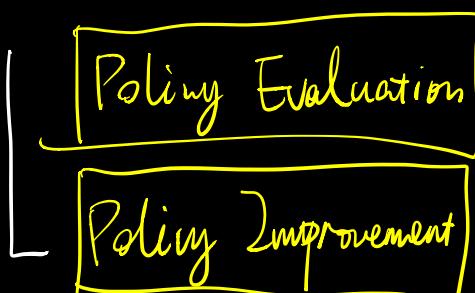


## Policy Iteration (cont.)

$$\pi_f(s) = \underset{a \in A}{\operatorname{argmax}} f(s, a)$$

Initialize  $\pi_0$ .

For  $k=1, 2, \dots$ ,



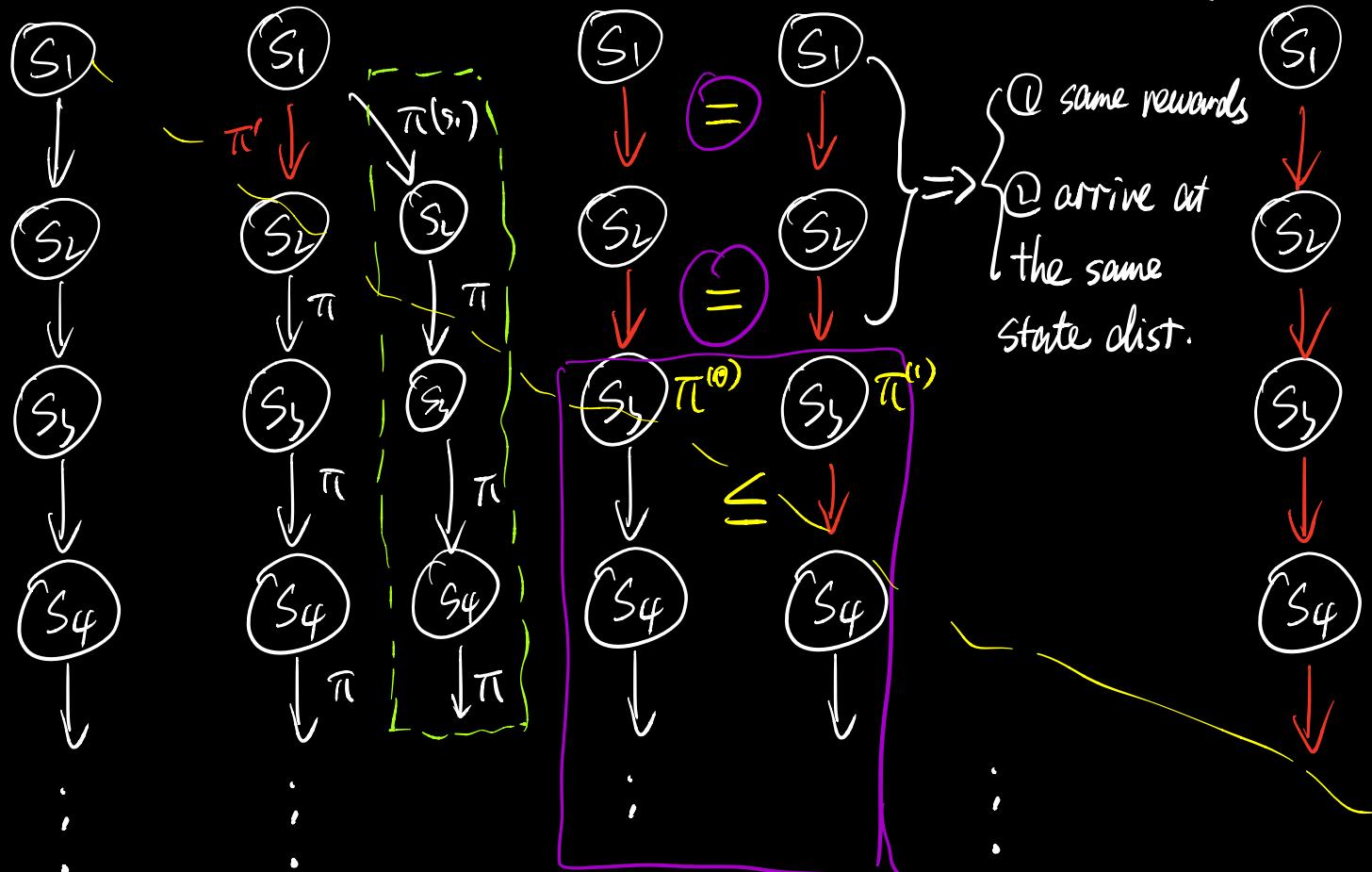
Compute  $Q^{\pi_{k-1}}$

$$\pi_k \leftarrow \pi \quad Q^{\pi_{k-1}}$$

$$\downarrow \pi' (= \pi_k)$$

$$\downarrow \pi (= \pi_{k-1})$$

$$\pi^{(0)} = \pi \leq \pi^{(1)} \leq \pi^{(2)} \leq \pi^{(3)} \leq \dots \leq \pi^{(\infty)} = \pi'$$



Q: if starting in  $S_1 = s$ , how good is  $\pi^{(1)}$ ?

A:  $Q^\pi(s, \pi') \geq Q^\pi(s, \pi) = V^\pi(s).$

b/c:  $\pi'$  greedily maximizes  $Q^\pi(s, \cdot)$

key:  
 $\pi^{(1)}$  better than  $\pi^{(0)}$   
 holds for all starting state

PD-Lemma:  $\forall \pi, \tilde{\pi} \in \mathcal{V}$ . more general.

$$\mathbb{V}_{\underline{s}}^{\pi} - \mathbb{V}^{\pi}(s) = \frac{1}{|S|} \sum_{s \sim d_s} [\mathbb{Q}^{\pi}(s, \tilde{\pi}) - \mathbb{V}^{\pi}(s)].$$

Linear Programming (LP)

Primal LP for solving  $\mathbb{V}^* \in \mathbb{R}^{|S|}$

Choose any  $d_0 \in \mathbb{R}^{|S|}$  s.t.  $\forall s, d_0(s) > 0$ .

$$\sum_s d_0(s) = 1$$

