

Monte Carlo Tree Search in AlphaZero

Junyeong Park

Hanyang University
Department of Computer Science

November 7, 2020

1 Introduction

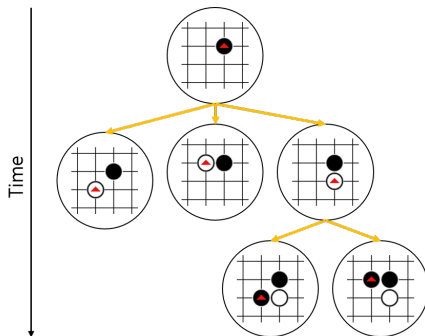
- Tree Structure
- Tree Search

2 AlphaZero

- Monte Carlo Simulation

How to represent the game

- We can represent the game with **tree** structure.
- Root represents current state.
- Depth is chronological order.



How to represent the game

Each node stores a set of information and a set of statistics.

- s_t is state at timestep t .
- a_t is action at timestep t .
- $N(s, a)$ denote that the total number of simulations.
- $W(s, a)$ denote that the number of wins.
- and so on (ex. $\text{Var}[V]$, $\sigma[V]$)

Tree Search

Our goal is to select the best move at state s_t .

→ We need to define the value of action.

How can represent the value of action?

- We have $W(s, a)$ and $N(s, a)$.
- So winning rate can be calculated by averaging the results of simulations.

$$\text{winrate}(s, a) = \frac{W(s, a)}{N(s, a)}$$

- winning rate is also used for state-action value.

Tree Search Algorithm - Mini-max

Definition

Mini-max is an algorithm that find way to **minimize** **maximum** loss.

```
function MINIMAX(node, depth, maximizingPlayer)
  if depth = 0 or node is a terminal node then
    return the heuristic value of node
  if maximizingPlayer then
    value  $\leftarrow -\infty$ 
    for all child of node do
      value  $\leftarrow \max(\text{value}, \text{Minimax}(\text{child}, \text{depth}-1, \text{false}))$ 
    return value
  else
    value  $\leftarrow +\infty$ 
    for all child of node do
      value  $\leftarrow \min(\text{value}, \text{Minimax}(\text{child}, \text{depth}-1, \text{true}))$ 
    return value
```

Tree Search Algorithm - Alpha-beta pruning

Definition

Alpha-beta pruning is a algorithm that seeks to **decrease the number of nodes** that are evaluated by the mini-max algorithm in its search tree.

- We don't need to search a whole tree.
- If one node is uninterested, do not search more.

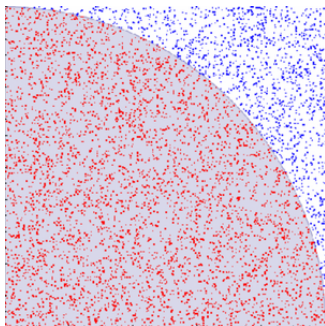
Limit of Mini-max and Alpha-beta pruning

- They need to investigate all nodes of the tree.
- However, a lots of environments' search space is huge.
(ex) Go, Atari games, etc.
- Needs a lots of times and memories.

Monte Carlo Method

Definition

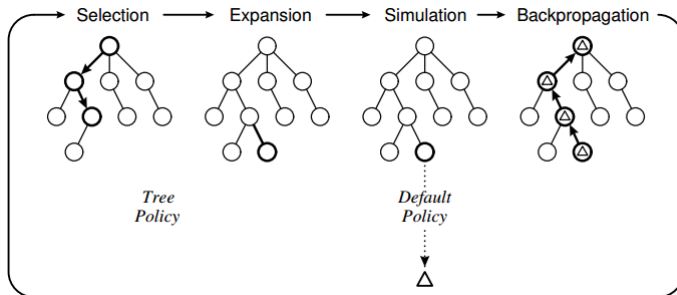
Monte carlo method is an algorithm that obtain numerical value with repeated random sampling.



Monte Carlo Tree Search

Definition

Monte carlo tree search is a tree search algorithm that adapts monte carlo method.



Monte Carlo Tree Search

■ *Selection*

Start from root R and select successive child nodes until a leaf node L is reached.

