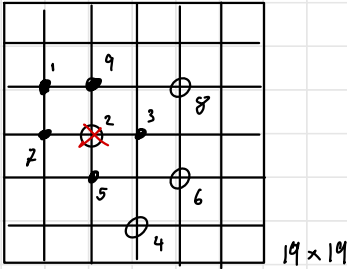


# Lecture 20



• Alpha Go :

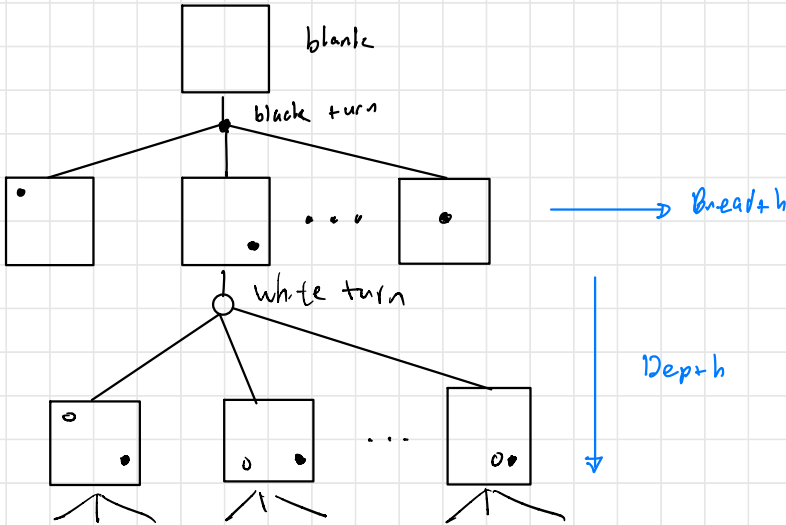
- Go: game



Player 1 : Black stone ●

Player 2 : White stone ○

• Game Tree :



• End Result :  $Z = \begin{cases} +1 & \text{black wins} \\ -1 & \text{white wins} \\ 0 & \text{draw} \end{cases}$

• The Game tree is extremely big

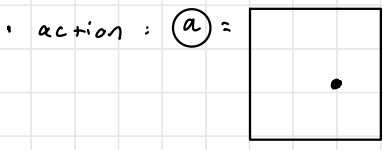
• state  $s = \left( \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 1 & & \\ \hline & 0 & 1 & 0 \\ \hline 0 & 1 & 0 & 0 \\ \hline \end{array} , \begin{array}{|c|c|c|c|} \hline 1 & 0 & 0 & 1 \\ \hline 0 & & & \\ \hline & 1 & 0 & \\ \hline 1 & 0 & 1 & 0 \\ \hline \end{array} , \begin{array}{|c|c|c|c|} \hline 1 & 1 & 1 & 1 \\ \hline 1 & 1 & 1 & 1 \\ \hline 1 & 1 & 1 & 1 \\ \hline \end{array} \right)$

$19 \times 19$

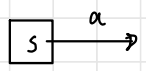
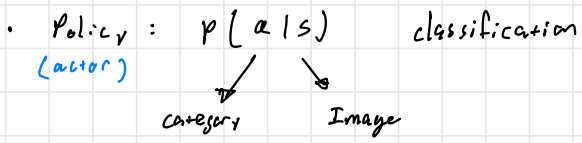
1: position of black stones

1: white stones

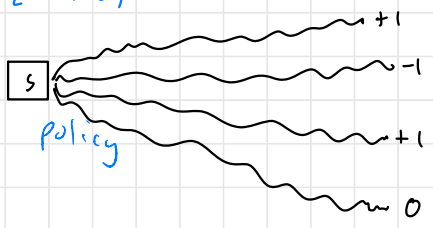
all 1's if black  
all 0's if white



$19 \times 19 + 1 = 362$   
↓  
pass

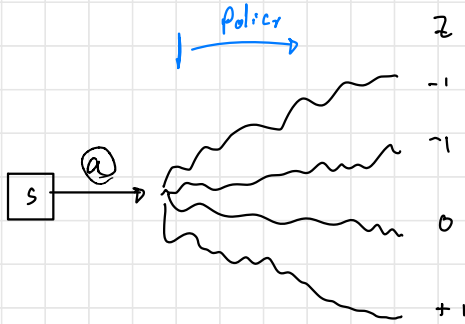


• Value:  $V_{\text{policy}}(s) = \mathbb{E}_{\text{policy}}(z|s)$  (critic)



