

Policy Iteration | Init  $\pi_0$  arbitrarily.

for  $k=1, 2, 3, \dots$ ,  $\pi_k \leftarrow \pi_{Q^{\pi_{k-1}}}$  "Policy Eval" step

$\pi_f(s) := \operatorname{argmax}_{a \in A} f(s, a)$ .  $\hookrightarrow$  [ 1. Compute  $Q^{\pi_{k-1}}$   
2. Take its greedy policy to be  $\pi_k$ .

Remarks:

1.  $\pi^* \xrightarrow{\text{eval}} Q^{\pi^*} = Q^* \xrightarrow{\text{improve}} \pi_{Q^*} = \pi^*$ .

2. Assume  $Q^{\pi}$  is calculated exactly.

"Policy Improvement" step.

Convergence? | Policy Improvement Theorem.

In PI,  $V^{\pi_k}(s) \geq V^{\pi_{k-1}}(s)$ ,  $\forall k \geq 1, s \in S$ .

(Further, if  $\pi_{k-1} \neq \pi^*$ , improvement is strict in at least 1 state)

Corollary: PI converges in  $|A|^{|S|}$  iterations.

(b/c, policies in  $\pi_0, \pi_1, \pi_2, \dots$ , never repeats before  $\pi^*$ .)

See standard proof in notes (using "montone property of  $T$ ")

Lemma: (Performance-Difference Lemma).  $\forall \pi, \pi', s \in S$ .

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

s-th row of matrix.

$$(I - \gamma P^{\pi})^{-1}$$

$\Leftrightarrow$  normalized discounted occupancy of  $\pi$  starting in  $s$ .

$$\downarrow \sum_{s'} \gamma P^{\pi'} + \gamma^2 (P^{\pi'})^2 + \gamma^3 (P^{\pi'})^3 + \dots$$

s-th row: if  $s_1 = s$ , dist of  $s_t$ .

where  $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$ . "advantage function".

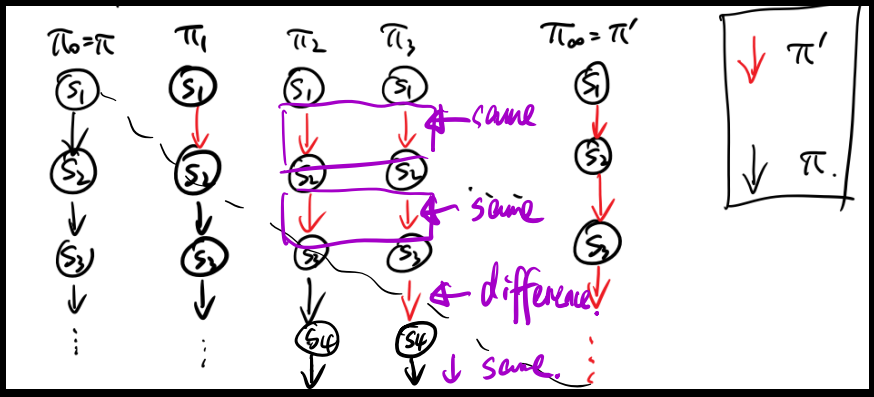
Construct  $\pi_i$ :

$$\pi_0 = \pi, \quad \pi_\infty = \pi'$$

For  $1 \leq i < \infty$ .

$\pi_i$  is non-stationary:

$$\begin{cases} a_t = \pi'(s_t) & \text{if } t \leq i. \\ a_t = \pi(s_t) & \text{if } t > i. \end{cases} \quad (s_1 = s)$$



$$V^{\pi'}(s) - V^\pi(s) = V^{\pi_\infty}(s) - V^{\pi_0}(s) = \sum_{i=0}^{\infty} (V^{\pi_{i+1}}(s) - V^{\pi_i}(s))$$

$$= \sum_{i=0}^{\infty} \gamma^i \sum_{s' \in S} \underbrace{P[s_{i+1} = s' | s_i = s, \pi']}_{\text{the distribution of } s_{i+1} \text{ for both } \pi_i \text{ \& } \pi_{i+1}} \cdot \underbrace{(Q^\pi(s', \pi') - Q^\pi(s', \pi))}_{\Delta}$$

$$= \sum_{i=0}^{\infty} \gamma^i \sum_{s' \in S} P[s_{i+1} = s' | s_i = s, \pi'] \cdot A^\pi(s', \pi')$$

$$= \sum_{s' \in S} \left( \sum_{i=0}^{\infty} \gamma^i P[s_{i+1} = s' | s_i = s, \pi'] \right) A^\pi(s', \pi')$$

$$= \sum_{s' \in S} \frac{1}{1-\gamma} \cdot d^{\pi', s}(s') \cdot A^\pi(s', \pi')$$

