**AlphaGo (2014-2017)**   DeepMind

Step 1: Supervised learning for convolutional <u>deep</u> neural network

3 weeks

use database from games
of expert players

— matched 55% of time

+ smaller (faster) 25% of time

13 layers

input:
19×19×48
locations features

output: a move (19×19 + 1)

hand-coded
features

black
white
empty
# opp captured
# own captured
liberties
ladder capture
ladder escape

Step 2: reinforcement learning for convolutional <u>deep</u> neural network

1 day

beat SL network   80% of time

Step 3: reinforcement learning for value network

using data from step 2 network
plays itself 30M times
samples 1 pos per game

+1 black win
0
-1 white win

Step 4: MCTS

default: use fast network from step 1
initialize new nodes' values using value net from step 3

tree policy:   $q(s,a)$ + $c \, P(s,a)$ $\dfrac{\sqrt{\text{\# times parent visited}}}{1 + \text{\# times child visited}}$

exploit
observed

from larger
step 1 network

Elo   3144 → 3739 → 5185
2015      2016      2017
(Fan Hui)  (Lee Sedol)  (retired)

Δ Elo   400 → higher rated player has 90 + % chance of winning

Dec 2^J 2021 Page 1

**AlphaGo Zero**

input :   $19 \times 19 \times \underline{17}$   current pos + last 7 positions
                                              + turn (all $\begin{array}{l}1 = \text{black}\\0 = \text{white}\end{array}$)

output :   move $(19 \times 19 + 1)$ and value $[-1, +1]$

play against self   $\longrightarrow$   (state, action probabilities, winner)
                                                                 from MCTS

                                                                 $\downarrow$

retrain  $\xleftarrow{\text{batch of 2048}}$  DB of last
                                                                 500K games

play against      every
best so far  $\xleftarrow{}$  1000 iterations

win 55%
to become new
best; else reset

# Deep Q Learning

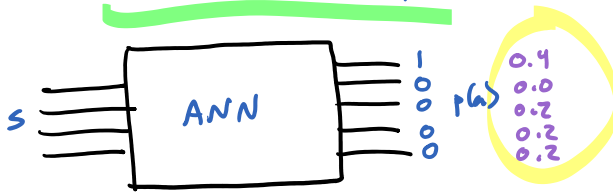$$s \equiv\!\!\!\equiv \boxed{ANN} \equiv\!\!\!\equiv \hat{Q}(s,a)$$

$$a \equiv\!\!\!\equiv$$

$$\Pi^*(s) = \underset{a}{argmax}\ q^*(s,a)$$

two networks: learning, target

play using learning observe $(s,a,s',r) \rightarrow$ add to replay database

compute error $\hat{Q}_\ell(s,a) \sim (r + \gamma \max_{a'} \hat{Q}_t(s',a))$

adjust learning accordingly $\qquad \hookrightarrow (s,a,s',r)$ sampled from replay database

periodically copy learning to target

# REINFORCE: learn policy directly

$$s \equiv\!\!\!\equiv \boxed{ANN} \equiv\!\!\!\equiv \begin{matrix}1\\0\\0\\0\\0\end{matrix}\ p(a) \begin{matrix}0.4\\0.0\\0.2\\0.2\\0.2\end{matrix}$$

objective: maximize $\quad E\left[\sum_{t=1}^{T} r_t\right]$

vector of weights in ANN

trajectory: $(s_1,a_1), (s_2,a_2), (s_3,a_3)\ldots$

objective to maximize $\quad J(\theta) = \int \Pi_\theta(\tau) \cdot r(\tau)\, d\tau$

prob ANN w/ weights $\theta$ yields $\tau$ — total reward over trajectory $\tau$

gradient ascent on weights $\quad \nabla_\theta J(\theta) = \int \boxed{\nabla_\theta \Pi_\theta(\tau) \cdot r(\tau)}\, d\tau$

how to adjust weights $\theta$ to increase $J(\theta)$ most?

$$\Pi_\theta(\tau) \cdot \nabla_\theta \log \Pi_\theta(\tau)$$
$$= \Pi_\theta(\tau) \cdot \frac{\nabla_\theta \Pi_\theta(\tau)}{\Pi_\theta(\tau)}$$

$$= \int \Pi_\theta(\tau) \nabla_\theta \log \Pi_\theta(\tau)\, r(\tau)\, d\tau$$

$$= E\left[\nabla_\theta \log \Pi_\theta(\tau) \cdot r(\tau)\right]$$

$$= E\left[\left(\sum_{t=1}^{T} \nabla_\theta \log \Pi_\theta(a_t|s_t)\right) \cdot \left(\sum_{t=1}^{T} r(s_t,a_t)\right)\right]$$

$\pi_\theta(\tau)$

$$= E\left[\left(\sum_{t=1}^{T} \nabla_\theta \log \pi_\theta(a_t | s_t)\right) \cdot \left(\sum_{t=1}^{T} r(s_t, a_t)\right)\right]$$

$\pi_\theta(\tau) = $ prob of trajectory $\tau$

$= \pi_\theta\left((s_1, a_1)(s_2, a_2) \cdots \right)$

$= p(s_1) \cdot \pi_\theta(a_1 | s_1) \cdot p(s_2 | s_1, a_1) \cdot \pi_\theta(a_2 | s_2) \cdots$

$= p(s_1) \cdot \prod_{i=1}^{T} \pi_\theta(a_t | s_t) \cdot p(s_{t+1} | s_t, a_t)$

$$\frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \left(\nabla_\theta \log \pi_\theta(a_{i,t} | s_{i,t})\right) \cdot \left(\sum_{t=1}^{T} r(s_{i,t}, a_{i,t})\right)$$

$\log \pi_\theta(\tau) = \log$

$= \log p(s_1) + \sum_{i=1}^{T} \left(\log \pi_\theta(a_t | s_t) + \log p(s_{t+1} | s_t, a_t)\right)$

$\nabla_\theta \log \pi_\theta(\tau) = \nabla_\theta \log p(s_1) + \sum_{i=1}^{T} \left(\log \pi_\theta(a_t | s_t) + \log p(s_{t+1} | s_t, a_t)\right)$

REINFORCE

                         (or N of them)

     get sample trajectory by running policy

     compute estimate of $\nabla_\theta J(\theta)$

     update $\theta$ : $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$