

Policy Iteration: $\pi_k \leftarrow \pi_{\underline{Q}^{\pi_{k-1}}}$

PD Lemma: $\forall \pi, \pi'$,

$$\underline{V^{\pi'}(s) - V^{\pi}(s)} = \frac{1}{1-\gamma} \mathbb{E}_{\tilde{s} \sim d^{\pi', s}} [Q^{\pi}(\tilde{s}, \pi') - V^{\pi}(\tilde{s})]$$

Lemma: $\forall f \in \mathbb{R}^{S \times A}$

$$\underline{\|V^* - V^{\pi_f}\|_{\infty}} \leq \frac{2 \cdot \|f - Q^*\|_{\infty}}{1-\gamma}$$

Proof: $V^*(s) - V^{\pi_f}(s) = \frac{1}{1-\gamma} \mathbb{E}_{\tilde{s} \sim d^{\pi_f, s}} [V^*(\tilde{s}) - Q^*(\tilde{s}, \pi_f)]$

$$\frac{1}{1-\gamma} \mathbb{E}_{d^{\pi_f, s}} [Q^*(\tilde{s}, \pi^*) - Q^*(\tilde{s}, \pi_f)]$$

$$\leq \underline{2 \cdot \|f - Q^*\|_{\infty}}$$

□

Then: $\|Q^* - Q^{\pi_{k+1}}\|_{\infty} \leq \gamma \cdot \|Q^* - Q^{\pi_k}\|_{\infty}$ ✓

Proof: Claim: (1) $\mathcal{T}^{\pi_{k+1}} Q^{\pi_k} \geq \mathcal{T}^{\pi} Q^{\pi_k} \forall \pi$.
 (2) $\mathcal{T}^{\pi_{k+1}} Q^{\pi_k} \leq Q^{\pi_{k+1}}$ ✓

For (1): $(\mathcal{T}^{\pi} Q^{\pi_k})(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [Q^{\pi_k}(s', \pi)]$

For (2): $(\mathcal{T}^{\pi_{k+1}} Q^{\pi_k})(s, a) = \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} V_t \mid s_1 = s, a_1 = a, a_2 \sim \pi_{k+1}, a_3: \infty \sim \pi_k \right]$
 $Q^{\pi_{k+1}}(s, a) = \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} V_t \mid s_1 = s, a_1 = a, a_2: \infty \sim \pi_{k+1} \right]$

$$\begin{aligned}
 Q^* - Q^{\pi_{k+1}} &= Q^* - \mathcal{T}^{\pi_{k+1}} Q^{\pi_k} + \mathcal{T}^{\pi_{k+1}} Q^{\pi_k} - Q^{\pi_{k+1}} \\
 &\leq \mathcal{T}^{\pi^*} Q^* - \mathcal{T}^{\pi_{k+1}} Q^{\pi_k} \leq 0 \\
 &\leq \mathcal{T}^{\pi^*} Q^* - \mathcal{T}^{\pi^*} Q^{\pi_k} \\
 &\leq \gamma \cdot \|Q^* - Q^{\pi_k}\|_{\infty}, \checkmark
 \end{aligned}$$

$$\begin{aligned}
 &Q^* - \mathcal{T}^{\pi_{k+1}} Q^{\pi_k} \\
 &= \mathcal{T} Q^* - \mathcal{T} Q^{\pi_k}
 \end{aligned}$$

$$\begin{aligned}
 \mathcal{T} \circ \mathcal{T} &= \mathcal{T} \\
 \mathcal{T} \circ \mathcal{T} &\neq \mathcal{T}
 \end{aligned}$$

$$T^{\pi} f = R(s, a) + \gamma \mathbb{E}_{\dots} [f(s', \pi f)].$$

$$\max_{a'} f(s', a')$$

$$= T f.$$

Linear Programming for MDPs.

Primal Form:

$$\begin{array}{l} \min_{V \in \mathbb{R}^S} \quad d_0^T V \\ \text{s.t.} \quad V \geq TV. \end{array}$$

Claim: if $d_0(s) > 0 \quad \forall s$.

then optimal sol. $V = V^*$.

Q1: why min? Q2: why linear?

$$Q1: \quad V \geq TV \Rightarrow V \geq V^* \quad \checkmark$$

Lemma: $\forall f \geq f' \in \mathbb{R}^S, T f \geq T f'$

$$(T f)(s) = \max_a (R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [f(s')]).$$

$$V \geq TV \Rightarrow TV \geq T(TV)$$

$$V \geq T^2V \geq T^3V \geq \dots \geq \underbrace{T^{\infty}V}_{= V^*}$$

$$V \geq TV : \forall s.$$

$$V(s) \geq \max_a (R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [V(s')])$$

