

# How to solve for $V^*/Q^*$ | Value Iteration (VI)

Definition: Bellman optimality operator  $T: \mathbb{R}^{S \times A} \rightarrow \mathbb{R}^{S \times A}$   
 $\forall f \in \mathbb{R}^{S \times A}, (Tf) \in \mathbb{R}^{S \times A}$

$$(Tf)(s,a) = R(s,a) + \gamma \mathbb{E}_{s \sim P(\cdot|s,a)} \left[ \max_{a'} f(s,a') \right]$$

Bellman opt. eq. for  $Q^*$ :  $Q^* = TQ^*$   
( $Q^*$  is the fixed point of operator  $T$ ).

## VI for solving $Q^*$

- Initialize  $f_0$  arbitrarily (e.g.  $f_0 = \vec{0}$ )
- for  $k=1, 2, 3, \dots$ ,  $f_k = Tf_{k-1}$ .

Intuition:  
when  $f_{k-1} = Q^*$   
 $f_k = TQ^* = Q^*$

Convergence of VI:

$$\|f\|_{\infty} = \max_{s,a} |f(s,a)|$$

Proposition:  $T$  is a  $\gamma$ -contraction under  $\|\cdot\|_{\infty}$ . That is,

$$\forall f, f' \in \mathbb{R}^{S \times A}, \|Tf - Tf'\|_{\infty} \leq \gamma \cdot \|f - f'\|_{\infty}$$

Use contraction of  $T$  to show convergence of VI:

$$\|f_k - Q^*\|_{\infty} = \|Tf_{k-1} - TQ^*\|_{\infty}$$

$$\leq \gamma \cdot \|f_{k-1} - Q^*\|_{\infty} \leq \gamma^2 \cdot \|f_{k-2} - Q^*\|_{\infty} \leq \dots \leq \gamma^k \|f_0 - Q^*\|_{\infty} \\ \leq \gamma^k \cdot \frac{R_{\max}}{1-\gamma} \quad (\text{for } f_0 = \vec{0}).$$

Proof of contraction:  $\forall s, a$ .

$$\begin{aligned}
 & |(Tf)(s, a) - (Tf')(s, a)| \\
 &= |(\cancel{R(s, a)} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [\max_{a'} f(s', a')]) - (\cancel{R(s, a)} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [\max_{a'} f'(s', a')])| \\
 &= \gamma \left| \mathbb{E}_{s' \sim P(\cdot | s, a)} [\max_{a'} f(s', a') - \max_{a'} f'(s', a')] \right| \\
 &\leq \gamma \cdot \max_{s' \in S} \left| \max_{a'} f(s', a') - \max_{a'} f'(s', a') \right| \\
 &= \gamma \cdot \max_{s \in S} \left| \max_a f(s, a) - \max_a f'(s, a) \right|.
 \end{aligned}$$

(want to show)  $\leq \gamma \max_{s \in S} \max_{a \in A} |f(s, a) - f'(s, a)|$ .

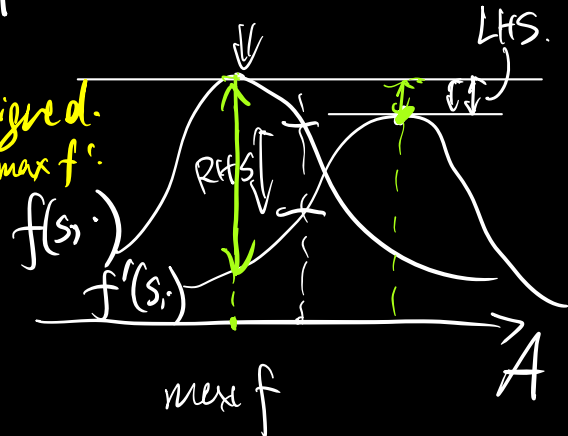
suffice to show:  $\forall s$ .

$$\left| \max_a f(s, a) - \max_a f'(s, a) \right|$$

actions not aligned.  
argmax  $f$  vs. argmax  $f'$ .

$$\leq \max_a |f(s, a) - f'(s, a)|.$$

aligned.

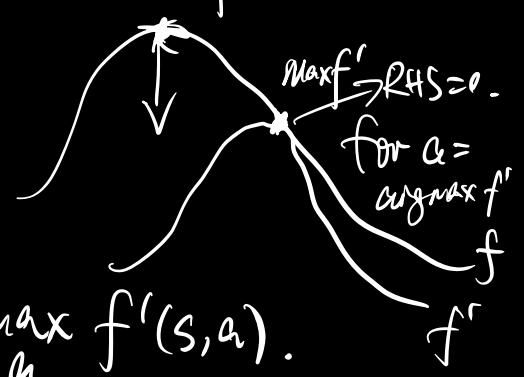


w.l.o.g.  $\max_a f(s, a) \geq \max_a f'(s, a)$ .

let  $\underline{a^*} = \text{argmax}_a f(s, a)$ .

$$\max_a f(s, a) - \max_a f'(s, a) = f(s, a^*) - \max_a f'(s, a).$$

$$\leq f(s, a^*) - f'(s, a^*) \leq \max_a |f(s, a) - f'(s, a)|.$$



Remarks:  $\|Q^* - f\|_{\infty} \leq \gamma^k \cdot \|Q^* - f\|_{\infty}$ . b/c  $T$  is contraction.

