

# MCTS

Robert Marc James-Stroud  
*School of Computer Science and Electronic Engineering*  
*University of Essex*  
Colchester, UK  
Email rj18801@essex.ac.uk

## Abstract

*To be written.*

## I. INTRODUCTION

Monte Carlo Tree Search (MCTS) has received much interest from researchers in multiple areas, including General Game Playing (GGP) [1], where it has been used to great success.

MCTS is an algorithm for returning a decision by creating a search tree of the domain, constructed by taking random samples of the search space. MCTS does not require domain knowledge to function although it may be helpful [2].

Research into Game AI is important, prior to MCTS databases were used to hold all possible game states. An example of database is presented by [6], a game of *American Checkers* in which there are  $3.9 \times 10^{13}$  entries, or states. Decision trees used in GGP and General Video Game Playing (GVGP) suffer from dimensionality.

A simple game, Noughts and Crosses has a branching factor of 4, and a full tree has 10 levels [3], compared with Chess which has a branching factor of 35 and rarely a full search tree [3]. Noughts and Crosses has unique states numbering in the thousands, it is clear that crafting rules for each different state is not feasible, especially when the branching factor in other games used in GGP/GVGP can be in the hundreds.

Firstly this paper will look at what MCTS is and previous work conducted in the field of Game AI, specifically MCTS. Secondly MCTS will be used in a Noughts and Crosses (also known as Os and Xs, OXO and Tic-Tac-Toe).

### A. Monte Carlo Tree Search

MCTS is a tree search algorithm that constructs an asymmetric search tree based on actions, the tree is a biased representation towards more promising areas of the search space [7].

Through self-play MCTS estimates the value of a node from that point until it reaches a terminal node [7]. A node is a state, and the action is the move made that results in that state.

Each node might be visited multiple times having its value adjusted accordingly. A node that looks promising on the first run through could look less promising as it is visited more times, or become less promising in relation to other branches.

- 1) *Selection* - Using a policy traverse the search tree. Starting at a root node, the current state, select child nodes with promising values until a leaf node is reached.

- 2) *Expansion* - Once *Selection* has selected a leaf node, expand the current node and select one of its children so long as the current node is not terminal.
- 3) *Simulation* - Using the node expanded in the *Expansion* phase and a policy if defined, playout the game. *Playout* in this context means until completion or a result is achieved.
- 4) *Backpropagation* - After *Simulation* is complete the search tree is traversed in reverse order, propagating up the tree the result, updating the value of the nodes, this continues until the current node is the root node. The new values are then used in the *Selection* process.

MCTS will continue to cycle through these phases until it is stopped, or an 'execution budget is reached' [8]. Three ways are presented to achieve this, each with their own detracting:

- Time - cull MCTS after a certain amount of time has lapsed - if stopped too quickly it may not evaluate nodes accurately.
- Depth Culling - Stop MCTS after hit reaches a certain depth limit in the tree - stopping with depth may stop the search from ever reaching a terminal node.
- Iteration Limit - Give it a predefined number of cycles, after which it will return a result. One cycle is all four phases of MCTS.

Depending on when or how MCTS is stopped will effect how good the estimation is of the action to take. Using depth culling or an iteration limit is useful when comparing algorithms across inconsistent hardware and programming languages.

### B. Previous Work

MCTS has received much interest from researchers, becoming somewhat of an umbrella term covering any implementation of MCTS. Browne *et al.* in [2] summarise the different MCTS implementations up to 2011.

MCTS has been applied in co-operative scenarios [7], Real-time games [2], Non-deterministic games [9] and numerous non-game applications [2].

Game AI research used to focus on two-player zero-sum games of perfect information with alternating turns [2], every two-player zero-sum game has a solution [3]. Two main ways of evaluating states in a minimax scenario is to use heuristics, or statistics. The statistics approach can utilise MCTS to run thousands of playouts to find the optimal strategy [3].

1) *Upper Confidence Bounds*: UCB1 is a bandit algorithm with a logarithmic regret, proposed by Auer *et al.* [4]. In a multi-armed bandit problem the ‘empirically best action’ should be taken as often as possible [4].

