

Alt. proof for monotone improvement.

$$PI: \pi_{k+1} \leftarrow \pi_{Q^{\pi_k}}$$

$$\text{Monotone improvement: } V^{\pi_{k+1}} \geq V^{\pi_k}$$

Performance difference Lemma. [Kakade & Langford '02]

$\forall \pi, \pi', s$. (P-D Lemma)

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

$$\begin{aligned} Q^{\pi'}(s', \pi') - V^{\pi}(s) &= \underbrace{Q^{\pi'}(s', \pi') - Q^{\pi}(s', \pi)}_{=: A^{\pi}(s', \pi') \text{ "advantage"}} \\ &=: A^{\pi}(s', \pi') \text{ "advantage"}. \end{aligned}$$

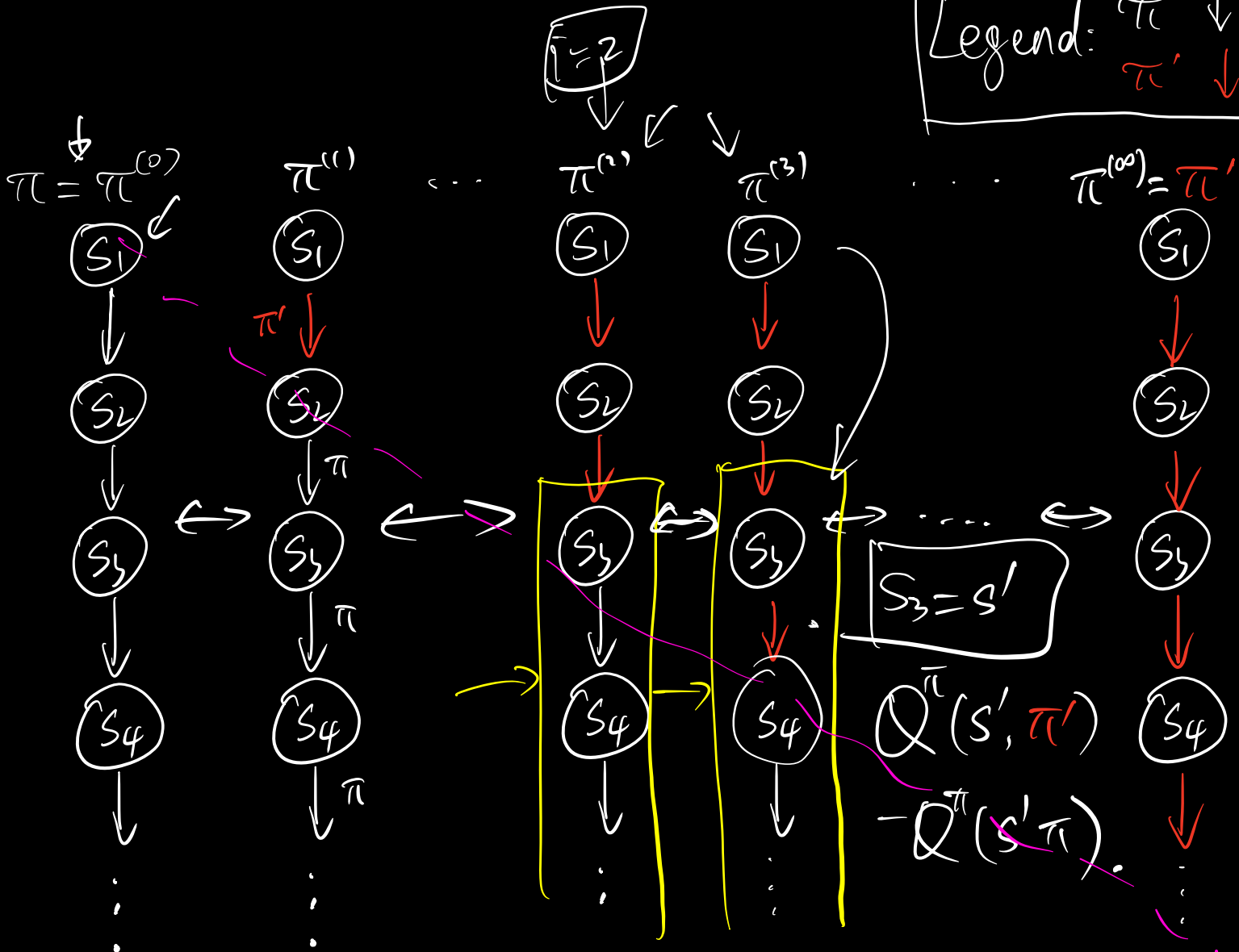
Use P-D Lemma to establish monotone impr.