- ① why this solves  $\mathbb{V}^*$ ?
- ② why this is LP?

$$\begin{array}{ll} \min_{V \in \mathbb{R}^{|S|}} & d_0^T V \\ \text{s.t.} & \mathcal{T}V \leq V \end{array}$$

nonlinear

linear  
(equality or  
inequality)  
constraints.

$$\begin{aligned} \text{recall:} \\ \mathcal{T}^{\pi} V = \\ R^{\pi} + \gamma P^{\pi} V \end{aligned}$$

①  $\mathcal{T}V \leq V \Rightarrow \underbrace{\mathcal{T}(\mathcal{T}V)}_{\text{monotone property of } \mathcal{T}} \leq \mathcal{T}V \leq V \Rightarrow \dots$

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

original problem  $\Leftrightarrow \boxed{\begin{array}{ll} \min & d_0^T V \\ \text{s.t.} & V \geq \mathbb{V}^* \end{array}}$   $\Rightarrow V = \mathbb{V}^*$ .

② Why linear? obj is linear ✓.

constraint is:  $\mathcal{T}V \leq V \Leftrightarrow \forall s, (\mathcal{T}V)(s) \leq V(s)$ .

$$\rightarrow \max_{a \in A} (R(s, a) + \gamma \langle P(\cdot | s, a), V \rangle) \leq V(s). \quad \boxed{\forall s}$$

$$R(s, a) + \gamma \langle P(\cdot | s, a), V \rangle \leq V(s).$$

$\int s_a$

Primal form:  $\min_{V \in \mathbb{R}^{|S|}}$   $d_0^T V$ . s.t.  $R + \gamma P V \leq V$ .