There are other "VZ" algorithms.

$$\begin{aligned} \cdot V^*(s) &= \max_a (R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [V^*(s')]) \\ &:= (T V^*)(s). \end{aligned}$$

$\xrightarrow{\text{abuse of notation}}$

VZ for  $V^*$ : init  $f_0 \in \mathbb{R}^S$ .  $f_k \leftarrow T f_{k-1}$ .

$$\cdot V^\pi = R^\pi + \gamma P^\pi V^\pi$$

Define  $T^\pi: \mathbb{R}^S \rightarrow \mathbb{R}^S$ .  $\forall f \in \mathbb{R}^S$ .

Bellman operator  
for policy  $\pi$ .

$$T^\pi f = R^\pi + \gamma P^\pi f. \Rightarrow V^\pi = T^\pi V^\pi.$$

$$\text{Also for } Q^\pi: Q^\pi = T^\pi Q^\pi. \quad \uparrow \text{ abuse of notation}$$

VZ for  $Q^*$ : as  $k \uparrow$ ,  $\|f_k - Q^*\|_{\infty} \downarrow$ .

We know:  $\pi_{\underline{Q^*}} = \pi^*$ .

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

But we only have  $f_k \approx Q^*$ .

How to induce a good policy? Naive:  $\pi_{\underline{f_k}}$ .

Q: Is this near-optimal?

A: YES!

Lemma:  $\forall f \in \mathbb{R}^{S \times A}$ .  $\|V^* - V^{\pi_f}\|_\infty \leq \frac{2 \|f - Q^*\|_\infty}{1 - \gamma}$ .

Proof:  $\forall s \in S$ .

$$V^*(s) - V^{\pi_f}(s) = \underbrace{Q^*(s, \pi^*) - Q^*(s, \pi_f)}_{(i)} + \underbrace{Q^*(s, \pi_f) - Q^{\pi_f}(s, \pi_f)}_{(ii)}$$

"real":  $V^\pi, Q^\pi, V^*, Q^*$   
 "arbitrary":  $f, f'$

$$\leq \underbrace{Q^*(s, \pi^*) - f(s, \pi^*)}_{(i)} + \underbrace{f(s, \pi_f) - Q^*(s, \pi_f)}_{\geq 0} - Q^*(s, \pi_f)$$

actions aligned

$$\left( \cancel{R(s, \pi_f)} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi_f)} [V^*(s')] \right) - \left( \cancel{R(s, \pi_f)} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi_f)} [V^{\pi_f}(s')] \right)$$

(ii)

$$\leq \underbrace{2 \|f - Q^*\|_\infty}_{(i)} + \underbrace{\gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi_f)} [V^*(s') - V^{\pi_f}(s')]}_{(ii)}$$

$$\leq 2 \cdot \|f - Q^*\|_\infty + \gamma \cdot \max_{s'} (V^*(s') - V^{\pi_f}(s'))$$

non-negative

$$= 2 \cdot \|f - Q^*\|_\infty + \gamma \cdot \|V^* - V^{\pi_f}\|_\infty$$

Therefore:

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

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

Remarks:

$$1. \|f_k - Q^*\|_\infty \leq \gamma^k \|f_0 - Q^*\|_\infty \leq \gamma^k \frac{R_{\max}}{1-\gamma} \quad \text{|| } V_{\max} \text{}$$

$\Rightarrow$  if we want  $\|f_k - Q^*\|_\infty \leq \epsilon$ , how large should  $k$  be?

$$\gamma^k V_{\max} \leq \epsilon \Rightarrow \gamma^k \leq \frac{\epsilon}{V_{\max}}$$

$$k \geq \log_{1/\gamma} \frac{V_{\max}}{\epsilon} \Rightarrow k \geq \frac{\ln \frac{V_{\max}}{\epsilon}}{\ln \frac{1}{\gamma}}$$

$$\frac{\ln \frac{1}{\gamma}}{\ln \frac{1}{\gamma}} \approx 1-\gamma = O\left(\frac{1}{1-\gamma}\right)$$

"effective horizon"

