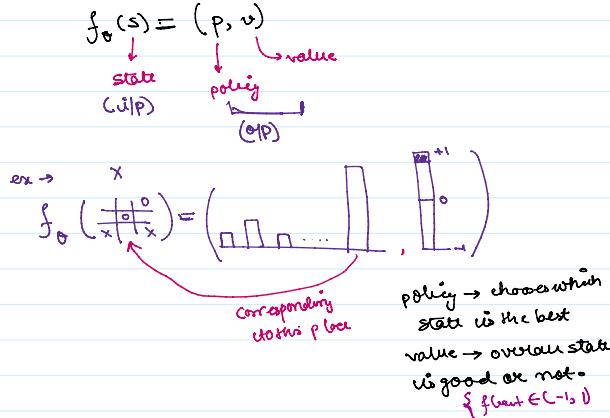
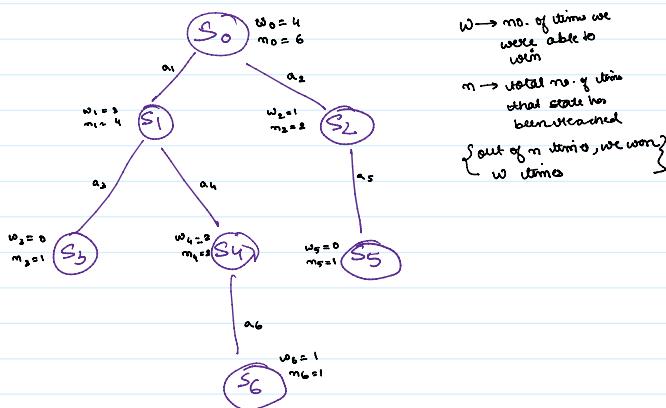
MODEL:-MONTE CARLO TREE SEARCH

ex → for  $a_1 \rightarrow a_1 = \frac{3}{4}$   
 $a_2 = \frac{1}{2}$

$$(a_2 >> a_1)$$

- 1) Selection | walk down with leaf mode
- 2) Expansion | create new mode
- 3) Simulation | play randomly
- 4) Backpropagation

which direction to walk down?

follow UCB formula

$$\frac{w_i}{m_i} + c \sqrt{\frac{\ln(N_i)}{m_i}} \quad | \quad c=2$$

highest winning ratio  
child mode has been visited the least no. of times

$$\text{ex } \Rightarrow a_1 = \frac{3}{4} + 2 \sqrt{\frac{\ln(6)}{4}}$$

$\approx 2.089$

$$a_2 = \frac{1}{2} + 2 \sqrt{\frac{\ln(1)}{2}} \\ = 2.393$$

so,  $a_2 >> a_1$

↓  
walk down via  $a_2$

- 1) include policy and value  
→ no random future

## i) include policy and value

$$UBC = \left( \frac{w_i}{m_i} \right) \rightarrow P_i C = \frac{\sqrt{N_i}}{1+m_i}$$

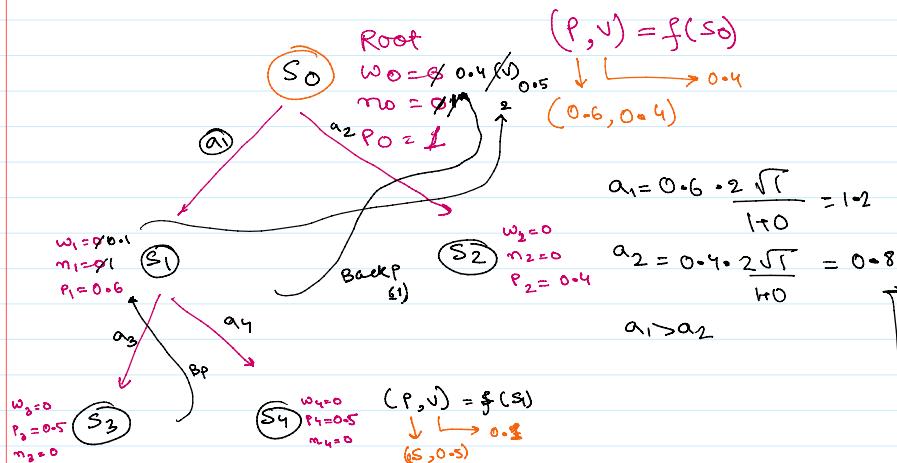
iff  $m_i > 0$

no random future iteration

children with higher policy gets selected

\* Simulation step skipped

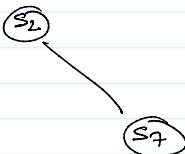
ii) Create all possible nodes rather than a single node in expansion phase.



$S_0, a_2 \gg a_1$

walk down via  $a_2$

Expansion:



Create new node  
 $w_7 = 0$   
 $m_7 = 0$

Simulation → play until game over.  
random play



Back propagation

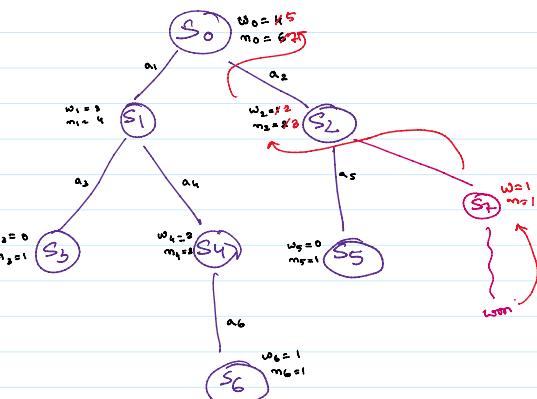
Since we won the game, we won 1 time out of 1 time we came in  
S4

## Self play

X: → Perform MCTS  $\pi$  [mix - distribution?]  
 $\{0.13, 0.09, 0.12, 0.09, 0.08, 0.10\}$   
 $\{0.16, 0.10, 0.13\}$   $Z=1 \rightarrow$  Reward / final result of player

O:   
 $\{0.14, 0.06, 0.11, 0.13, 0.23, 0.07, 0.0, 0.12, 0.14\}$   $Z=-1$

Draw / one player win



## Training

1)  $S, R, Z = \text{Sample}$  | Take sample from training data

2)  $f_\theta(s) = (P, V)$  - | Get off p from model

3) minimize the difference between the policy  $P$  and MCTS distribution  $\pi$

minimizing variance and final reward difference.

$$J = (z - v)^T - \pi^T \log p + \alpha \| \theta \|^2$$

Minnimize loss by  
backpropagation