■ *Expansion*

In L , if the game has not ended yet, create one (or more) child nodes and choose node C from them.

■ *Simulation*

Complete one random playout from node C .

■ *Backpropagation*

Update statistics from C to R .

$$\begin{aligned}N(s, a) &\leftarrow N(s, a) + 1 \\W(s, a) &\leftarrow W(s, a) + z\end{aligned}$$

AlphaZero

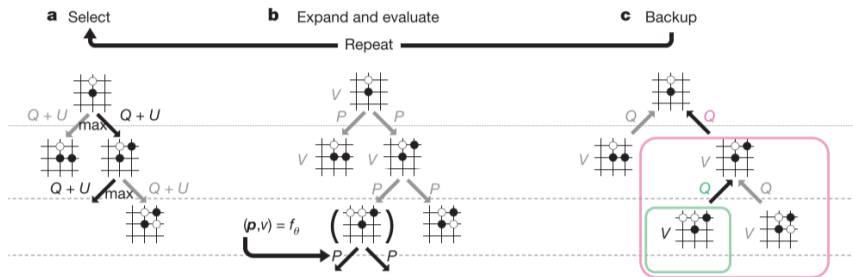
Definition

AlphaZero is reinforcement learning architecture to solve board games (Go, Chess, Shogi).

- It has neural network, $f_{\theta}(s) = (\mathbf{p}, v)$.
- Policy head produce move probabilities (\mathbf{p}).
: reduce the breadth of tree.
- Value head produce winning probability (v).
: reduce the depth of tree.

MCTS in AlphaZero

AlphaZero use modified monte carlo tree search.



MCTS in AlphaZero

Selection

- Use PUCT-algorithm.

- First, calculate values following

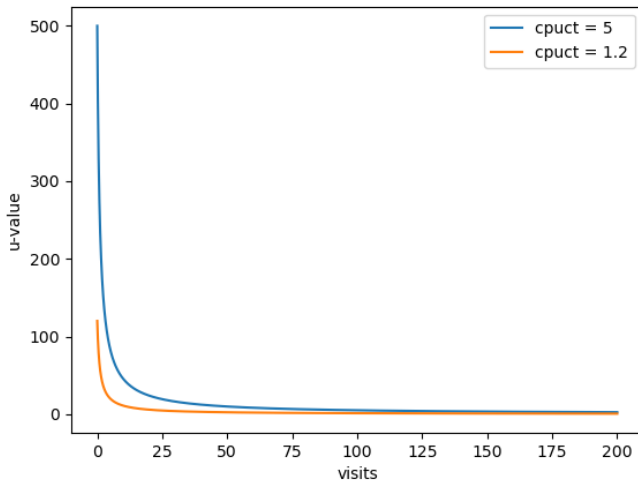
$$Q(s, a) = \frac{W(s, a)}{N(s, a)}$$
$$u(s, a) = c_{\text{puct}} \times P(s, a) \times \frac{\sqrt{\sum_b N(s, b)}}{N(s, a) + 1}$$

- Then, select action according to

$$a_t = \operatorname{argmax}_a (Q(s, a) + u(s, a))$$

until reach to leaf node s_L .

MCTS in AlphaZero



MCTS in AlphaZero

- The higher c_{puct} , the greater the breath of tree.
- The lower c_{puct} , the greater the depth of tree.
- AlphaZero apply heuristic for effective search.
(except training)

$$\text{PUCT}(s) = c_{\text{init}} + \log \left(\frac{\sum_b N(s,b) + c_{\text{base}+1}}{c_{\text{base}}} \right)$$

MCTS in AlphaZero

Expand and Evaluate

- Add the leaf node s_L to neural network evaluation queue.
- Expand the leaf node and each edge (s_L, a) is initialized to

$$N(s_L, a) = 0, W(s_L, a) = 0, P(s_L, a) = p_a$$

MCTS in AlphaZero

Backup (backpropagation)

- Update edge statistics from edge to root.

- $N(s, a) \leftarrow N(s, a) + 1$

- $W(s, a) \leftarrow W(s, a) + v$