① Recap / Value Iteration

$\Rightarrow$ compute $\pi$ , $\varepsilon > 0$ :

$v^{\pi} \geq v^* - \varepsilon \mathbb{1}$ .

---

② Policy Iteration $\boxed{\tilde{O}\left(\frac{SA}{1-\gamma}\right)}$

New!

---

③ Comp. complexity of

planning in finite MDPs?

Lower bound

---

① $M = (S, A, P, r, \gamma)$  $0 \leq \gamma < 1$

$\pi <$ general
    memoryless

Fundamental Theorem
---

1.     $\pi$    greedy   w.r.t. $v^*$

$\left[ \begin{array}{l} \boxed{T_{\pi} v^* = T v^*} : \\ T_{\pi} v = r_{\pi} + \gamma P_{\pi} v \\ T_a v = r_a + \gamma P_a v \\ (Tv)(s) = \max\limits_{a} (T_a v)(s) \end{array} \right.$

$\Rightarrow v^{\pi} = v^*$

2.   $T v^* = v^*$

$$\tilde{v}^*(s) = \sup_{\pi \in ML} v^\pi(s)$$

Part 1:   $v^* \leftarrow \tilde{v}^*$, prove thm

Part 2:   $v^* = \tilde{v}^*$

---

Part 1:   $\pi$ is **greedy**

w.r.t. $\tilde{v}^* \Rightarrow v^\pi = \tilde{v}^*$. ?

$v^\pi \leq \tilde{v}^*$ ✓

$\tilde{v}^* \leq_? v^\pi$

$\boxed{\tilde{v}^* \leq T\tilde{v}^*}$ ?

$\forall \pi \; ML: \quad v^\pi \leq^? T\tilde{v}^* \quad /\sup_{\pi ML}$

$\Rightarrow \tilde{v}^* \leq T\tilde{v}^*$

---

$\pi \quad ML$

$v^\pi = T_\pi v^\pi$

$v^\pi \leq \tilde{v}^* \qquad / T_\pi$

$T_\pi v^\pi \leq T_\pi \tilde{v}^* \leq T\tilde{v}^*$

$\Rightarrow v^\pi \leq T\tilde{v}^*$ ✓

---

Take $\pi$ greedy wrt $\tilde{v}^*$:

$$T_\pi \tilde{v}^* = T\tilde{v}^* \geq \tilde{v}^* \qquad / T_\pi$$

$$T_\pi^2 \tilde{v}^* \geq T_\pi \tilde{v}^* \geq \tilde{v}^*$$

$\boxed{a_k \geq a \Rightarrow \lim_{k \to \infty} a_k \geq a}$

Banach's fixed point theorem $\begin{cases} \end{cases}$

$\underline{T_\pi^k \tilde{v}^*} \geq \tilde{v}^*$

$\downarrow k \to \infty$

$v^\pi \geq \tilde{v}^*$

// Qu.e.d.

part 1 fT

$$v^* = \tilde{v}^*, \quad \tilde{v}^* \leq T\tilde{v}^*$$
$$v^* \leq Tv^*$$

$\pi$ grady w.r.t. $v^*$:

$$v_\pi = T_\pi v_\pi = T_\pi v^* = Tv^*$$
$$\parallel$$
$$v^*$$

// av. ed.

# Value-iteration

$$v_0 \in \mathbb{R}^S ; \quad v_0 = 0$$
$$v_{k+1} = T v_k \quad \Rightarrow O(S \times SA)$$

Banach's FP Thm $\Rightarrow$

$$\| v_k - v^* \|_\infty \leq \gamma^k \| v_0 - v^* \|$$

$v_0 = 0$

$$= \gamma^k \| v^* \|_\infty \leq \frac{\gamma^k}{1-\gamma} \leq \varepsilon$$

$\uparrow$

$r_a(s) \in [0,1]$

$$\Rightarrow 0 \leq \underbrace{\sum \gamma^t r_{A_t}(s_t)}_{1} \leq \frac{1}{1-\gamma}$$

$$k \geq \frac{\log\left(\frac{1}{\varepsilon(1-\gamma)}\right)}{\log\left(\frac{1}{\gamma}\right)}$$

$$k \geq \boxed{\frac{\log\left(\frac{1}{\varepsilon(1-\gamma)}\right)}{1-\gamma}} \geq \frac{\log\left(\frac{1}{\varepsilon(1-\gamma)}\right)}{\log\left(\frac{1}{\gamma}\right)}$$

$$H_{\gamma,\varepsilon}$$

$$\log x \leq x - 1$$
$$x > 0$$

$$k: \quad v_k \geq v^* - \varepsilon 1$$

$$\boxed{\|v_k - v^*\|_\infty \leq \varepsilon} \Rightarrow v_k \geq v^* - \varepsilon 1)$$

