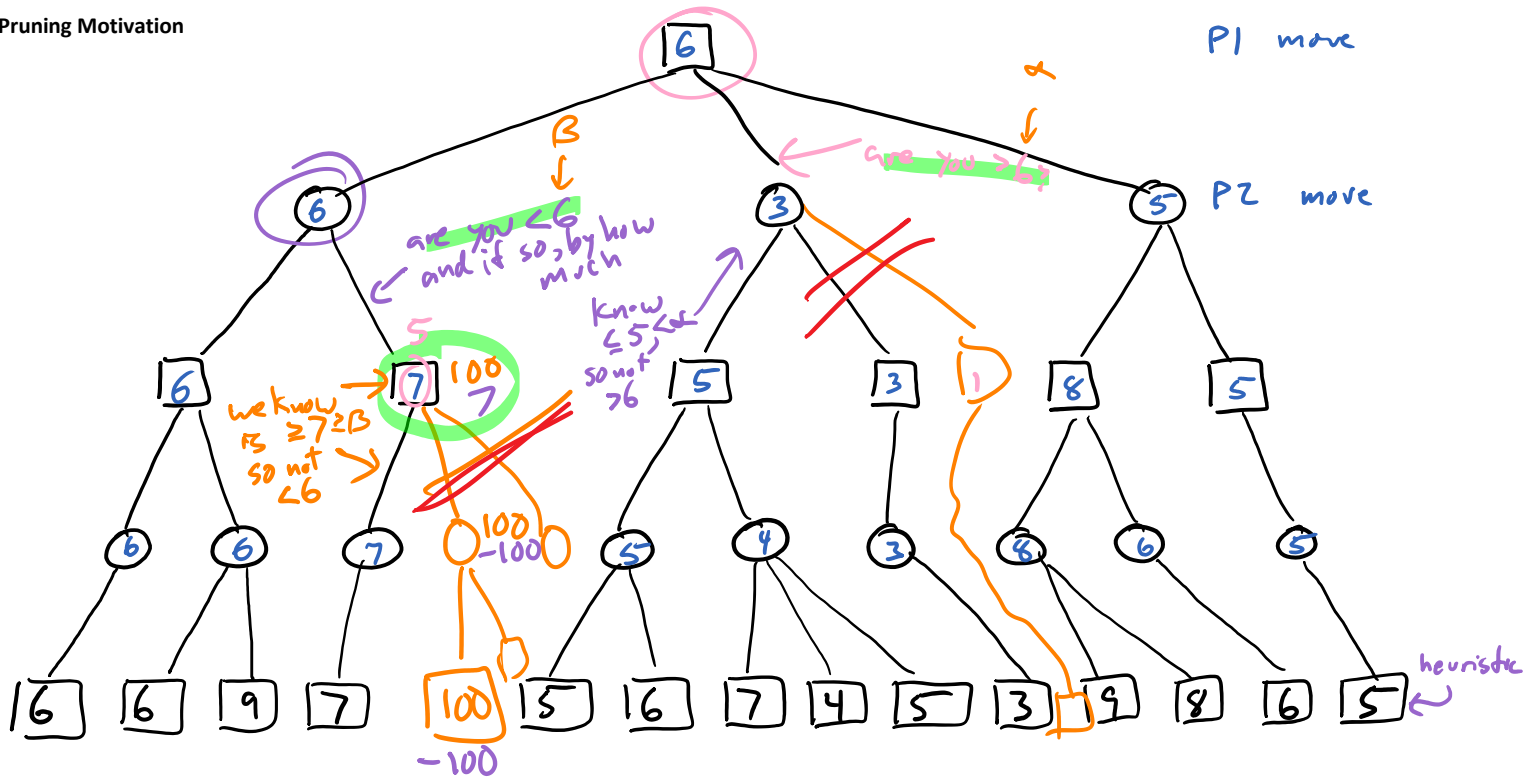


Pruning Motivation



Modified example from [http://en.wikipedia.org/wiki/Alpha%E2%80%93beta\\_pruning](http://en.wikipedia.org/wiki/Alpha%E2%80%93beta_pruning)

