

Review on playing GO

1. Supervised learning of policy networks

Goal: predict what a master would do given a board state

method

$$s \rightarrow \sigma \rightarrow P_\sigma(a|s) \quad (a, s) = \text{training data}$$

gradient ascent

2. Reinforcement learning of policy networks

Goal: make moves that win games

method: start from prior network. Play self games

sample $i \rightarrow z_i$ outcome = ± 1

$$\Delta \sigma = \alpha \frac{1}{m} \sum_{i=1}^m \left[\sum_{t=1}^{T_i} \nabla \log P_\sigma(a_i^t | s_0^t) \right] z_i$$

Policy Gradient Theorem

→ first order step to maximize expected reward

MCTS

Reinforcement learning of value networks

Predicts ^{value} win or lose, quizzily!

- goal: make it very fast
does not need to be super accurate

method: use network trained to play games

goal: predict game value

gradient descent

both players use
network

$$s \rightarrow \begin{bmatrix} \theta \end{bmatrix} \rightarrow v_{\theta}(s)$$

$$\min E(z - v_{\theta})^2$$

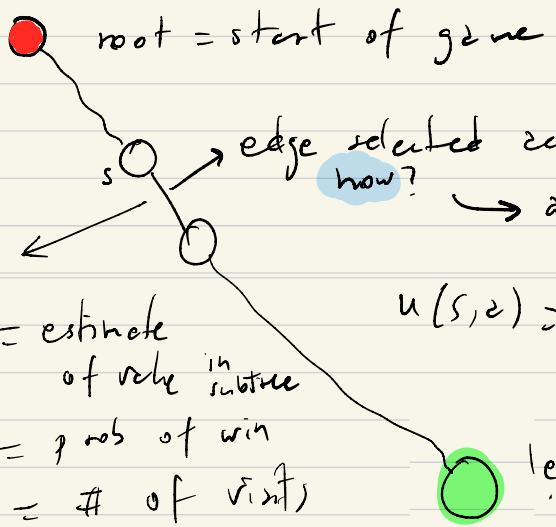
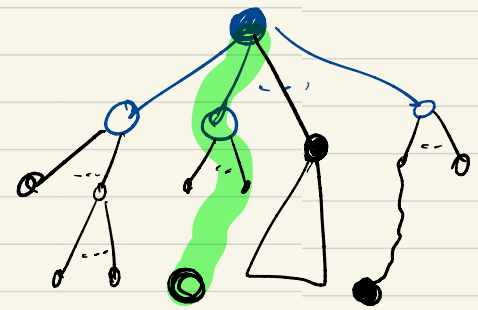
$$E(z_t | s_t = s, a_t, \dots, T \text{ VP})$$

$$\Delta_{\theta} \propto \nabla_{\theta}(v) \text{ MSE} (z - v_{\theta})$$

skipping value network

Monte-Carlo Tree Search (MCTS)

- Balance between exploration exploitation



$$a = \text{argmax} \{Q(s, a) + u(s, a)\}$$

$Q(s, a)$ = estimate of value in subtree

$$u(s, a) = \frac{P(s, a)}{1 + N(s, a)}$$

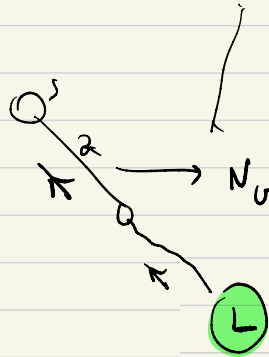
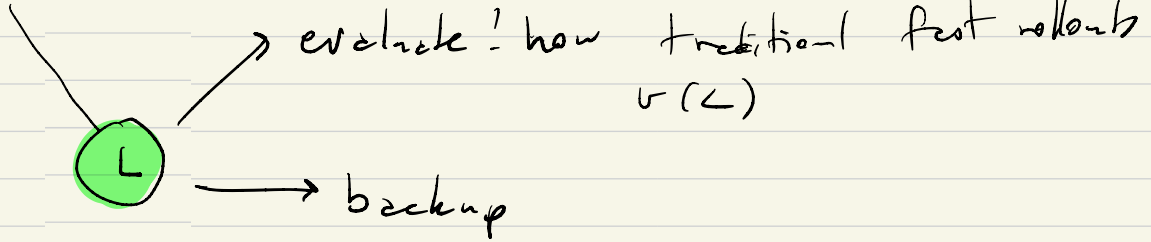
$P(s, a)$ = prob of win

$N(s, a)$ = # of visits

- leaf of current tree
 - evaluate
 - backup
 - expand

Generic

MCTS



$N_v(s, \alpha) = N_v(s, \alpha) + 1$

$Q(s, \alpha) = \frac{1}{N_v(s, \alpha)} \sum_{i=1}^m v_i$

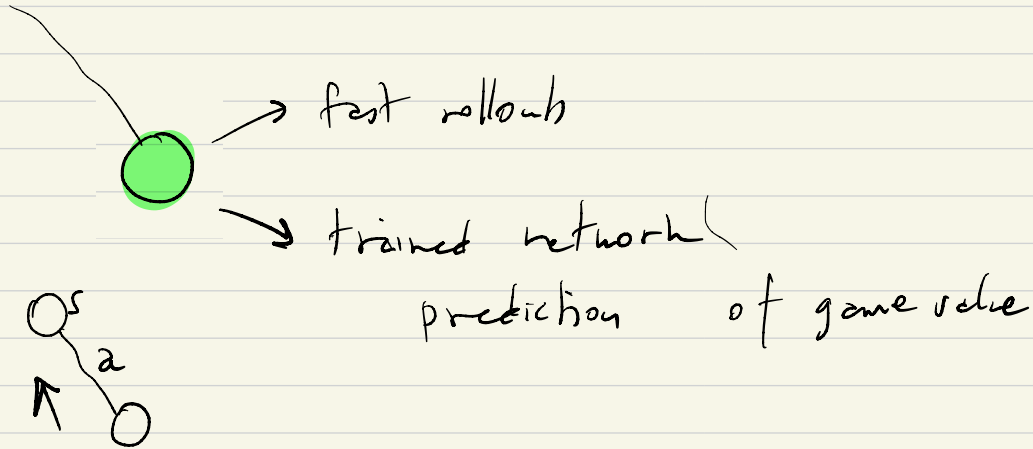
↓
define

Expand

L

when visit count is "large enough"

Alpha GO



$$Q(s, a) = \lambda \frac{W_r(s, a)}{N_r(s, a)} + (1 - \lambda) \frac{W_o(s, a)}{N_o(s, a)}$$

Expand

40 threads 48 CPUs 8 GPUs \rightarrow for $\sigma_2(L)$ leaf evaluation

40 threads 1202 CPUs 176 GPUs

" methods require several orders of magnitude more computation than traditional heuristics)

- Alpha 60 "many times stronger than any other previous program"
- Defeated (2015) Fan Hui 5-0