Proof of policy improvement:

$$V^{\pi_k}(s) - V^{\pi_{k-1}}(s) = \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d^{\pi_k, s}} [A^{\pi_{k-1}}(s, \pi_k)]$$

$$A^{\pi_{k-1}}(s, \pi_k) = Q^{\pi_{k-1}}(s, \pi_k) - \underbrace{V^{\pi_{k-1}}(s)}_{Q^{\pi_{k-1}}(s, \pi_{k-1})}$$

Claim:  $\downarrow \geq 0$ .

recall:  $\pi_k$  is greedy w.r.t.  $Q^{\pi_{k-1}}$ .

$$\therefore \underbrace{Q^{\pi_{k-1}}(s, \pi_k)}_{\Delta} \geq Q^{\pi_{k-1}}(s, \pi_{k-1}) = V^{\pi_{k-1}}(s).$$

**Strict improvement.**  $\pi_{k-1} \neq \pi^*$ , then  $\exists s, V^{\pi_k}(s) > V^{\pi_{k-1}}(s)$

$$(1) \exists s, A^{\pi_{k-1}}(s, \pi_k) > 0.$$

Proof: if  $\pi_{k-1} \neq \pi^*$ .  $\exists s_0, V^{\pi^*}(s_0) - V^{\pi_{k-1}}(s_0) > 0$ .

$$\Rightarrow \underbrace{\mathbb{E}_{s' \sim d^{\pi^*, s_0}} [A^{\pi_{k-1}}(s', \pi^*)]}_{\Delta} > 0.$$

$$\exists \tilde{s} : A^{\pi_{k-1}}(\tilde{s}, \pi^*) > 0.$$

$$\underbrace{A^{\pi_{k-1}}(\tilde{s}, \pi_k)}_{\Delta} \geq A^{\pi_{k-1}}(\tilde{s}, \pi^*) > 0.$$

$$\begin{aligned} V^{\pi_k}(\tilde{s}) - V^{\pi_{k-1}}(\tilde{s}) &= \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d^{\pi_k, \tilde{s}}} [A^{\pi_{k-1}}(s', \pi_k)] \\ &\geq A^{\pi_{k-1}}(\tilde{s}, \pi_k) > 0. \end{aligned}$$

# Exponential convergence to $\epsilon$ -optimality.

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

Proof: Key facts. (to be established). <sup>"pointwise"</sup>

(a)  $J^{\pi_{k+1}} Q^{\pi_k} \geq J^{\pi} Q^{\pi_k} \quad \forall \pi.$

(b)  $J^{\pi_{k+1}} Q^{\pi_k} \leq Q^{\pi_{k+1}}$

$$\underbrace{Q^* - Q^{\pi_{k+1}}}_{\geq 0} = \underbrace{Q^* - J^{\pi_{k+1}} Q^{\pi_k}}_{\geq 0} + \underbrace{J^{\pi_{k+1}} Q^{\pi_k} - Q^{\pi_{k+1}}}_{\leq \epsilon}$$

$$\leq \underbrace{J^{\pi^*} Q^* - J^{\pi^*} Q^{\pi_k}}_{\geq 0}$$

$$\|Q^* - Q^{\pi_{k+1}}\|_{\infty} \leq \|J^{\pi^*} Q^* - J^{\pi^*} Q^{\pi_k}\|_{\infty} \leq \gamma \cdot \|Q^* - Q^{\pi_k}\|_{\infty}$$

To establish (a) & (b).

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

$$J^{\pi} Q^{\pi_k} \leq J^{\pi_{k+1}} Q^{\pi_k} \quad (Q^{\pi_k}(s, \pi_{k+1}) \geq Q^{\pi_k}(s, \pi) \quad \forall \pi)$$

better ↑

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

$$(J^{\pi} f)(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [f(s', \pi)]$$