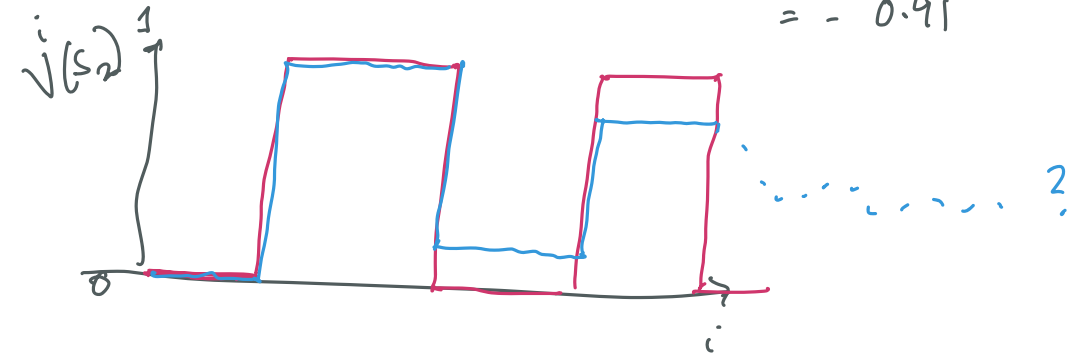


$$V(s_1) = -1 + \gamma V(s_2)$$

$$V(s_2) = 1 + \gamma V(s_1)$$

$$\gamma \approx 0.9$$



$$V(s_1) = 1, 1-\gamma, 1-\gamma+\gamma^2, 1-\gamma+\gamma^2-\gamma^3, \dots$$

$0 < \gamma < 1$

$$\frac{1}{1+\gamma}$$

$$= \frac{1}{1+0.9} \approx 0.5$$

$$V(s_2) = -1 + \gamma - \gamma^2 + \gamma^3 - \dots = \frac{-1}{1+\gamma}$$

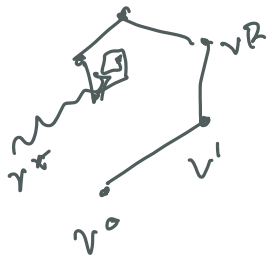
<u>Iter</u>	<u><math>V(s_1)</math></u>	<u><math>V(s_2)</math></u>
1	0	0
2	-1	1
3	0	0
4	-1	1
5	0	0

---

1.	0	0
2.	-1	1
3.	-0.1	0.1
4.	$-1 + 0.9 \times 0.1$ $-1 + 0.09$ $= -0.91$	0.91

# CONVERGENCE OF VALUE ITERATION

$\gamma < 1$



BELLMAN OPERATOR

$$v^{k+1} = \max_a \left( C(s,a) + \gamma \sum_{s'} P(s'|s,a) v^k \right)$$

B

Piece 1 Bellman operator is a contraction.

$$\|B U - B V\|_{\infty} \leq \gamma \|U - V\|_{\infty}$$

$v^{k+1} = B v^k$

Piece 2

$$\|V^* - v^{k+1}\|_{\infty} = \|B V^* - B v^k\|_{\infty}$$

$$\leq \gamma \|V^* - v^k\|_{\infty}$$

(A)

$$\|V^* - v^k\|_{\infty} \leq \gamma \|V^* - v^{k-1}\|_{\infty}$$

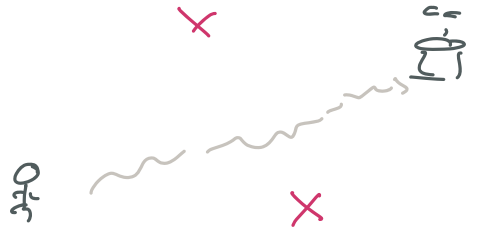
$$\leq \gamma^2 \|V^* - v^{k-1}\|_{\infty}$$

⋮

$$\leq \gamma^{k+1} \|V^* - v^0\|_{\infty}$$

# ATTRIBUTES OF A HARD MDP

- ① PLAN DEEP INTO THE FUTURE
- ② LARGE STATE SPACE / ACTION SPACE  
(CONTINUOUS ACTS)
- ③ CYCLES
- ④ STOCHASTIC MDP ARE HARDER  
(NOISE)



( R - S.AT vs SAT )

GOAL IS TO FIND A POLICY  $\hat{\pi}$  S.T.

$$\| J(\hat{\pi}) - J(\pi^*) \|_{\infty} \leq \epsilon$$

WITH MINIMUM PLANNING EFFORT

PERFORMANCE DIFFERENCE LEMMA

$\pi$

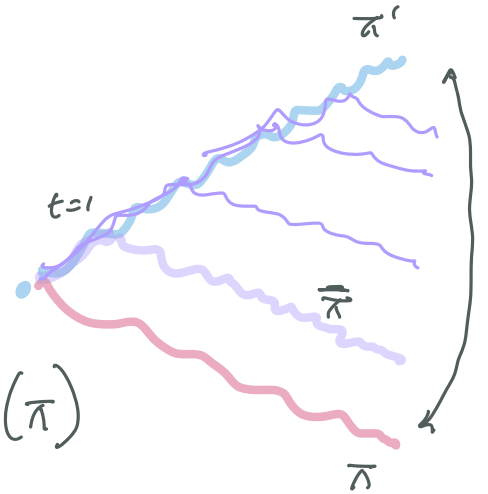
$\pi'$

TELESCOPING TRICK

$$J(\pi') - J(\pi)$$

PERFORMANCE OF  $\pi'$

PERF OF  $\pi$



$$= J(\pi') - J(\bar{\pi}) + J(\bar{\pi}) - J(\pi)$$

$$= V^{\pi'}(s_0) - V^{\pi}(s_0)$$

$$= V^{\pi'}(s_0) - \left( c(s_0, a_0) + \gamma V^{\pi}(s_1) \right) + \left( c(s_0, a_0) + \gamma V^{\pi}(s_1) \right) - V^{\pi}(s_0)$$

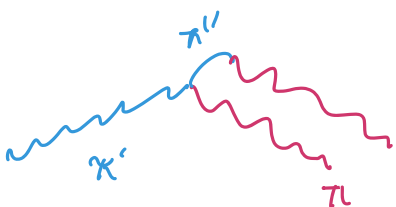
$a_0 \sim \pi'$                        $a_0 \sim \pi$

$$= \left( c(s_0, a_0) + \gamma V^{\pi'}(s_1) \right) - \left( c(s_0, a_0) + \gamma V^{\pi}(s_1) \right) + \dots$$

$a_0 \sim \pi'$                        $a_0 \sim \pi$

$$= \underbrace{\gamma \left( V^{\pi'}(s_1) - V^{\pi}(s_1) \right)}_{A^{\pi}(s_0, \pi'(s_0))} + \underbrace{Q^{\pi}(s_0, \pi'(s_0)) - V^{\pi}(s_0)}_{A^{\pi}(s_0, \pi'(s_0))}$$

$$J(\pi') - J(\pi) = \sum_{t=0}^{T-1} \mathbb{E}_{s_t \sim \pi'} \left[ \gamma^t A^{\pi}(s_t, \pi'(s_t)) \right]$$



$$\pi' = \underset{a}{\operatorname{argmax}} \left( c(s, a) + \gamma V^{\pi}(s_1) - V^{\pi}(s) \right)$$