$$V^{\pi_{k+1}}(s) - V^{\pi_k}(s) = \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_s^{\pi_{k+1}}} \left[Q^{\pi_k}(s', \pi_{k+1}) - Q^{\pi_k}(s', \pi_k) \right]$$

$$\text{Recall: } \pi_{k+1} \leftarrow \pi_{Q^{\pi_k}}$$

$$\text{i.e. } \forall s, \pi_{k+1}(s) = \operatorname{argmax}_a Q^{\pi_k}(s, a)$$

Legend: $\pi \downarrow$
 $\pi' \downarrow$



$$V^{\pi'}(s) - V^\pi(s) = \sum_{i=0}^{\infty} V^{\pi'}(s) - V^\pi(s)$$

$$= \sum_{i=0}^{\infty} \gamma^i \mathbb{E}_{s' \sim d_{s, i+1}^{\pi'}} [Q^{\pi'}(s', \pi') - Q^\pi(s', \pi)]$$

$s_i = s$

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

Linear Programming

$$\begin{aligned} \max & \quad x_1 + x_2 \\ \text{s.t.} & \quad x_1 - 2x_2 \leq 1. \end{aligned}$$

Solve for V^* (primal).

Choose any $\mu \in \mathbb{R}^S$ s.t. $\forall s, \mu(s) \geq 0$.

$$\begin{aligned} \min_{V \in \mathbb{R}^S} & \quad \mu^T V \\ \text{s.t.} & \quad \mathcal{T}V \leq V. \end{aligned}$$

$$\sum_s \mu(s) = 1.$$

① why is opt. V^* ?

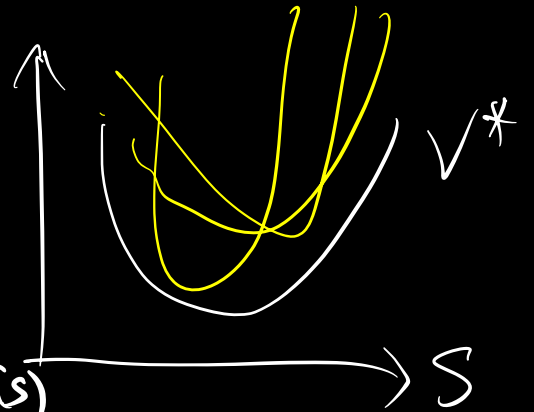
② why linear?

①: $\mathcal{T}V \leq V \Rightarrow \mathcal{T}(\mathcal{T}V) \leq \mathcal{T}V \leq V.$

$$V^* = \mathcal{T}^\infty V \leq V.$$

\therefore any feasible V must. $V \geq V^*.$

$\&$ V^* is feasible.



② Constraint: $\forall s.$

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

$\forall s, a$

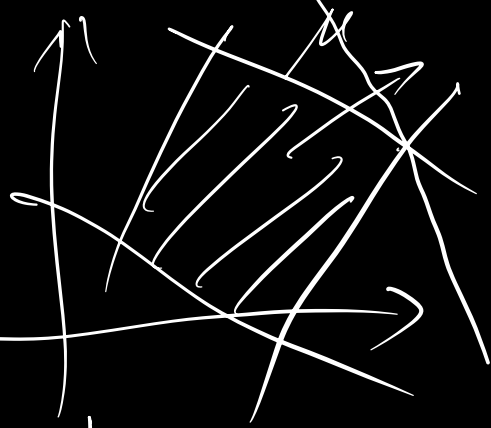
$$R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [V(s')] \leq V(s).$$