# Alpha-Beta Pruning

range of values to distinguish between

Alpha-Beta ( $p, \alpha, \beta, h, \text{depth}$ ) returns

postconditions

$\text{value}(p)$  if  $p$  is terminal  
 $h(p)$  if  $\text{depth} = 0$   
 $MM(p, h, d)$  if  $-\infty < MM(p, h, d) < \infty$   
 lower bound  $\geq \beta$  on  $MM(p, h, d)$  if  $MM(p, h, d) \geq \beta$   
 upper bound  $\leq \alpha$  on  $MM(p, h, d)$  if  $MM(p, h, d) \leq \alpha$

if  $\text{depth} = 0$  then return heuristic( $p$ )  
 if  $p$  is terminal then return value( $p$ )

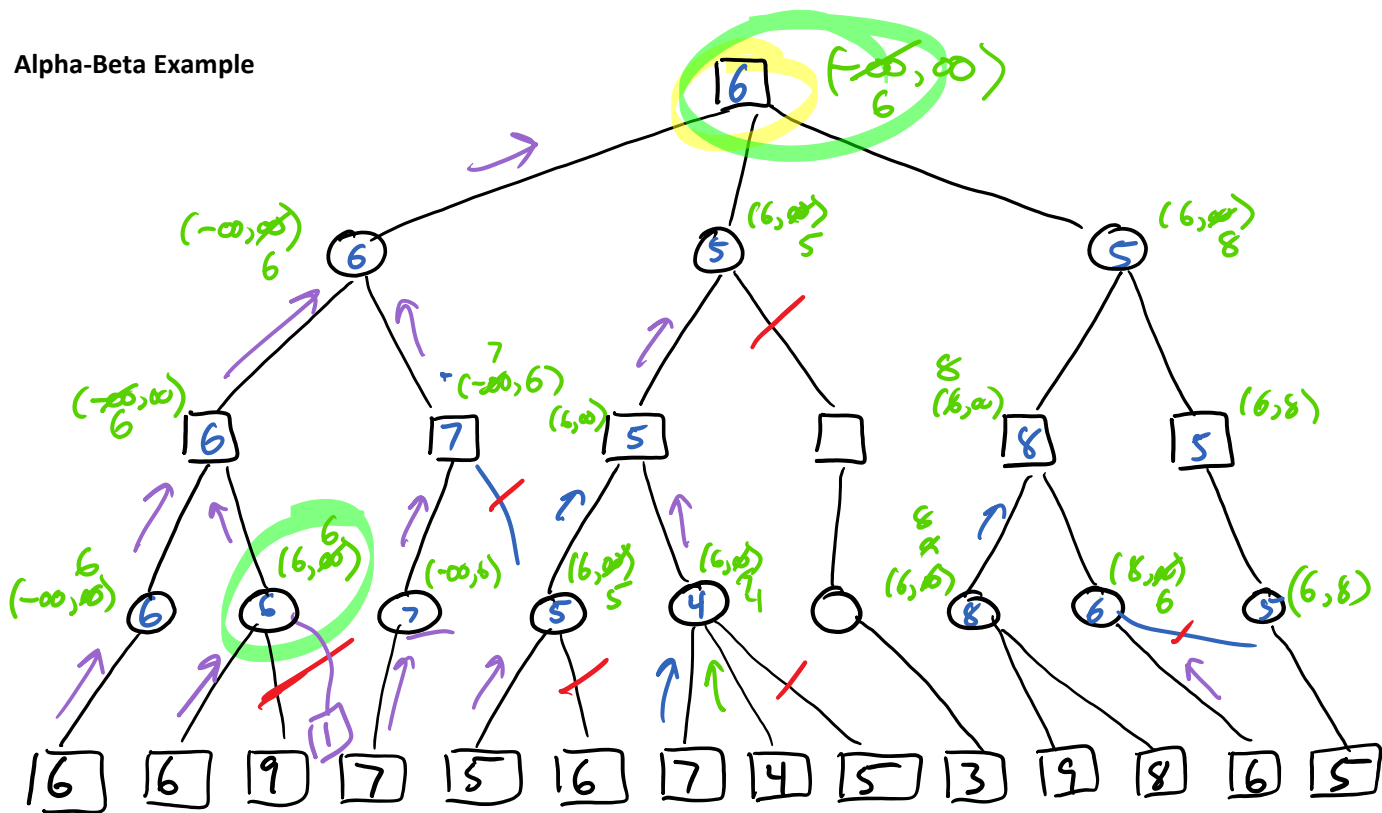
can use transposition table too → keys include  $\alpha, \beta$  values specify  $\leq, \geq, \text{or} ?$

