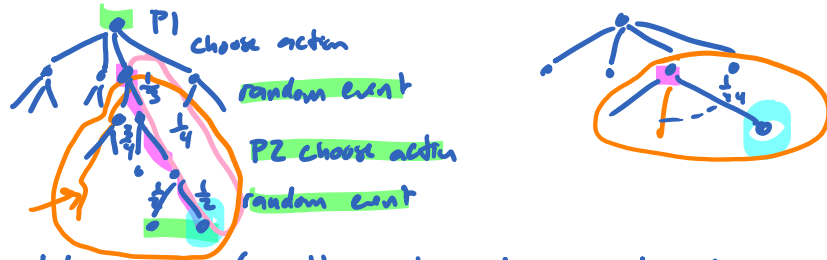


Two-Player Games

$p(s', r | s, a)$  satisfies  $p(s', r | s, a) = 0$  if  $s'$  nonterminal  
 or  $s'$  is P1 win and  $r \neq 1$   
 $s'$  is P2 win and  $r \neq -1$   
 $s'$  is draw and  $r \neq 0$

Can you model 2<sup>nd</sup> player? and P1 want adjust to P1

If so:



If not: find equilibrium (neither player has incentive to change policy)

$$v^*(s) = \max_a \sum_{s', r} p(s', r | s, a) \cdot [r + \gamma \cdot v^*(s')] \quad \text{for zero-sum if } s \text{ is P1 turn}$$

$g(s, a)$

$$v^*(s) = \min_a \sum_{s', r} p(s', r | s, a) \cdot [r + \gamma \cdot v^*(s')] \quad \text{if } s \text{ is P2 turn}$$

Two-player Yahtzee

anchors: P1

$2^{12}$  used / unused  
 3 Yahtzee unused / 0 / 50  
 106 upper total

P2

used / unused  
 Yahtzee unused / 0 / 50  
 upper total

Score difference -1500 ... +1500

$2^{12}$   
 3  
 106  
 3000<sup>4</sup>

3 billion times longer than solitaire

$$v_{\pi}(s) = \sum_a \pi(a|s) \cdot \sum_{s', r} p(s', r | s, a) \cdot (r + \gamma v_{\pi}(s'))$$

$$= \sum_{s'} \sum_r p(s', r | s, \pi(s)) \cdot (r + \gamma v_{\pi}(s'))$$

for deterministic  $\pi$

substitute  $v[s']$  for  $v_{\pi}(s')$

Initialize  $v[s] \leftarrow 0$  for terminal  
and arbitrarily for other

$\Delta \leftarrow 0$   
repeat

for each state  $s$   
 $v_{old} \leftarrow v[s]$

let  $v[s] \leftarrow$

until  $\Delta \leftarrow \max(\Delta, |v[s] - v_{old}|)$   
small enough

## Policy Iteration

Initialize  $\pi[s]$  arbitrarily

repeat

evaluate  $\pi$  to get  $v_\pi$

change  $\leftarrow$  false

for each state  $s$

$a_{old} \leftarrow \pi[s]$

$\pi[s] \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) \cdot (r + \gamma v_\pi[s'])$

if change  $\leftarrow$  false and  $\pi[s] \neq a_{old}$   
change  $\leftarrow$  true

until change = false