

Markov Decision Process: (S, A, P, R)

states (positions) S , actions (moves) $A(s)$, rewards (scores) R

$$P: S \times A \times S \times R \rightarrow [0, 1]$$

current state S , action A , next state S , reward R

$P(s', r | s, a)$ = prob. of reaching s' w/ reward r given current state s , action a

$$P(s' | s, a) = \sum_r P(s', r | s, a)$$

expected reward $r(s, a) = \sum_r r \cdot \sum_{s'} P(s', r | s, a)$

$$r(s, a, s') = \frac{\sum_r r \cdot P(s', r | s, a)}{P(s' | s, a)}$$

initial state s , terminal state s'

Episode: $(s_0, a_0, s_1, R_1), (s_1, a_1, s_2, R_2), \dots, [(s_{T-1}, a_{T-1}, s_T, R_T)]$

$G = R_1 + R_2 + R_3 + \dots$ goal: maximize sum of rewards

$$= R_1 + \gamma R_2 + \gamma^2 R_3 + \dots = \sum_{t=1}^{\infty} \gamma^{t-1} R_t$$

$0 < \gamma \leq 1$
for finite process, can take $\gamma = 1$

Policy: $\pi: A \times S \rightarrow [0, 1]$ s.t. $\sum_a \pi(a|s) = 1$

prob action a taken in state s
for deterministic policy $\pi(a|s) = 1$ for one a for any given s
write $\pi(s) = a$ when a is

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

$v_{\pi}(s)$: value of states under policy π
expected future reward from state s following π

$\sum_{s'} \sum_r P(s', r | s, a) \cdot r$: all new states s' rewards r (immediate reward)
 $\gamma v_{\pi}(s')$: future rewards from new state s'

$$g_{\pi}(s, a_1) < g_{\pi}(s, a_2)$$

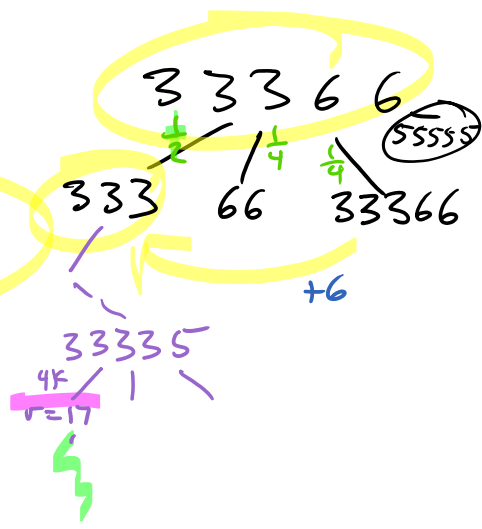
0.5	0.5
1.0	0.0
0.0	1.0

optimal value

$$v^*(s) = \max_{\pi} v_{\pi}(s)$$

$$g^*(s, a) = \max_{\pi} g_{\pi}(s, a)$$

Bellman Eq $\pi^*(s) = \operatorname{argmax}_a g^*(s, a)$



$$g^*(s,a) = \max_{\pi} g_{\pi}(s,a)$$

$$\pi^*(s) = \operatorname{argmax}_{\pi} g^{\pi}(s,a)$$

$$v^*(s) = \max_{a \in A(s)} g^*(s,a) = \max_a \max_{\pi} g_{\pi}(s,a)$$

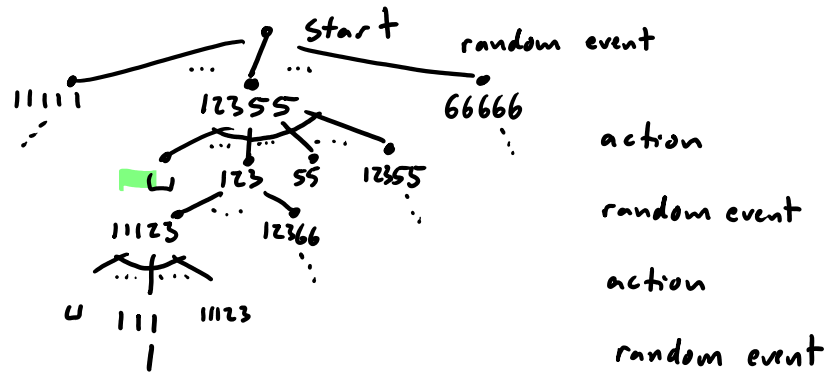
$$= \max_a \sum_{s',r} p(s',r|s,a) \cdot [r + \gamma v^*(s')]$$

Bellman Eq

which π gives max?
 π that gives max $v_{\pi}(s')$;
 which gives $v^*(s')$

Some deterministic
 π gives max
 (for g^* too)

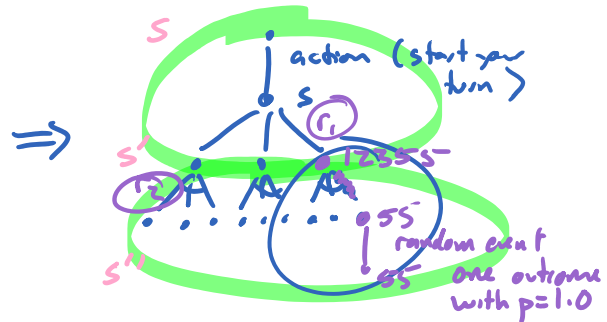
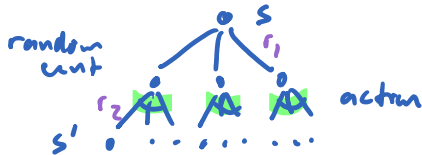
for a finite process, this is a recurrence (base cases for terminal states)



Traditional Markov Decision Process: action / random event

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

Many Typical Game: random event / action



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

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

for player action

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

assume deterministic rewards

$$= \max_a [r(s',a) + \gamma v^*(n(s',a))]$$

state s'' that results from taking action a in state s'

anchor : start of turn

component : states between start of one turn and next

number of anchors: attempt 1
(all in 6)

ones... sixes	7^6
3K, 4K	28^2
↙	27
FH, LS, SS, Y	3^4
Yahtzee bonus	13

≥ 2 trillion anchors

- 1800 states / anchor
- ≈ 4 quadrillion states

50 days @ 1 billion / sec
very fast

ones... sixes	used/unused	2^6
upper total	0...63	64
3K, ... C	used/unused	2^6
Y	used 10/50	3

$\frac{3}{4}$ million

≈ 1 minute (Java code)

Yahtzee Graph

