

Advanced Machine Learning: Assignment 3

Seung-Hoon Na

December 05, 2017

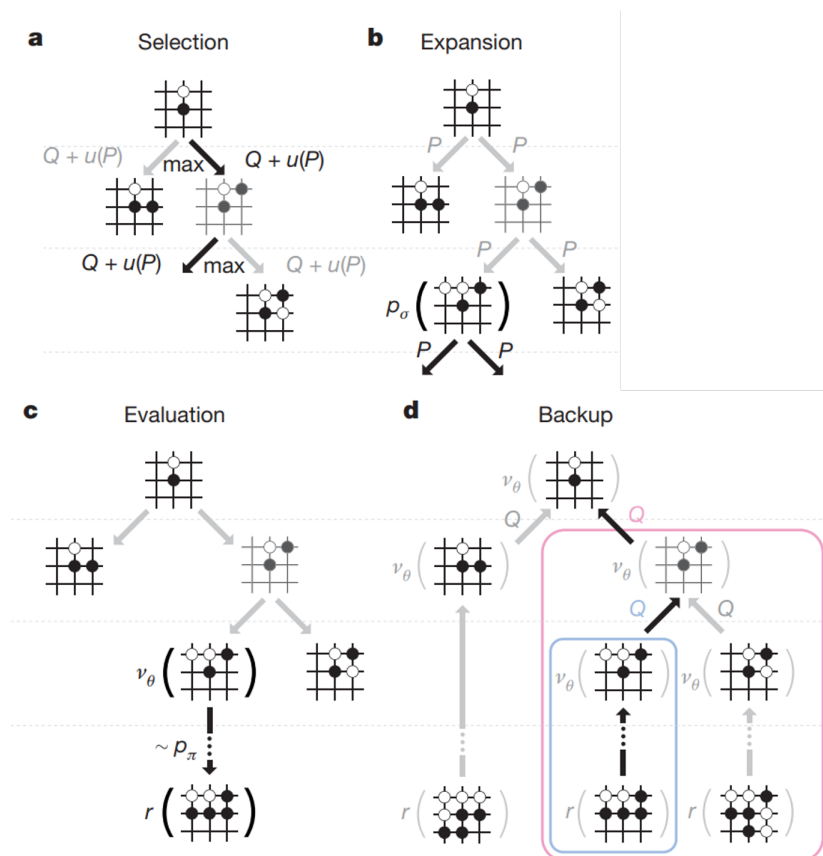
1 Monte Carlo Tree Search

Monte-Carlo tree search (MCTS) consists of the four steps – Selection, expansion, evaluation, and backup.

Read the following paper and summarize each step of MCTS and how MCTS is applied for selecting an action in a board game.

<http://www.cs.utexas.edu/~pstone/Courses/394Rspring13/resources/mcrave.pdf>

For reference, the following is the description of MCTS presented in the DeepMind's AlphaGo paper¹.



¹<https://www.nature.com/articles/nature16961>

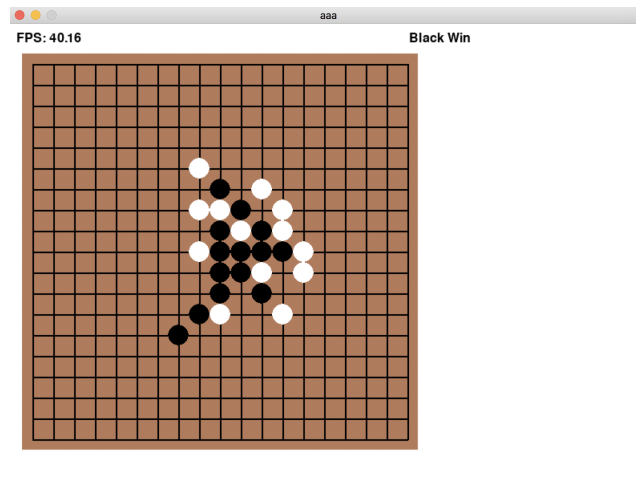
2 Self-play reinforcement learning for Gomoku

In this problem, you should implement an actor-critic model based on MCTS for mastering Gomoku.

Once a Gomoku AI is trained, you need to provide an additional Gomoku game code which can play with the trained Gomoku AI. To create a gomoku board game, you can refer to the following pygame gomoku (your own code is welcome):

<https://github.com/HackerSir/PygameTutorials/tree/master/Lesson04/Gomoku>

Here is the screen shot when playing the Gomoku game.



For concrete details, refer to the DeepMind's AlphaGo Zero paper.

<http://nature.com/articles/doi:10.1038/nature24270>

2.1 Actor-critic model for Gomoku

Actor-critic model consists of two networks – a value network $v(s, \mathbf{w})$ and a policy network $\pi(s, \theta)$. In Gomoku game, the state s is represented by the raw board representation – the black and white stones from the 19×19 (or 9×9) board (i.e., a binary image).

