

Understanding AlphaGo

Machine Learning papers

Bitdefender ML Team
eburceanu@bitdefender.com

June 6, 2016

Intro

NN - Concepts

AlphaGo

Conclusions

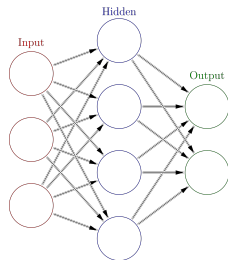
- ▶ brief NN introduction
- ▶ what are the components of AlphaGo (Nature magazine paper, January 2016)
- ▶ how do they link
- ▶ speculate about mistakes in games 3, 4

- ▶ NN purpose: given some data, approximate a good function of the input
 - ▶ $y_{nn} = output_{nn} = f_{W_f}(g_{W_g}(h_{W_h}(...(input))))$
 - ▶ $y_{correct} = output_{correct} = \text{known from data}$
 - ▶ define how different is y_{nn} from the $y_{correct}$

- ▶ NN purpose: given some data, approximate a good function of the input

- ▶ $y_{nn} = output_{nn} = f_{W_f}(g_{W_g}(h_{W_h}(...(input))))$
- ▶ $y_{correct} = output_{correct} = \text{known from data}$
- ▶ define how different is y_{nn} from the $y_{correct}$

- ▶ 150 vs 300
- ▶ cat vs dog
- ▶ $(y_{nn} - y_{correct})^2$
- ▶ small vs big penalties



- ▶ minimize cost = "difference" between y_{nn} and $y_{correct}$
- ▶ find proper parameters $f = f_{W_f}(X) = W_f * X + b_f$
- ▶ each time we make a mistake, find what caused it and make small adjustments in the network (fancy name: back-propagation)
- ▶ **Q:** How W influences the *Cost*?

- ▶ minimize cost = "difference" between y_{nn} and $y_{correct}$
- ▶ find proper parameters $f = f_{W_f}(X) = W_f * X + b_f$
- ▶ each time we make a mistake, find what caused it and make small adjustments in the network (fancy name: back-propagation)
- ▶ **Q:** How W influences the *Cost*?
 - ▶ iteratively, $W = W - \alpha * \nabla_W Cost$ (Taylor)
 - ▶ end: W will have a value for which *Cost* is minimum
 - ▶ **Q:** How else?

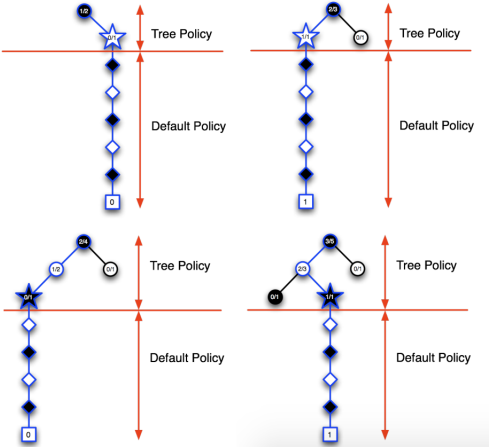
- ▶ perfect information game and zero-sum game
 - ▶ V_{optim} exists (state value, under optim play by all players)
- ▶ solving it:
- ▶ exhaustive
 - ▶ $branches^{depth}$ sequences of moves
 - ▶ chess: 35^80
 - ▶ go: 250^{150} (3^{361} table states)
- ▶ search tree
 - ▶ truncate (alpha-beta pruning)
 - ▶ sampling (Monte Carlo Tree Search - MCTS)
 - ▶ **Q**: Pro and cons?

- ▶ **Q:** Who read the paper?

- ▶ **Q:** Who read the paper?
- ▶ "Mastering the game of Go with deep neural networks and tree search"
- ▶ David Silver - first author, maybe the best RL...human expert? PhD in AI on Go
- ▶ in: dataset of expert moves (KGS)
- ▶ out: first program that defeats a professional player
- ▶ components:
 - ▶ Monte Carlo Tree Search (MCTS)
 - ▶ NN for $value(state) : V_{\theta}(s)$
 - ▶ NN for "strategy" (policy), $policy(action|state) : p_{\sigma}, p_{\pi}, p_{\rho}$

- ▶ $p_{\sigma}(action|state)$ policy
 - ▶ learns to do expert human moves
- ▶ $p_{\pi}(action|state)$ policy
 - ▶ learns to do **fast** expert human moves
- ▶ $V_{\theta}(state)$ value function
 - ▶ evaluates a state
- ▶ tree search algorithm
- ▶ **Q**: To do what?