2. Stopping criterion:  $\|f_k - f_{k+1}\|_\infty \leq \epsilon'$

Why?  $\|f_k - f_{k+1}\|_\infty = \|f_k - \mathcal{T}f_k\|_\infty$

Bellman error/residual of  $f_k$ .  
 $Q^* = \mathcal{T}Q^*$

$$\|f_k - Q^*\|_\infty$$

$$= \|f_k - \mathcal{T}f_k + \mathcal{T}f_k - \mathcal{T}Q^*\|_\infty$$

$$\leq \|f_k - \mathcal{T}f_k\|_\infty + \|\mathcal{T}f_k - \mathcal{T}Q^*\|_\infty$$

$$\leq \epsilon' + \gamma \cdot \|f_k - Q^*\|_\infty$$

$$\Rightarrow \|f_k - Q^*\|_\infty \leq \frac{\epsilon'}{1-\gamma}$$

3. Similar algorithm/analyses for  $V^*$ ,  $V^\pi$ ,  $Q^\pi$

$$Q^\pi = \mathcal{T}^\pi Q^\pi \rightarrow \text{"VI" alg: } f_k \leftarrow \mathcal{T}^\pi f_{k-1}$$

# Alternative Proof for convergence of VI

"truncated value-func"  $V^{\pi, H}(s) = \mathbb{E} \left[ \sum_{t=1}^H \gamma^{t-1} r_t \mid s_1 = s, \pi \right]$ .  
(in fact,  $V^{\pi}(s) = V^{\pi, \infty}(s)$ ).

Consider VI for solving  $V^*$ .

Init  $f_0 = \vec{0} \in \mathbb{R}^S$ .  $f_k \leftarrow \mathcal{T} f_{k-1}$ .

Claim:  $f_k(s) = \max_{\pi} V^{\pi, k}(s) =: V^{*, k}(s)$ .  
 $\pi \leftarrow$  "all policies, possibly non-stationary."

Examples:  $f_k(s) = \max_a (R(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} [f_{k-1}(s')])$

$k=0$ .  $f_0 = V^{*, 0} = \vec{0}$ .

$k=1$ .  $f_1(s) = \max_a R(s, a) = V^{*, 1}(s)$

$k=2$ .  $f_2(s) = V^{*, 2}(s) =$

$\max_a (R(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} [\max_{a'} R(s', a')])$

$k=3, \dots$

Remark: optimal policy for  $V^{*, H}$  is non-stationary.  
depends on time step.

Want to upper bound:  $\|f_k - V^*\|_\infty$

$$f_k = V^{*,k} \quad V^{*,k} - V^*$$

Proof: (1)  $V^{*,k} \leq V^*$

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

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

$V^{*,k} \leq V^*$  follows directly from  $r_t \geq 0$ .

$$(2) \quad V^{*,k} \geq V^* - \gamma^k V_{\max} = \frac{R_{\max}}{1-\gamma}$$

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

→ optimal for infinite horizon

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

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

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

$$= V^*(s) - \gamma^k \mathbb{E} \left[ \sum_{t'=k+1}^{\infty} \gamma^{t'-k-1} r_{t'} \mid s_1 = s, \pi^* \right]$$

↓  $t' := t+k$

$$\geq V^*(s) - \gamma^k \cdot V_{\max}$$

$$\sum_{t'=1}^{\infty} \gamma^{t'-1} r_{t'+k} \mid s_{1-k} = s, \pi^*$$

ignore

$$V^*(s) - \gamma^k V_{\max} \leq V^{*,k}(s) \leq V^*(s).$$

$$f_k(s)$$

$$\Rightarrow \|V^* - f_k\|_{\infty} \leq \gamma^k V_{\max}$$

(same guarantee as before)  $\leq \frac{R_{\max}}{1-\gamma}$

$R_{\max}$   
 $\sum_{t=1}^{\infty} \gamma^{t-1} R_{\max}$   
 $= \frac{1}{1-\gamma} R_{\max}$

Remark: for  $Q^*$ , we have.

$$\|V^* - V^{\pi_{f_k}}\|_{\infty} \leq \frac{2 \cdot \|f_k - Q^*\|}{1-\gamma} \leq \frac{2 \cdot \gamma^k V_{\max}}{1-\gamma}$$

Turns out you can save  $\frac{2}{1-\gamma}$  factor!

Trick: instead of outputting  $\pi_{f_k}$ .

output instead  $\pi_{NS} = \pi_{f_k} \circ \pi_{f_{k-1}} \circ \pi_{f_{k-2}} \circ \dots \circ \pi_{f_1}$

"non-stay"

$\Leftrightarrow a_1 \sim \pi_{f_k}, a_2 \sim \pi_{f_{k-1}}, a_3 \sim \pi_{f_{k-2}}, \dots$

$a_k \sim \pi_{f_1}, a_{k+1} = \infty \sim \text{arbitrary}.$

Claim: this policy optimizes  $V^{*,k}$ .

$$V^{\pi_{NS}} \geq V^{\pi_{NS,k}} = V^{*,k} \geq V^* - \gamma^k R_{\max}$$

$$\Rightarrow \|V^* - V^{\pi_{NS}}\|_{\infty} \leq \gamma^k R_{\max}. \quad \triangleleft$$



In practice,  $\gamma$  is often introduced for convenience.

Given finite-horizon undiscounted problem (w/  $H$ )  $\leftarrow$   
can turn into inf-horizon discounted,  $H \approx \underline{O\left(\frac{1}{1-\gamma}\right)}$