The self-play reinforcement learning consists of two parts – 1) self-play, 2) neural network training.

1) Self-play step

In the self-play step, n games of self-play are generated using MCTS. Here, the last value/policy networks – $v(s, \mathbf{w})$ and $\pi(s, \theta)$ – are used to guide the simulations of MCTS. As in the DeepMind's paper², more specifically, each edge (s, a) in the search tree stores a visit count $N(s, a)$ and an action value $Q(s, a)$.

Then, take the four steps of MCTS as follows:

1. Selection: For selecting a node in a tree in MCTS, the UCT algorithm is applied by treating each state of the search tree as multi-armed bandit, in which each action corresponds to an arm of the bandit. Given the

²Different from the DeepMind's paper, we use different setting for $U(s, a)$

UCT algorithm, each simulation starts from the root state and iteratively selects moves that maximize an upper confidence bound (UCB), until a leaf node s_L is encountered, as follows:

$$a^* = \operatorname{argmax}_a Q(s, a) + U(s, a).$$

where $U(s, a)$ is defined as:

$$U(s, a) = c \sqrt{\frac{\log N(s)}{1 + N(s, a)}}$$

where $N(s)$ is defined as $\sum_{a'} N(s, a')$ and c is a constant determining the level of exploration.

2. Expansion: The leaf node is expanded by using the last policy network $\pi(s, a, \theta)$.
3. Evaluation: The value of the leaf node s_L is evaluated in two ways – 1) by using the last value network $v(s, \mathbf{w})$, 2) by running a rollout to the end of the game using the fast rollout policy network $\pi(s, a, \theta)$, resulting in the final return z_L . These two evaluated values are merged using a mixing parameter λ , giving the combined value $V(s_L)$ as follows:

$$V(s_L) = (1 - \lambda)v(s_L, \mathbf{w}) + \lambda z_L$$

4. Backup: For each edge (s, a) traversed in the simulation, increment its visit count $N(s, a)$ and update action values to the mean valuation over all these simulations ³.

$$Q(s, a) = \frac{1}{N(s, a)} \sum_{s_L | s, a \rightarrow s_L} V(s_L)$$

where $s, a \rightarrow s_L$ indicates that a simulation eventually reached s_L from (s, a) .

2) Network training step

A single game is automatically generated by using N simulations for each MCTS (N is 1,600 in AlphaGo Zero setting). Along this way, suppose that we generate n games of self-play based on MCTS, and i -th simulation is given as $s_1^i, a_1^i, s_2^i, a_2^i, \dots, s_{T^i}^i, a_{T^i}^i$ with a final reward $z_{T^i}^i$. Let z_t^i be $\pm r(s_{T^i})$ from each player's perspective.

The simulations of all n games are converted to the list of minibatches (without using the game identity).

Without loss of generality, we represent each minibatch as $D_k = \{(s_t, a_t, z_t)\}$

First, given minibatch, a policy network is updated to minimize the following criterion:

³According to the DeepMind's AlphaGo 2016 paper, $Q(s, a)$ is equivalently rewritten as:

$$Q(s, a) = \frac{1}{N(s, a)} \sum_{i=1}^n 1(s, a, i) V(s_L^i)$$

where s_L^i the leaf node from the i -th simulation, and $1(s, a, i)$ indicates whether an edge (s, a) was traversed during the i th simulation.

$$J(\theta) = \mathbb{E}_{(s,a,z) \sim D_k} [\log \pi(s, a, \theta)(z - v(s, \mathbf{w}))]$$

Second, a value network is updated to minimize the following criterion:

$$J(\mathbf{w}) = \mathbb{E}_{(s,a,z) \sim D_k} [(z - v(s, \mathbf{w}))^2]$$

2.2 Actor-critic model for Gomoku: Implementation

Design a value network and a policy network based on deep convolutional neural networks for mastering Gomoku.

Implement the self-play reinforcement learning of a actor-critic model that train $v(S, \mathbf{w})$ and $\pi(s, a, \mathbf{w})$ for Gomoku using tensorflow (python code)

2.3 Actor-critic model: Evaluation

Find an existing Gomoku AI program and evaluate the winning rate of your trained Gomoku AI against the existing game.

<https://github.com/lingz/gomokuai>

Draw a performance curve which shows the average winning rate against the previous Gomoku AI game (python code)

2.4 Actor-critic model: Simulation of a game of self-play

Write an additional code that simulates a game of self-play generated by two of your Gomoku AIs (using python)

2.5 Actor-critic model: Demo

Write an Gomoku game code (computer vs. human) in which a user can play with your Gomoku AI (using python)

3 Self-play reinforcement learning for Tetris

Implement an actor-critic model based on MCTS for mastering Tetris.

As in assignment 2, you need to revise the following pygame MaTris codes:

<https://github.com/SmartViking/MaTris>

Compare a MCTS-based model with the previously implemented DQN and AC models for Tetris.