Monte Carlo Tree Search I



- ☆ New node in the tree
- Node stored in the tree
- ◇ State visited but not stored
- Terminal outcome
- Current simulation
- Previous simulation

- ▶ simulate and find the best move
- ▶ statistics
 - ▶ $P(s, a)$ - p_σ prior probability
 - ▶ $N_v(s, a), N_r(s, a)$ - count
 - ▶ $W_v(s; a), W_r(s, a)$ - additive score for all simulations ($V_\theta(s)$ and $r(s)$)
 - ▶ $Q(s, a)$ - mixed score
 - ▶ for all legal actions. **Q**: How?
- ▶ +asynchronous, +PUCT (exploration factor)

► Steps

1. Selection

- statistics $t = \operatorname{argmax}_a (Q(s_t, a) + u(s_t, a))$ (best edge)
- $u(s, a) = c_{\text{puct}} P(s, a) \frac{\sqrt{\sum_b N_r(s, b)}}{1 + N_r(s, a)}$
- until no more expanded states (L - we need to evaluate it)

2. Evaluation

- $V_\theta(\text{state}_L)$ or
- reward of rollout with $p_\pi(a|\text{state}_L)$, for both players

3. Backup

- $N_r + 1, N_v + 1, W_v + V_\theta(\text{state}_L), W_r + \text{reward}$, async, n_{VL}
- $Q(s, a) = (1 - \lambda) \frac{W_v(s, a)}{N_v(s, a)} + \lambda * \frac{W_r(s, a)}{N_r(s, a)}$

4. Expansion

- $N_r(s, a) > n_{thr}$, init all to 0

- run multiple times to build the tree

- ▶ $p_{\sigma}(\text{action}|\text{state})$ policy
 - ▶ supervised learning (55%, SotA 44%, eval in 3 ms)
 - ▶ expert dataset (30 mil, in: Go board, out: move = action)
 - ▶ 48 planes - 19x19, locally preprocessed Go features
 - ▶ 13 layers, conv 5, 3 + ReLU + softmax legal moves
 - ▶ used as prior for $P(s, a)$ in MCTS ($u(s,a)$ component)
 - ▶ $\delta p = \frac{\partial \log(p(a|s))}{\partial p}$
- ▶ $p_{\pi}(\text{action}|\text{state})$ policy
 - ▶ supervised, feature-engineered input
 - ▶ conv + relu + softmax (24%)
 - ▶ simple and very fast (eval in 2 μs)
 - ▶ used for rollout

- ▶ $p_\rho(\text{action}|\text{state})$ policy
 - ▶ same structure, init with p_σ
 - ▶ 80% vs p_σ
 - ▶ play against itself (reinforcement learning, ± 1 reward)
 - ▶ not exactly itself, but sampling from previous versions
 - ▶ $\delta p = \frac{\partial \log(p(a|s))}{\partial p} * r(s_T)$
- ▶ $V_\theta(\text{state})$ value function
 - ▶ expected value of reward
 - ▶ supervised learning on p_ρ (playing history)
 - ▶ sample 1 state per game, 30 mil (**Q**: Why?)
 - ▶ similar architecture with p_σ (without softmax)
 - ▶ regression, RL, unbiased estimator of the reward
 - ▶ $\delta \theta = \frac{\partial \log(V_\theta(s))}{\partial \theta} * (z - V_\theta(s))$
 - ▶ average for all symmetries (8)

- ▶ **Q:** Why it is not maximizing the margin? (game 3)
- ▶ **Q:** Time management?
- ▶ **Q:** Why did it figure out so late that it is losing? (game 4)
- ▶ **Conclusions**
 - ▶ great things are simple
 - ▶ lots of engineering challenges (not addressed today)
 - ▶ insights and experience in the field (David Silver, Aja Huang)

Other Questions?



- ▶ **AlphaGo paper**, payed from Nature, just request it from the ML team, <https://deepmind.com/alpha-go.html>
- ▶ **Reinforcement Learning course**, <http://www0.cs.ucl.ac.uk/staff/d.silver/web/Teaching.html>
- ▶ **Neural Nets course**, <http://cs231n.stanford.edu/>
- ▶ **Reinforcement Learning Book**, Algorithms for Reinforcement Learning, Csaba Szepesvari
- ▶ **MCTS survey**, <http://www.cameronius.com/cv/mcts-survey-master.pdf>