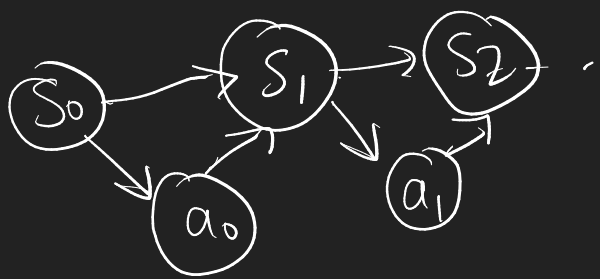


## 5) State-Action Distribution

Trajectory of MDP up to  $t$ ,  
 $(s_0, a_0, s_1, a_1, \dots, s_t, a_t)$

Consider possible stochastic policies

$$\mathbb{P}^\pi(s_0, a_0, \dots, s_t, a_t) = \pi(a_0 | s_0) P(s_1 | s_0, a_0) \\ \times \pi(a_1 | s_1) P(s_2 | s_1, a_1)$$



$$\times P(s_t | s_{t-1}, a_{t-1}) \\ \times \pi(a_t | s_t)$$

$$\mathbb{P}_t^\pi(s, a; s_0) = \sum_{\substack{a_{0:t-1} \in \mathcal{A}^t \\ s_{1:t-1} \in \mathcal{S}^{t-1}}} \mathbb{P}^\pi(s_0, a_0, \dots, s_t, a_t)$$

# Discounted Average State-Action Distribution

$$d_{s_0}^{\pi}(s, a) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t P_t^{\pi}(s, a; s_0)$$

HW0: is this a valid distribution?

$$V^{\pi}(s_0) = \frac{1}{1-\gamma} \sum_{\substack{s \in \mathcal{S} \\ a \in \mathcal{A}}} d_{s_0}^{\pi}(s, a) r(s, a)?$$

## 2) Optimal Policies & Bellman Optimality

$A^S$  policies!

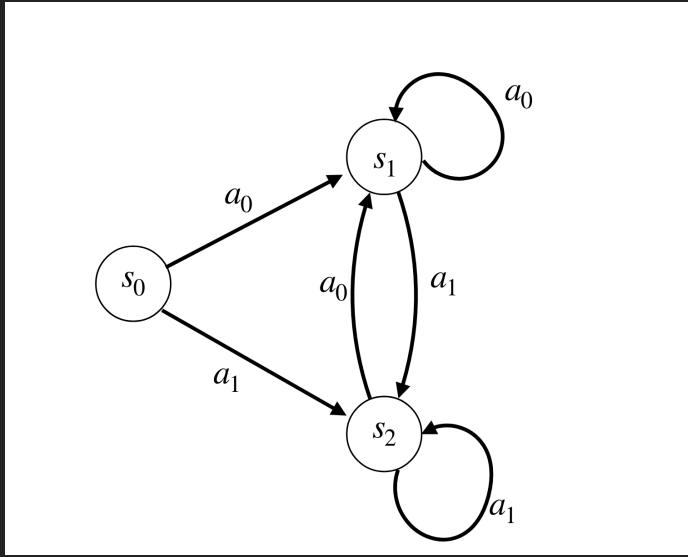
$$\pi^* = \underset{\pi}{\operatorname{argmax}} \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \begin{array}{l} s_{t+1} \sim P(s_{t+1}, a_t) \\ a_t = \pi(s_t) \end{array} \right]$$

Fact: there always exist  $\pi^* : \mathcal{S} \rightarrow \mathcal{A}$  s.t.  $V^{\pi^*}(s) \geq V^{\pi}(s)$