$$\frac{|S \times A| \times 1}{|S \times A| \cdot |S|} \quad \frac{|S|}{|S \times A| \cdot |S|}$$

Dual form:  $\max_{\mu \geq 0} \underbrace{\mu^T R}_{\text{Bellman "flow" eq.}} \rightarrow |S \times A| \times 1 \text{ reward vector}$

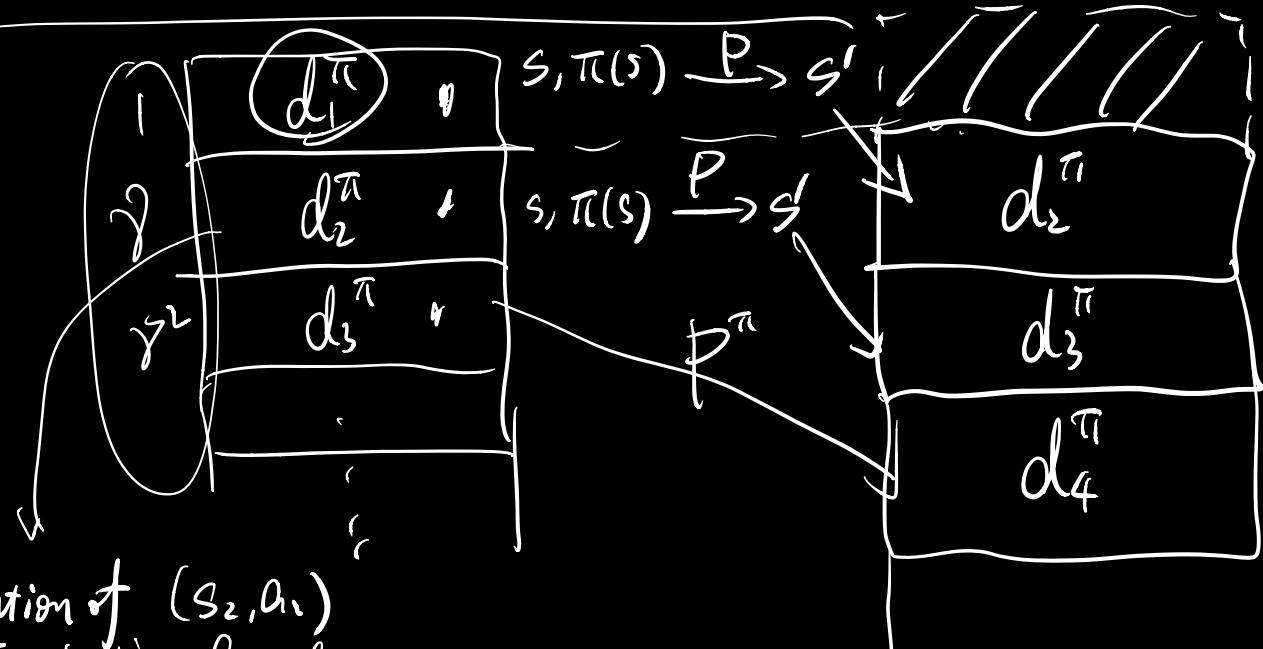
s.t.  $\sum_{a'} \mu(s', a') = d_0(s') + \gamma \sum_{s, a} \mu(s, a) \cdot P(s' | s, a), \forall s'$

Interpretation: any feasible  $\mu = \tilde{d}^\pi$  for some stochastic policy  $\pi$ .

$\langle \tilde{d}^\pi, R \rangle = \mathbb{E}_{s \sim d_0} [V^\pi(s)]$

$\downarrow$  state-action occupancy  
 $\uparrow$  unnormalized

$\pi(a|s) = \frac{\mu(s, a)}{\sum_{a'} \mu(s, a')}$ .



distribution of  $(s_2, a_1)$  under  $\pi$ , starting from  $d_0$