However less explored options should not be ignored in favour of exploitation. Disregarding exploration may leave better actions unexplored and thus increase regret.

UCB1 policy dictates that arm  $j$  maximises the average reward [4]:

$$UCB1 = \overline{X}_j + \sqrt{\frac{2 \ln n}{n_j}}$$

where  $\overline{X}_j$  is the average reward from arm  $j$ ,  $n_j$  is the number of times arm  $j$  was played and  $n$  is the total number of plays.

Reward  $\overline{X}_j$  encourages exploitation of high-reward choices [4] and  $\sqrt{\frac{2 \ln n}{n_j}}$  encourages exploration of less visited actions [2].

2) *Upper Confidence Bounds for Trees*: The result of combining UCB1 and MCTS is UCT. Upper Confidence Bounds for Trees (UCT) is a well known and popular algorithm in the MCTS family [2]. Kocsis and Szepesvari [5] proposed applying UCB1 to MCTS, using UCB1 as the tree policy [2].

UCT is efficient and simple [2] and a ‘promising candidate to address the exploration-exploitation dilemma in MCTS’ [2].

Treating the selection of a child node as a multi-armed bandit problem, using Monte Carlo simulations the ‘expected reward can be approximated’. Just as a node can be viewed as a multi-armed bandit problem so to can an action to be selected. Treated as an independent multi-armed bandit problem the child node  $j$  is selected to maximise [2]:

$$UCT = \overline{X}_j + 2C_p \sqrt{\frac{2 \ln n}{n_j}}$$

where  $C_p$  is a constant. The value of  $C_p$  can be adjusted to increase or decrease the amount of exploration MCTS does, however a default value of  $C_p = \frac{1}{\sqrt{2}}$  is commonly used [2].

## II. METHODOLOGY

Noughts and Crosses is a two-player zero sum game, as mentioned previously it will always have a solution. A search tree for UCT can be constructed using a state-action model, where the state is the game board and the action is the move made. This will result in a new state.

## III. EXPERIMENTS

## IV. RESULTS

## V. CONCLUSION

## VI. PLAN

## REFERENCES

[1] C. F. Sironi, J. Liu, D. Perez-Liebana, R. D. Gaina, I. Bravi, S. M. Lucas, M. H. M. Winands, “Self-Adaptive MCTS for General Video Game Playing”.

[2] C. Browne, E. Powley, D. Whitehouse, S. Lucas, P. I. Cowling, P. Rohlfshagen, S. Taverner, D. Perez, S. Samothrakis, S. Colton, “A Survey of Monte Carlo Tree Search Methods,” *IEEE Transactions on Computational Intelligence and AI in Games*, Vol. 4, 2012.

[3] R. A. Bartle “Game Theory,” *CE810 Game Design Lecture 5*, 6 November 2018.

[4] P. Auer, N. Cesa-Bianchi, P. Fischer, “Finite-Time Analysis of the Multiarmed Bandit Problem,” *Machine Learning*, Vol. 47, 2002.

[5] L. Kocsis, C. Szepesvari, “Bandit Based Monte-Carlo Planning”

[6] C. Browne, E. Powley, D. Whitehouse, S. Lucas, P. I. Cowling, P. Rohlfshagen, S. Taverner, D. Perez, S. Samothrakis, S. Colton, “EvoMCTS: A Scalable Approach for General Game Learning,” *IEEE Transactions on Computational Intelligence and AI in Games*, Vol. 6, 2014.

[7] P. R. Williams, J. Walton-Rivers, D. Perez-Liebana, S. M. Lucas, “Monte Carlo Tree Search Applied to Co-operative Problems,” *2015 7th Computer Science and Electronic Engineering Conference (CEECE)*.

[8] R. D. Gaina, S. M. Lucas, D. Perez-Liebana, “Rolling Horizon Evolution Enhancements in General Video Game Playing,” *IEEE Conference on Computational Intelligence and Games*, 2017.

[9] P. I. Cowling, E. J. Powley, D. Whitehouse, “Information Set Monte Carlo Tree Search” *IEEE Transactions on Computational Intelligence and AI in Games*, Vol. 4, 2012.