→ dominates for all  $\pi$  and  $s \in \mathcal{S}$

notation:  $V^{\pi^*} = V^*$      $Q^{\pi^*} = Q^*$

# Example: deterministic MDP



$$r(s_1, a_0) = 1$$

0 otherwise

$$V^{\pi_0}(s) = \begin{cases} \frac{\gamma}{1-\gamma} & s_0, s_2 \\ \frac{1}{1-\gamma} & s_1 \end{cases}$$

$$\pi_0(s) = a_0 \quad \forall s$$

PollEV.com / sarahdean011

## Bellman Optimality

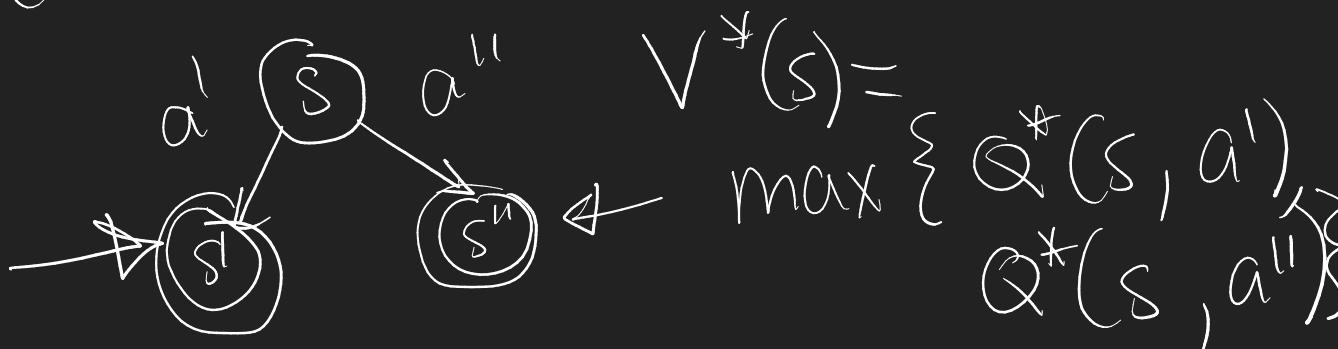
$$Q^*(s, a)$$

Theorem 1:

$$V^*(s) = \max_{a \in \mathcal{A}} \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} [V^*(s')] \right]$$

$$\forall s \in \mathcal{S}$$

If  $V^*(s')$  is known,  
can compute  $V^*(s)$



Proof: We show

$$\hat{\pi}(s) = \arg \max_{a \in \mathcal{A}} Q^*(s, a)$$

$$\text{that } V^{\hat{\pi}}(s) = V^*(s)$$

a) by definition

$$V^*(s) \geq V^{\hat{\pi}}(s) \quad \forall s.$$

b) now show  $V^*(s) \leq V^{\hat{\pi}}(s) \quad \forall s.$

$$\begin{aligned} \underline{V^*(s)} &= r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} [V^*(s')] \\ &\leq \max_a r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} [V^*(s')] \\ &= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} [V^*(s')] \end{aligned}$$

$$\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} [r(s', \hat{\pi}(s'))] + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} [V^*(s'')] + \dots$$

$$\leq \mathbb{E}_{\substack{s' \sim P(s, \hat{\pi}(s)) \\ s'' \sim P(s', \hat{\pi}(s'))}} (r(s, \hat{\pi}(s)) + \gamma r(s', \hat{\pi}(s')) + \gamma^2 r(s'', \hat{\pi}(s'')) + \dots)$$

$$= V^{\hat{\pi}}(s)$$

□

This means that  $\hat{\pi}$  achieves optimal value, so

$$\pi^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q^*(s, a)$$

is an optimal policy

Now, show Bellman optimality is sufficient to characterize  $V^*(s)$ .

Theorem 2: for any  $V: \mathcal{S} \rightarrow \mathbb{R}$   
if  $V(s) = \max_{a \in \mathcal{A}} r(s, a) + \gamma \mathbb{E}[V(s')] ]$   
 $s' \sim P(s, a)$

for all  $s \in \mathcal{S}$ , then  
 $V(s) = V^*(s)$ .