if  $p$  is a max position  
 $a \leftarrow -\infty$  (value of best child so far - a lower bound on value of  $p$ )  
 for each position  $p'$  reachable in one move from  $p$  and while  $\alpha < \beta$   
 $a \leftarrow \max(a, AB(p', \alpha, \beta, h, \text{depth}-1))$   
 $\alpha \leftarrow \max(\alpha, a)$   
 return  $a$

else  
 $b \leftarrow \infty$   
 for each position  $p'$  reachable in one move from  $p$  and while  $\alpha < \beta$   
 $b \leftarrow \min(b, AB(p', \alpha, \beta, h, \text{depth}-1))$   
 $\beta \leftarrow \min(\beta, b)$   
 return  $b$

$\alpha = -\infty$   
 1st child makes  $\alpha = 6$   
 ? to 2nd child is value  $\geq 6$   
 $\alpha$

Alpha-Beta Example



Modified example from [http://en.wikipedia.org/wiki/Alpha%E2%80%93beta\\_pruning](http://en.wikipedia.org/wiki/Alpha%E2%80%93beta_pruning)

# Multi-Armed Bandit

Given unknown probability distributions  $R_1, \dots, R_k$   
with means  $\mu_1, \dots, \mu_k$

Choose indices  $i_1, i_2, \dots$  to optimize payout play arm w/ highest  $\mu_i$

Regret = diff between observed reward and best possible expected reward

$$P_T = T \cdot \underbrace{\mu^*}_{\text{highest } \mu_i} - \sum_{t=1}^T \underbrace{\hat{r}_t}_{\text{reward obtained at time } t}$$

"optimal" means zero average regret

$$P\left(\lim_{T \rightarrow \infty} \frac{P_T}{T} = 0\right) = 1$$

Ex:

Arm 1	Arm 2	Arm 3
prob $\frac{1}{3}$	prob $\frac{1}{4}$	prob $\frac{1}{100}$
payout 2	payout 3	payout 200
$\frac{2}{3}$	$\frac{1}{4}$	$\frac{99}{100}$
0	$\frac{1}{2}$	0
$\frac{2}{3}$	0	
$\mu_1 = \frac{2}{3}$	$\mu_2 = \frac{7}{8}$	$\mu_3 = 2$

uniform rotation: cycle through each arm 1 2 3 1 2 3 ...

$\frac{1}{3}$  of time average  $(2 - \frac{1}{3})$  regret

$$\lim_{T \rightarrow \infty} \frac{P_T}{T} \geq \frac{4}{9}$$

greedy: play each machine once, then always play one with highest reward seen

greedy : play each machine once, then always play one w/ highest reward from 1st round

if pick best machine after 1 round,  $\lim_{T \rightarrow \infty} \frac{R_T}{T} = \mu^*$

but  $P(\text{pick best machine}) < 1$

$\epsilon$ -greedy : play one round, then best observed mean reward w/p  $1-\epsilon$  and randomly w/p  $\epsilon$

balances exploration / exploitation mean regret from arm 1

still has  $\lim_{T \rightarrow \infty} \frac{R_T}{T} \geq \underbrace{\frac{\epsilon}{2}}_{P(\text{choose arm 1})} \cdot \frac{4}{3}$

zero regret  
UCB

Choose arm  $j$  that maximizes

$\underbrace{\bar{r}_j}_{\text{exploitation}} + \sqrt{\frac{2 \ln T}{\underbrace{n_j}_{\text{exploration}}}}$   
 mean reward for arm  $j$       # plays of arm  $j$

