_table = {}
42         self.name = "Q-Learning Agent"
43         self.save_path = save_path or f"q_table_player{player_id}.pkl"
44
45         # keep track of the current game for learning
46         self.current_state = None
47         self.current_action = None
48
49         if load_qtable and os.path.exists(self.save_path):
50             self.load_q_table()
51
52     def state_to_tuple(self, game) -> tuple:
53         """
54         Convert a game state to a tuple that can be used as a dictionary key.
55         Uses a simplified representation to keep the state space manageable.
56         """

```

```

54     Args:
55         game: The game object
56
57     Returns:
58         A tuple representation of the state
59     """
60     if hasattr(game, 'rows') and hasattr(game, 'cols'):
61         # use a more compact representation for Connect4
62         # we can't use the raw board as a key because it's too large
63         # instead, we can create a tuple of tuples where each inner tuple represents a
column
64         columns = []
65         for col in range(game.cols):
66             column_pieces = []
67             for row in range(game.rows):
68                 if game.board[row, col] != 0:
69                     # store (row, player_id) for each piece in the column
70                     column_pieces.append((row, int(game.board[row, col])))
71             columns.append(tuple(column_pieces))
72
73         # include the current player in the state
74         return (tuple(columns), game.current_player)
75     else:
76         # use a flattened board for TicTacToe
77         # convert board to tuple of tuples to hash it
78         board_tuple = tuple(tuple(row) for row in game.board)
79         return (board_tuple, game.current_player)
80
81 def get_q_value(self, state: tuple, action: Union[Tuple[int, int], int]) -> float:
82     """
83     Get the Q-value for a state-action pair.
84
85     Args:
86         state: The state tuple
87         action: The action
88
89     Returns:
90         The Q-value
91     """
92     if state not in self.q_table:
93         self.q_table[state] = {}
94
95     # convert action to a hashable type if it's not already
96     if isinstance(action, np.ndarray):
97         action = tuple(action)
98
99     if action not in self.q_table[state]:
100         self.q_table[state][action] = 0.0
101
102     return self.q_table[state][action]
103
104 def update_q_value(self, state: tuple, action: Union[Tuple[int, int], int],
105                   reward: float, next_state: tuple, done: bool) -> None:
106     """
107     Update the Q-value for a state-action pair.
108
109     Args:
110         state: The current state tuple
111         action: The action taken
112         reward: The reward received
113         next_state: The resulting state tuple
114         done: Whether the episode is done
115     """
116     # convert action to a hashable type if it's not already
117     if isinstance(action, np.ndarray):
118         action = tuple(action)
119
120     q_value = self.get_q_value(state, action)

```

```

121
122     # if the episode is done, there is no next state
123     # else, get max Q-value for next state
124     if done:
125         max_next_q = 0
126     else:
127         if next_state in self.q_table and self.q_table[next_state]:
128             max_next_q = max(self.q_table[next_state].values())
129         else:
130             max_next_q = 0
131
132     # update the Q-learnign formula
133     new_q_value = q_value + self.learning_rate * ( \
134         reward + self.discount_factor * max_next_q - q_value)
135
136     if state not in self.q_table:
137         self.q_table[state] = {}
138
139     self.q_table[state][action] = new_q_value
140
141 def select_action(self, game, training: bool = False) -> Union[Tuple[int, int], int]:
142     """
143     Select an action based on the current game state.
144
145     Args:
146         game: The game object
147         training: Whether the agent is in training mode
148
149     Returns:
150         The selected action
151     """
152     state = self.state_to_tuple(game)
153     valid_actions = game.get_valid_actions()
154
155     if not valid_actions:
156         raise ValueError("No valid actions available")
157
158     # save current state for learning update
159     if training:
160         self.current_state = state
161
162     # let agent explore with random action
163     if training and random.random() < self.exploration_rate:
164         action = random.choice(valid_actions)
165         if training:
166             self.current_action = action
167         return action
168
169     # let agent take best known action
170     q_values = {action: self.get_q_value(state, action) for action in valid_actions}
171
172     # find actions with the maximum Q-value
173     max_q = max(q_values.values()) if q_values else 0
174     best_actions = [action for action, q_value in q_values.items() if q_value == max_q]
175
176     # if multiple actions have the same Q-value, choose randomly among them
177     action = random.choice(best_actions) if best_actions else random.choice(
178         valid_actions)
179
180     if training:
181         self.current_action = action
182
183     return action
184
185 def learn(self, game, reward: float, done: bool) -> None:
186     """
187     Learn from the most recent action.

```

```

188     Args:
189         game: The game object after the action was taken
190         reward: The reward received
191         done: Whether the episode is done
192     """
193     if self.current_state is None or self.current_action is None:
194         return
195
196     next_state = self.state_to_tuple(game) if not done else None
197
198     self.update_q_value(self.current_state, self.current_action, reward, next_state,
199 done)
200
201     self.current_state = next_state
202     self.current_action = None
203
204     if done:
205         self.exploration_rate = max(self.min_exploration,
206                                     self.exploration_rate * self.exploration_decay)
207
208 def save_q_table(self, path: str = None) -> None:
209     """Save the Q-table to a file."""
210     save_path = path or self.save_path
211     with open(save_path, 'wb') as f:
212         pickle.dump(self.q_table, f)
213     print(f"Q-table saved to {save_path}")
214     print(f"Q-table size: {len(self.q_table)} states")
215
216 def load_q_table(self, path: str = None) -> None:
217     """Load the Q-table from a file."""
218     load_path = path or self.save_path
219     if os.path.exists(load_path):
220         with open(load_path, 'rb') as f:
221             self.q_table = pickle.load(f)
222         print(f"Q-table loaded from {load_path}")
223         print(f"Q-table size: {len(self.q_table)} states")
224     else:
225         print(f"No Q-table found at {load_path}")
226
227 def reset(self) -> None:
228     """Reset agent state for a new episode."""
229     self.current_state = None
230     self.current_action = None

```