Dual form. $\mathbb{R}^{S \times A}$

$\max d^T R. \quad (d_\mu^\pi)^T R$

$d \geq 0$

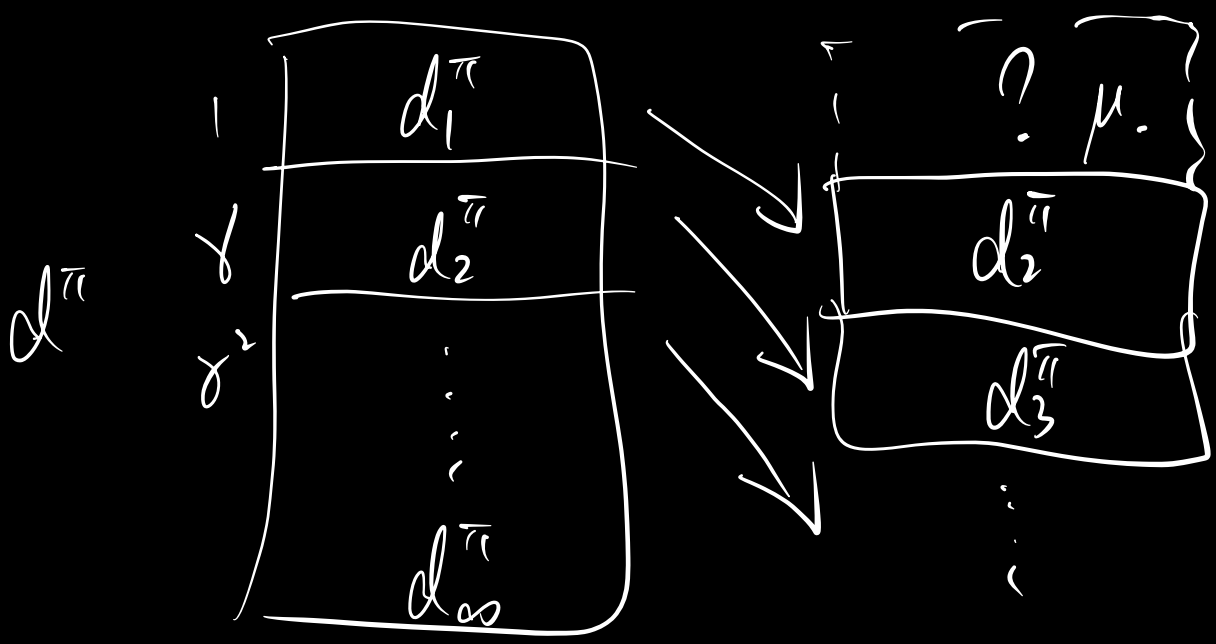
$= (1-\gamma) V^\pi(\mu)$



st. $\forall s'$.

$\sum_{a'} d(s', a') = \boxed{\mu(s')} + \gamma \sum_{s, a} d(s, a) P(s' | s, a)$

any feasible $d = d_\mu^\pi. \forall \pi.$



$$V^{\pi^{(3)}}(s_1) - V^{\pi^{(2)}}(s_1)$$

$$= \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid s_1, \pi^{(3)} \right]$$

$$- \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid s_1, \pi^{(2)} \right]$$

$$= \mathbb{E} \left[\sum_{t=3}^{\infty} \gamma^{t-1} r_t \mid s_1, \pi^{(3)} \right]$$

$$- \mathbb{E} \left[\sum_{t=3}^{\infty} \gamma^{t-1} r_t \mid s_1, \pi^{(2)} \right]$$

$\gamma^2 \cdot Q^{\pi}(s', \pi)$

$$= \mathbb{E} \left[\mathbb{E} \left[\sum_{t=3}^{\infty} \gamma^{t-1} r_t \mid s_1, s_3 = s', \pi^{(3)} \right] \right]$$

$$- \mathbb{E} \left[\mathbb{E} \left[\sum_{t=3}^{\infty} \gamma^{t-1} r_t \mid s_1, s_3 = s', \pi^{(2)} \right] \right]$$

$$s_3 = s', \pi^{(2)}$$

$$Q^\pi(s', \pi).$$