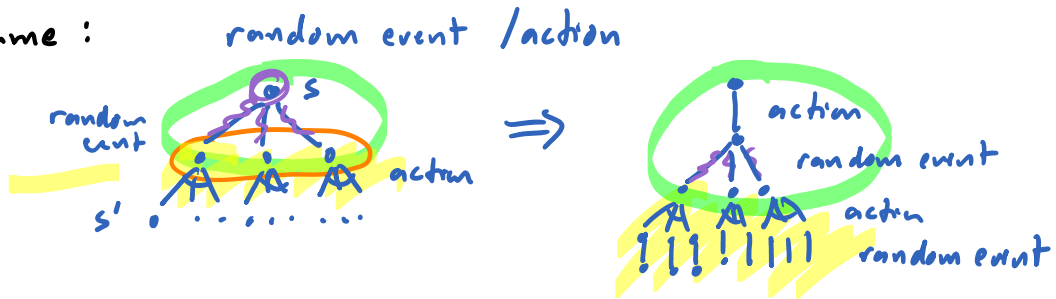


Traditional Markov Decision Process:

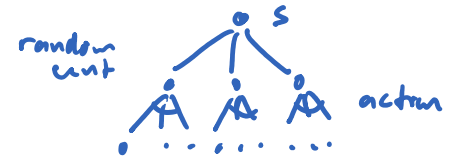
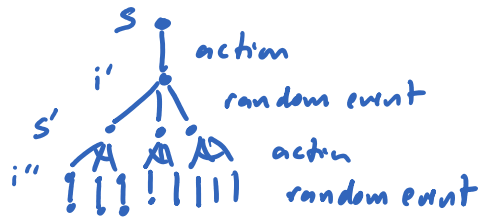
$$v^*(s) = \max_a \sum_{s',r} p(s',r | s, a) \cdot [r + \gamma \cdot v^*(s')]$$

Typical Game:



for random event:  $v^*(s) = \sum_{s',r} p(s',r | s) \cdot [r + \gamma \cdot v^*(s')]$

for player action:  $v^*(s') = \max_a [r + \gamma \cdot v^*(s'')]$

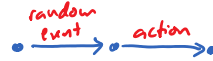


$$V^*(s) = \sum_{s', r} p(s' | i'') \cdot v^*(s')$$

$$\Rightarrow V^*(s) =$$

$$v^*(s') = \max_{i''} (r + v^*(s''))$$

alternative view



for  $P_2$  :  $V_2(s) =$  expected wins for  $P_2$

$$V_2(S, r, t) = \begin{cases} 1.0 & \text{if no } S' \subseteq S \text{ has } \text{sum}(S') = r \text{ and } \text{sum}(S) < t \\ 0.5 & \text{if no } S' \subseteq S \text{ has } \text{sum}(S') = r \text{ and } \text{sum}(S) = t \\ 0.0 & \text{if no } S' \subseteq S \text{ has } \text{sum}(S') = r \text{ and } \text{sum}(S) > t \end{cases}$$

$\max_{S' \subseteq S, \text{sum}(S')=r} V_2(S-S', t)$

$$V_2(S, t) = \sum_{r=\min(S)}^{\max(S)} p(\text{roll } r) \cdot V_2(S, r, t)$$

-- move outputs  $\text{argmax}_{S' \subseteq S, \text{sum}(S')=r} V_2(S-S', t)$

- $V_2(\{1,3,5,8,9\}, 9, 1)$
  - $V_2(\{1,3,5,8,9\}, 9, 2)$
  - $\vdots$
  - $V_2(\{1,3,5,8,9\}, 9, 13)$
- actions always  $\{\{1,8\}, \{9\}\}$   
 $\{1,3,5\}$

for  $P_1$  :  $V_1(s) =$  expected wins for  $P_1$

$$V_1(S, r) = \begin{cases} 1 - V_2(\{1, \dots, 9\}, \text{sum}(S)) & \text{if no } S' \subseteq S \text{ has } \text{sum}(S') = r \\ 1.0 & \text{if } \text{sum}(S) = r \\ \max_{S' \subseteq S, \text{sum}(S')=r} V_1(S-S') & \end{cases}$$

-- expect

$$V_1(S) = \sum_{r=\min(S)}^{\max(S)} p(\text{roll } r) \cdot V_1(S, r)$$

anchor: start of turns

component: start of one turn and the next

number of anchors:

aces... sixes	$7^6$
3K, 4K	$2 \cdot 8^2$
C	$2^7$
FH, SS, LS, T	$3^4$
Yahtzee bonus	<u>13</u>

$\approx 2$  trillion anchors  
 • 1800 states/anchor  
 $\approx 4$  quadrillion states  
 50 days @ 1 billion states/sec

aces... sixes	$2^6$
3K, 4K	$2^2$
C	2
FH, SS, LC	$2^3$
upper total 0...63	64
i	<u>3</u>
used 50	
used 0	
unused	

$\frac{3}{4}$  million  
 $\approx 1.3$  billion state  
 $\approx 1$  min in Java code



## Evaluating a Policy

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

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

for deterministic policy  $\pi$

current estimate of  $v_{\pi}(s)$

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

$\Delta \leftarrow 0$   
→ repeat (in any order)  
for each state  $s$

$v_{old} \leftarrow v[s]$

let  $v[s] \leftarrow$

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