$\forall a.$

$$V(s) \geq R(s,a) + \gamma \mathbb{E} [V(s')]$$

Dual form: $\max_{d \in \mathbb{R}^{S \times A}, d \geq 0} \underline{d^T R}$

s.t. $\forall s'.$

$$\sum_{a'} d(s', a') = (1-\gamma) d_0 + \gamma \sum_{s,a} P(s'|s,a) d(s,a)$$

distribution over S .
 $s' \Leftrightarrow s, a \sim d.$
 $s \sim P(\cdot|s,a)$

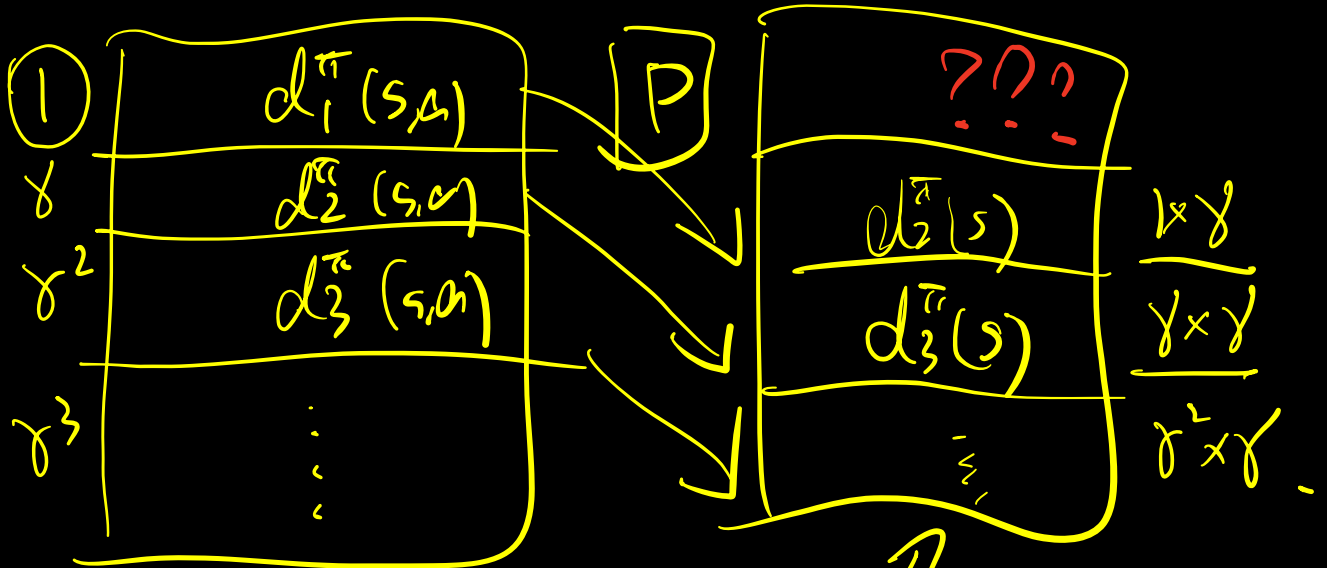
Claim: feasible space.

$$= \{ \underline{d}^\pi(s,a) : \forall \pi \}$$

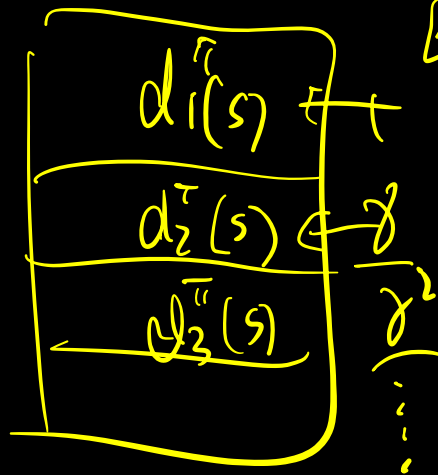
$$d^{\pi, s}, \quad d^{\pi}$$

$$d^{\pi}(s, a) = (1-\gamma) \sum_{t=1}^{\infty} \gamma^{t-1} d_t^{\pi}$$

distribution of (s_t, a_t) .



$$\sum_{a'} d^{\pi}(s', a') =$$



$$Q^{\bar{\pi}} = T^{\bar{\pi}} Q^{\bar{\pi}}$$

$$Q^* = \underline{T} Q^*$$

$$Q^* = Q^{\pi^*}$$

$$u \in \mathbb{R}^n, v \in \mathbb{R}^n.$$

$$\langle u, v \rangle$$

$$|u^T v| \leq \|u\|_1 \|v\|_\infty.$$

$$\mathbb{E}_{a \sim \pi} [Q(s, a)]$$

$$Q(s, \pi) - Q(s, \pi')$$

$$= \langle \underline{Q}(s, \cdot), \pi(\cdot | s) \rangle$$

$$- \langle \underline{Q}(s, \cdot), \pi'(\cdot | s) \rangle.$$