

$$Q^* = \mathcal{T}Q^*. \quad \forall f \in \mathbb{R}^{S \times A}, \quad (\mathcal{T}f)(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} f(s')$$

Alg: repeatedly apply \mathcal{T} to arbitrary f_0 $f(s', \arg \max_{a'} f(s', a))$

$$\text{Contraction: } \forall f, f' \in \mathbb{R}^{S \times A}, \quad \|\mathcal{T}f - \mathcal{T}f'\|_\infty \leq \gamma \|f - f'\|_\infty.$$

$$\text{Type 1: } Q^\pi, Q^* \quad \sum_{t=1}^{\infty} \gamma^{t-1} r_t \in [0, \frac{\max}{1-\gamma}].$$

$$\text{Type 2: } \forall f \in \mathbb{R}^{S \times A}.$$

Fact: $\exists f \in \mathbb{R}^{S \times A}$, where f is not Q^π for any π .

Example: MDP M, $R \equiv 0$. $\Rightarrow Q^\pi \equiv 0, Q^* \equiv 0$.

Value Iteration for V^* , Q^π , V^π .

$$V^* = \mathcal{T}V^*. \quad \forall f \in \mathbb{R}^S, \quad (\mathcal{T}f)(s) = \max_{a \in A} (R(s, a)$$

$$+ \gamma \mathbb{E}_{s' \sim P(s, a)} [f(s')]).$$

$$Q^\pi = \mathcal{T}^\pi Q^\pi. \quad \forall f \in \mathbb{R}^{S \times A}, \quad (\mathcal{T}^\pi f)(s, a) = R(s, a) +$$

$$\gamma \mathbb{E}_{s' \sim P(s, a)} [f(s', \pi(s'))]$$

$$\forall f, f' \in \mathbb{R}^{S \times A}, \quad \|\mathcal{T}^\pi f - \mathcal{T}^\pi f'\|_\infty \leq \gamma \|f - f'\|_\infty.$$

$$V^\pi = \mathcal{T}^\pi V^\pi, \quad \forall f \in \mathbb{R}^S, \quad (\mathcal{T}^\pi f)(s) = R(s, \pi(s)) +$$

$$\gamma \mathbb{E}_{s' \sim P(s, \pi(s))} [f(s')].$$

Several ways to calculate $V^\pi(s)$.

- $V^\pi(s) = \underbrace{e_s^\top}_{\text{Eq}} (I - \gamma P^\pi)^{-1} R^\pi.$
- Init $f_0 \in \mathbb{R}^S$. $f_k \leftarrow \mathcal{T}^\pi f_{k-1}$. $\xrightarrow{\text{converge to}} V^\pi$.
- (Monte-Carlo): start from s , rollout actual trajectories. take average over random returns. $\sum_{t=1}^{\infty} \gamma^{t-1} r_t$.

Alt. proof for convergence of Value Iter.

$$V^* = \mathcal{T} V^*, \quad f_0 \in \mathbb{R}^S, \quad f_k \leftarrow \mathcal{T} f_{k-1}.$$

$$f_0 = \vec{0}.$$

$$f_1 = \mathcal{T} \vec{0} \Rightarrow f_1(s) = \max_a R(s, a).$$

$$f_2 = \mathcal{T} f_1 \Rightarrow f_2(s) = \max_a \left(\underbrace{R(s, a)}_{\Delta} + \gamma \mathbb{E}_{s' \sim p(s, a)} \left[\max_a R(s', a) \right] \right)$$

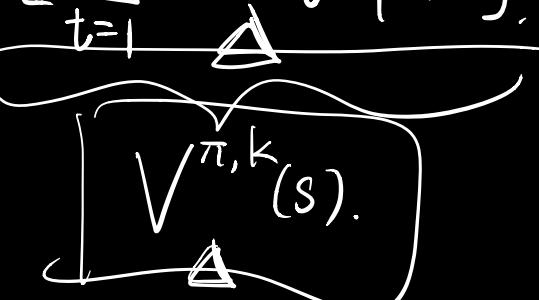
Generally: $f_k(s) = \max_{\text{non-stationary } \pi} \mathbb{E} \left[\sum_{t=1}^k \gamma^{t-1} r_t \mid \pi \right]$.

Claim: $\|f_k - V^*\|_\infty \leq \frac{\gamma^k R_{\max}}{1-\gamma}$.

Proof sketch: $f_k \leq V^*$.

$$f_k(s) = \max_{\text{non-stationary } \pi} V^{\pi, k}(s) \geq V^{\pi^*, k}(s).$$

$$= \mathbb{E} \left[\sum_{t=1}^k \gamma^{t-1} r_t \mid \pi^*, s_1 = s \right].$$



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

$$- \mathbb{E} \left[\sum_{t=k+1}^{\infty} \gamma^{t-1} r_t \mid \pi^*, s_1 = s \right]. \leftarrow$$

$\overbrace{\quad\quad\quad}^{\pi} \quad \underbrace{[0, R_{\max}]}_{\mathbb{R}}.$

$$\geq V^*(s) - \underbrace{\sum_{t=k+1}^{\infty} \gamma^{t-1} R_{\max}}_{\mathbb{R}}$$

$$= V^*(s) - \gamma^k \underbrace{\sum_{t=k+1}^{\infty} \gamma^{t-k-1} R_{\max}}_{(-\gamma)} = \frac{R_{\max}}{1-\gamma}.$$

$$= V^*(s) - \frac{\gamma^k R_{\max}}{1-\gamma}.$$

Together: $\forall s, V^*(s) - \frac{\gamma^k R_{\max}}{1-\gamma} \leq f_k(s) \leq V^*(s)$

$$\Rightarrow \|V^* - f_k\|_{\infty} \leq \frac{\gamma^k R_{\max}}{1-\gamma}.$$

if $f = Q^*$, $\pi_f(s) = \arg \max_{a \in A} f(s, a)$. \rightarrow optimal.

In value iter, we get $f_k \approx Q^*$.

output π_{f_k} . Q: any guarantee for the optimality

A: $\forall f \in \mathbb{R}^{S \times A}$. [Singh & Yee '94]. \downarrow of π_{f_k} ?

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

$Q^{\pi_f} ? f$