## Action Value

- $Q_{\text{policy}}(s, a) = \mathbb{E}_{\text{policy}}(Z | s, a)$



- Optimal Policy :

$$\max_{\text{policy}} V_{\text{policy}}(s) = V^*(s)$$

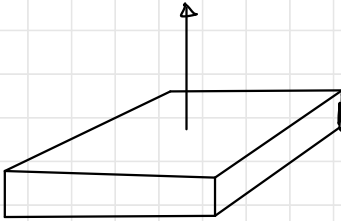
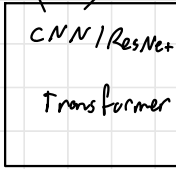
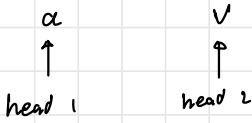
$$\max_{\text{policy}} Q_{\text{policy}}(s, a) = Q^*(s, a)$$

• Step 1 = Supervised Learning

Training Data

	1		
	2		
	⋮		
human →	i	$s_i$	$a_i$
	⋮		
	n		

Classification      regression



$19 \times 19 \times 3$

Supervised Policy Network  $P_{\sigma}(a|s)$

$$\max_{\sigma} \frac{1}{n} \sum_{i=1}^n \log P_{\sigma}(a_i | s_i)$$

MLE

$$\text{SGD: } \Delta \sigma \propto \frac{\partial}{\partial \sigma} \log P_{\sigma}(a|s)$$

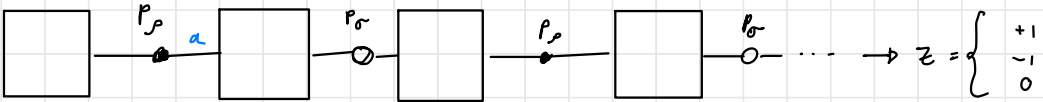
Imitation Learning  
Behavior Cloning

↑  
generated  
by Teacher

• Step 2: Reinforcement Learning:  $P_p(a|s)$

↓  
reinforcement

- initialize  $p = \sigma$
- improve  $p$  while fixing  $\sigma$

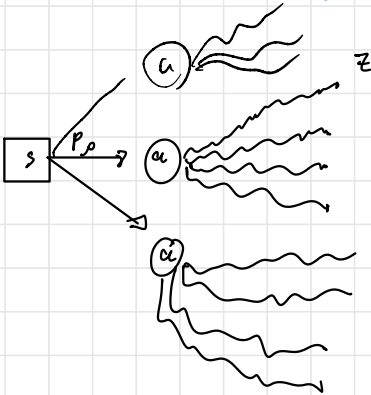


RL:

$$\Delta p \propto \frac{\partial}{\partial p} \log P_p(a|s) \stackrel{\text{self}}{\downarrow} z \begin{matrix} \text{consequence} \\ \text{reinforcement} \end{matrix}$$

- If  $a$  leads to a win  
increase probability

- If  $a$  leads to a loss  
decrease probability

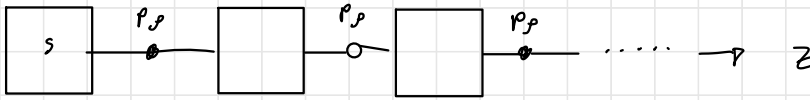


$$\max_p \mathbb{E}_p [z], \quad \mathbb{E}_p [z] = \sum_a z P_p(a|s)$$

$$\begin{aligned} \frac{\partial}{\partial \rho} \mathbb{E}_{\rho}(z) &= \sum_a z \left[ \frac{\partial}{\partial \rho} p_{\rho}(a|s) \right] \\ &= \sum_a z \left[ \frac{\partial}{\partial \rho} \log p_{\rho}(a|s) \right] \times p_{\rho}(a|s) \\ &= \mathbb{E}_{\rho} \left[ z \frac{\partial}{\partial \rho} \log p_{\rho}(a|s) \right] \end{aligned}$$

Policy Gradient

• Step 3:  $V_{\theta}(s)$



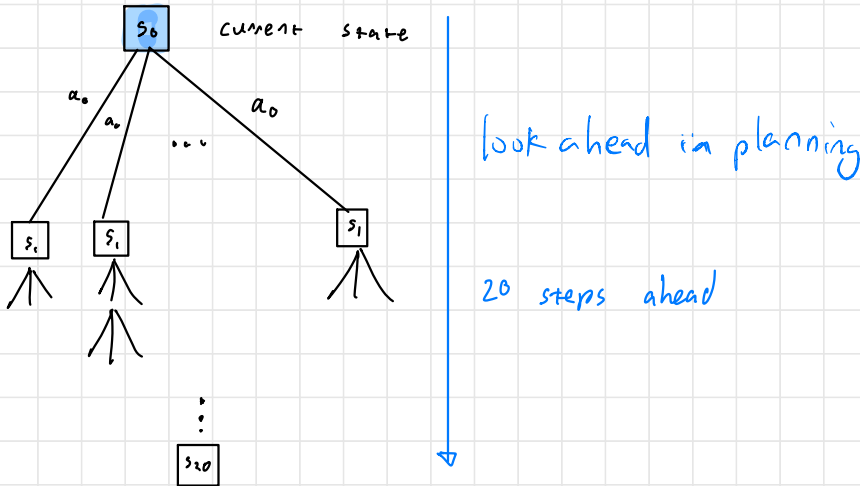
$$\Delta \theta \propto - \frac{\partial}{\partial \theta} (z - V_{\theta}(s))^2$$