This means we can confirm that  $V = V^*$  by checking

$$|V(s) - \max_{a \in \mathcal{A}} r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')| = 0 \quad \forall s.$$

Proof:

$$|V(s) - V^*(s)| =$$

$$|\max_a r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s') -$$

$$\left[ \max_a r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right]|$$

(basic inequality)

HWO  
(Jensens)

$$\leq \max_a |r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s') -$$

$$\max_a |r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')|$$

$$|V(s) - V^*(s)| \leq \max_{a, a'} \gamma \mathbb{E}_{s' \sim P(s, a)} |V(s') - V^*(s')|$$

$$\leq \max_{a, a'} \gamma^2 \mathbb{E}_{s' \sim P(s, a)} \mathbb{E}_{s'' \sim P(s', a')} |V(s'') - V^*(s'')|$$

iterate  $k$  times

limit to infinity  
 $|V(s) - V^*(s)| \rightarrow 0$

$$V(s) = V^*(s) \quad \square$$

Example: check  $\pi(s) = a_0 \forall s$   
is optimal using  
Bellman optimality.

### 3) Value Iteration

How to find optimal policy?  
( $\pi^*$ ,  $V^*$ ,  $Q^*$ )

enumeration:  $\mathcal{O}(\mathcal{A}^S \cdot S^3)$   
policy evaluator

Define Bellman Operator  $\mathcal{T}$

given  $Q: S \times \mathcal{A} \rightarrow \mathbb{R}$ , Bellman operator

$\mathcal{T}Q: S \times \mathcal{A} \rightarrow \mathbb{R}$ , defined as

$$\mathcal{T}Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \left[ \max_{a' \in \mathcal{A}} Q(s', a') \right]$$

Think of representing tabular  
 $Q \in \mathbb{R}^{SA}$      $\mathcal{J}: \mathbb{R}^{SA} \rightarrow \mathbb{R}^{SA}$

## Fixed Point Iteration

Bellman optimality

$$\star Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(s, a)} \left[ \max_{a'} Q^*(s', a') \right]$$

$$\mathcal{J}Q^* = Q^* \quad (\text{fixed point})$$

Alg: Value Iteration (VI)

Initialize  $Q^0$

for  $t=0, 1, 2, \dots$

$$Q^{t+1} \leftarrow \mathcal{J}Q^t$$

## Convergence

Lemma (contraction) for any  $Q, Q'$

$$\|\mathcal{J}Q - \mathcal{J}Q'\|_\infty \leq \|Q - Q'\|_\infty$$

remember  $\|Q\|_\infty = \max_{s, a} Q(s, a)$



Proof:

$$\begin{aligned} & |JQ(s, a) - JQ'(s, a)| \\ &= \cancel{V(s, a)} + \gamma \mathbb{E}_s \left[ \max_{a'} Q(s', a') \right] \\ &\quad - \left( \cancel{V(s, a)} + \gamma \mathbb{E}_{s'} \left[ \max_{a'} Q'(s', a') \right] \right) \end{aligned}$$

(basic  
inequalities)  
HW 0

$$f(s) \leq \max_s f(s)$$

$$\leq \gamma \mathbb{E}_s \left[ \max_{a'} Q(s', a') \right]$$

$$\leq \gamma \mathbb{E}_{s'} \left[ \max_{a'} |Q(s', a') - \max_{a'} Q'(s', a')| \right]$$

$$\leq \gamma \max_s \max_{a'} |Q(s, a) - Q'(s, a)|$$

$$= \gamma \|Q - Q'\|_\infty$$

Lemma (convergence)

$$\|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$$

Recall:  $\pi^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} \underline{Q^*(s, a)}$

misc notes from office hours:

$$\mathbb{E}(f(x)) \leq \max_{x \in X} f(x)$$

$x \sim \Delta(X)$        $\uparrow$

$$f(x) \leq \max_{y \in X} f(y)$$

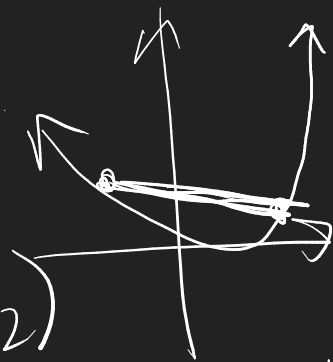
---

Jensen's inequality

$$f(\mathbb{E}(x)) \leq \mathbb{E}[f(x)]$$

$x \sim \mathcal{D}$        $x \sim \mathcal{D}$

$f$  is convex.



$$f(tx_1 + (1-t)x_2)$$

$$\geq tf(x_1) + (1-t)f(x_2)$$