

Policy iteration

init π_0 .

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

policy eval. $Q^{\pi_{k-1}}$

policy improve $\pi_k \leftarrow \pi_{Q^{\pi_{k-1}}}$

$$\pi_f(s) = \operatorname{argmax}_a f(s, a)$$

$$\pi_{k+1} = \pi^* \rightarrow Q^{\pi_{k+1}} = Q^* \rightarrow \pi_k = \pi_{Q^*} = \pi^*.$$

Monotone improvement: (1) $\forall k, \underline{V^{\pi_k} \geq V^{\pi_{k-1}}}$

(2) as long as $\pi_{k-1} \neq \pi^*, \exists s.$

$$V^{\pi_k}(s) > V^{\pi_{k-1}}(s).$$

Corollary: $\pi_k = \pi^* \quad \forall k \geq |A|^{|S|}$

Performance difference lemma: $\forall \pi, \pi'.$

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

$$(1-\gamma) \sum_{t=1}^{\infty} \gamma^{t-1} \underbrace{d_t^{\pi', s}}_{\text{advantage}}$$

$$\mathbb{P}[s_t = \cdot \mid s_1 = s, \pi'].$$

$A^{\pi}(\tilde{s}, \pi')$
"advantage"

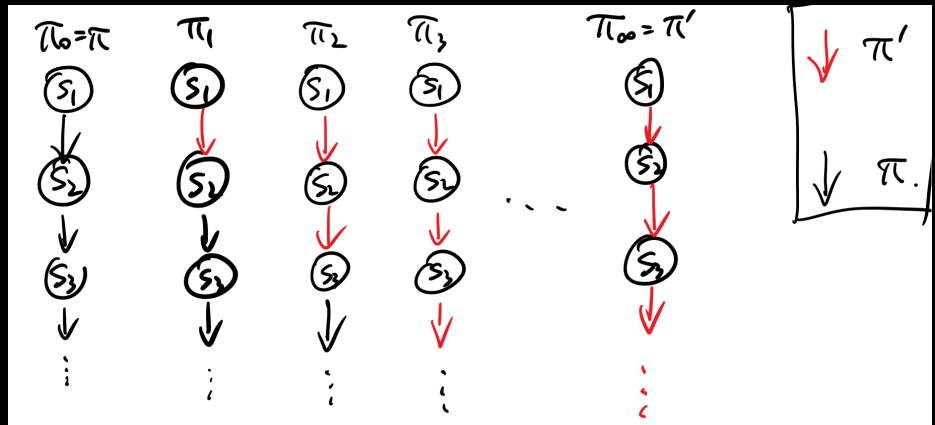
Given P-D Lemma:

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

$$\pi_k \leftarrow \pi Q^{\pi_{k-1}} \quad \boxed{\geq 0}$$

$$V^{\pi_{k-1}}(\tilde{s})$$

$$V^{\pi}(s) - V^{\pi'}(s) = \sum_{i=0}^{\infty} V^{\pi_{i+1}}(s) - V^{\pi_i}(s)$$



$$V^{\pi_2}(s) = \mathbb{E}[r_1 + \gamma r_2 | s_1=s, \pi'] + \mathbb{E}\left[\sum_{t=3}^{\infty} \gamma^{t-1} r_t \mid s_1=s, a_{1:2} \sim \pi', a_{3:\infty} \sim \pi\right]$$

$$V^{\pi_3}(s) = \mathbb{E}[r_1 + \gamma r_2 | s_1=s, \pi'] + \mathbb{E}\left[\sum_{t=3}^{\infty} \gamma^{t-1} r_t \mid s_1=s, a_{1:3} \sim \pi', a_{4:\infty} \sim \pi\right]$$

$$\gamma^2 \mathbb{E}\mathbb{E}\left[\sum_{t=3}^{\infty} \gamma^{t-3} r_t \mid s_1=s, a_{1:2} \sim \pi', s_3=s', a_3=a', a_{4:\infty} \sim \pi\right]$$

$$- \gamma^2 \mathbb{E} \left[Q^\top(s_3, \underline{\pi}) \mid s_1 = s_3, a_{1:2} \sim \underline{\pi}' \right]$$

$$+ \gamma^2 \mathbb{E} \left[Q^\top(s_3, \underline{\pi}) \mid s_1 = s, a_{1:2} \sim \underline{\pi}' \right]$$

$$= \gamma^2 \mathbb{E}_{s_3 \sim d_3^{\pi', s}} \left[Q^\top(s_3, \underline{\pi}') - Q^\top(s_3, \underline{\pi}) \right]$$

$$\sum_{i=0}^{\infty} \gamma^i \left(V^{\pi_{i+1}}(s) - V^{\pi_i}(s) \right)$$

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

$$\mathbb{E}_\phi[f] + \mathbb{E}_\theta[f]$$

$$= \langle \phi, f \rangle + \langle \theta, f \rangle$$

$$= 2 \langle \frac{\phi + \theta}{2}, f \rangle = 2 \mathbb{E}_{\frac{\phi + \theta}{2}}[f]$$

$$\boxed{\begin{aligned} \text{If } \pi_{k-1} \neq \pi^* \quad \exists \bar{s}, \quad Q^{\pi_{k-1}}(\bar{s}, \pi_{k-1}) < \max_a Q^{\pi_{k-1}}(\bar{s}, a) \\ \exists \bar{s} \quad Q^{\pi_{k-1}}(\bar{s}, \pi_k) > Q^{\pi_{k-1}}(\bar{s}, \pi_{k-1}) \end{aligned}}$$

$$\text{(o.w., } \forall s, \quad Q^{\pi_{k-1}}(s, \pi_k) = Q^{\pi_{k-1}}(s, \pi_{k-1}).$$

$$\rightarrow \sqrt{\pi^*}(s) - \sqrt{\pi_{k-1}}(s) = \frac{1}{1-\delta} \left[\underbrace{Q^{\pi_{k-1}}(s, \pi^*)}_{\pi^*} - Q^{\pi_{k-1}}(s, \pi_{k-1}) \right]$$

$$\rightarrow \sqrt{\pi_k}(s) - \sqrt{\pi_{k-1}}(s) = \frac{1}{1-\delta} \cdot \underbrace{\left[\underbrace{Q^{\pi_{k-1}}(\bar{s}, \pi_k)}_{\bar{s} \text{ d } \pi_k, s} - \underbrace{Q^{\pi_{k-1}}(\bar{s}, \pi_{k-1})}_{> 0} \right]}_{> 0}$$

$$\Rightarrow \bar{s} = \bar{s} \checkmark$$

$$\Rightarrow d^{\pi_k, \bar{s}}(\bar{s}) > 0$$

$$(\geq (1-\delta))$$