Monte Carlo Tree Search  
tree ← root

Until out of time

traverse tree root → leaf

expand if leaf <sup>non-terminal and non-zero visits</sup> expandable, add its children

simulate play to terminal pos (from arb. child if newly expanded)  
can be random

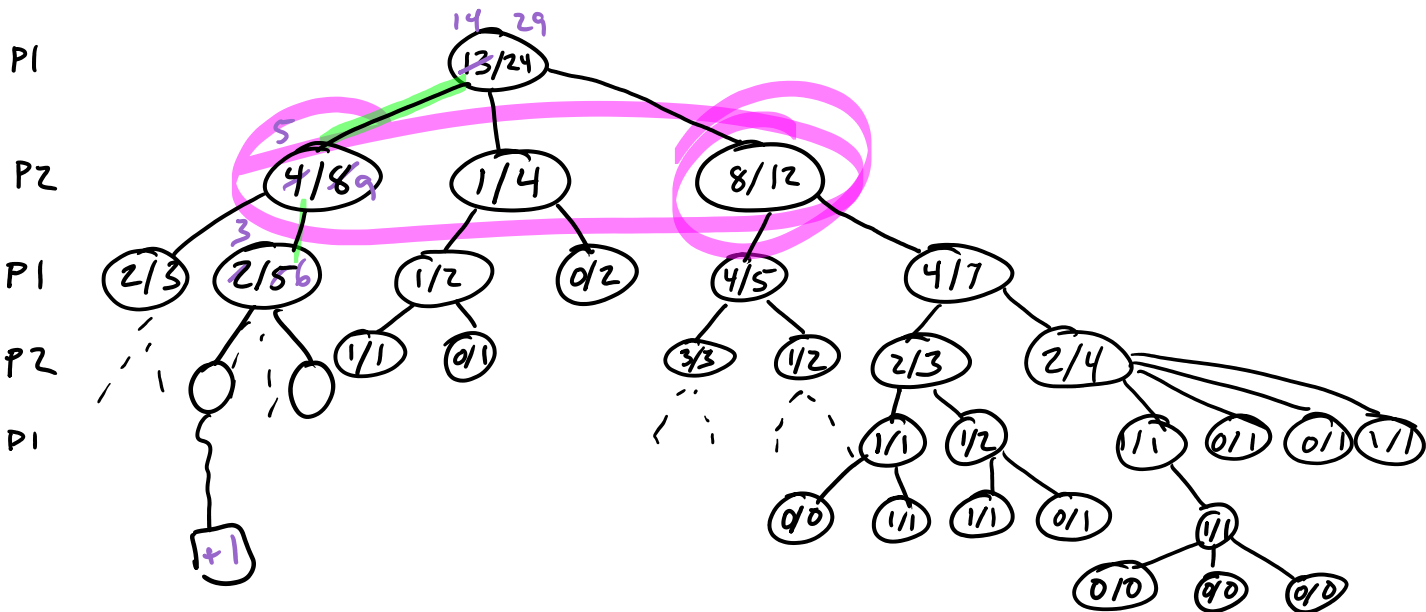
update backup result up tree from start of simulation → root

return max to child w/ best stats (highest mean reward or highest visit count)

$$UCT = MCTS + UCB$$

asymmetric tree growth  
no heuristic needed

UCB (choosing a child w/o visits when one exists)



$$UCT = MCTS$$

advantages: convergent

anytime

no domain knowledge

easily parallelized

disadvantages: no domain knowledge  
some games not amenable

MCTS adapted for games that aren't trees