- Learned  $V_{\theta}(s)$  more accurate for  $s$  close to the end

• Go to South Korea:

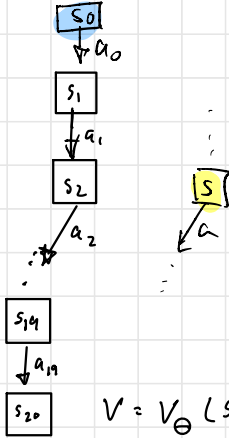
• Step 4: Monte Carlo Tree Search (MCTS)

- Real game v.s. Top Player





- Keep Sampling Branches:



$V = V_{\theta}(s_{20})$  more accurate than  $V_{\theta}(s_0)$

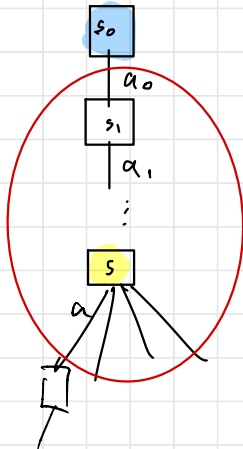
- For each  $(s, a)$  on branch

- visit #:  $N(s, a) \leftarrow N(s, a) + 1$
- total score:  $W(s, a) \leftarrow W(s, a) + V_{\theta}(s_{20})$

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

- How to sample branch?

real



Thinking (Imagination)

policy: reduce search breadth  
value: reduce search depth

$$\max_a \left[ \overset{\text{exploitation}}{Q(s, a)} + \overset{\text{exploration}}{c \cdot \frac{P_{\sigma}(a|s)}{N(s, a) + 1}} \right]$$

- After sampling many branches,  
choose  $a_0 : \max_{a_0} Q(s_0, a_0)$

or

$$\Pi(a_0 | s_0) \propto \mathcal{N}(s_0, a_0)^{\frac{1}{\beta}}$$

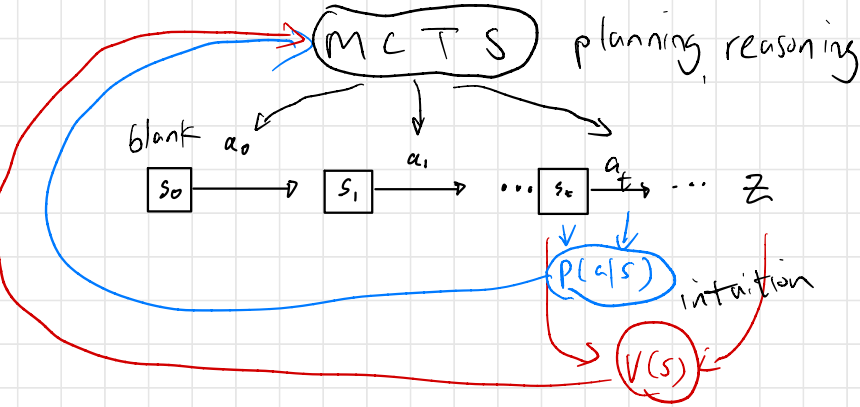
Temperature

throw away tree.

- Alpha Go Zero

↓  
No human

- In England: (Self-play)



Markov decision Process

$$\left\{ \begin{array}{l} \text{dynamics} \quad P(s_{t+1} | s_t, a_t) \\ \text{reward} \quad P(r_t | s_t, a_t, s_{t+1}) \end{array} \right.$$

$$s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} \dots \xrightarrow{a_t} s_t \xrightarrow{a_t} s_{t+1} \xrightarrow{a_{t+1}} \dots$$

model-based  
model-free

$$G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$$

max  $E_{\text{policy}}(G_t)$

recursive

$$\begin{aligned} G_t &= r_t + \gamma G_{t+1} \\ &= r_t + \gamma r_{t+1} + \gamma^2 G_{t+2} \end{aligned}$$

temporal difference