---

$$\underline{v \geq v^* - 1\varepsilon} \qquad 1(s) = 1 \ \forall s$$

Greedify!

$$\pi: \quad \boxed{T_\pi v = Tv}$$

Is $\pi$ good policy?

$$T_\pi v \geq T_\pi(\underbrace{v^* - 1\varepsilon})$$

$$T_\pi(\underline{v} + \underline{a1}) \qquad a \in \mathbb{R}$$

Dead-end...

$$= r_\pi + \gamma P_\pi(\overrightarrow{v + a1})$$

$$= \underbrace{r_\pi + \gamma P_\pi v} + \gamma a \underbrace{P_\pi 1}_{1}$$

$$= T_\pi v + \underline{\gamma a 1}$$

$$T(v + a1) = Tv + \underline{\gamma a 1}$$

---

$$v \geq v^* - 1\varepsilon \qquad /T$$

$$Tv \geq T(v^* - \varepsilon 1) = Tv^* - \gamma\varepsilon 1$$

$$\| \qquad = \underline{v^* - \gamma\varepsilon 1}$$

$$T_\pi v$$

---

$$\underline{T_\pi v \geq v^* - \gamma\varepsilon 1}$$

$$\geq v - (\gamma+1)\varepsilon 1 \ / T_\pi$$

$$\uparrow \qquad \varepsilon(\gamma^2 + \gamma + 1)1$$

$$v^* \geq v - 1\varepsilon$$

$$T_\pi^2 v \geq T_\pi v - \underline{\gamma(1+\gamma)\varepsilon 1} = Tv - \cdots$$

$$\geq v^* - \cdots - \gamma\varepsilon 1$$

$$T_\pi^2 v \geq v^* - \varepsilon(1 + \gamma + \gamma^2)\mathbf{1}$$

$$\vdots$$

$$T_\pi^k v \geq v^* - \varepsilon(1 + \gamma + \dots + \gamma^k)\mathbf{1}$$

$$/k \to \infty$$

$$\downarrow$$

$$v^\pi \geq v^* - \frac{\varepsilon}{1-\gamma}\mathbf{1}$$

---

Prop: $\pi$ is greedy w.r.t. $v$:

$$\|v - v^*\|_\infty \leq \varepsilon$$

$$\Rightarrow v^\pi \geq v^* - \frac{\textcircled{\varepsilon}}{1-\gamma}\mathbf{1}$$

'

$$\boxed{k \geq H_{\gamma,\ \varepsilon(1-\gamma)}}$$

$$\|v_k - v^*\| \leq \underline{\varepsilon(1-\gamma)}$$

$\pi_k$ greedy w.r.t. $v_k$

$$\Rightarrow v^{\pi_k} \geq v^* - \frac{\varepsilon(1-\gamma)}{1-\gamma}\cdot\mathbf{1}$$

Prop.

---

$$O\left(S^2 A \underbrace{H_{\gamma,\ \varepsilon(1-\gamma)}}_{\frac{\log\left(\frac{1}{\varepsilon}(1-\gamma)^2\right)}{1-\gamma}}\right)$$

$$\log\left(\frac{1}{\varepsilon}\right) ; \quad \frac{1}{1-\gamma}$$

Yinyu Ye
P.I.

$$\boxed{O\left(\text{poly }(S,A)/_{1-\gamma}\right)}$$

outputs $\pi$: $v^\pi = v^*$ !

No $\varepsilon$ !

---

## Policy Iteration

← determinstic

$\pi_0$   ML   arbitrarily

$k = 0, 1, \dots$

$\pi_{k+1}$ : $T_{\pi_{k+1}} v^{\pi_k} = T v^{\pi_k}$

$\overline{\pi_0}$

SA-A

S



---

Computation in round $\ell = 0, 1, \dots$ :

$v^{\pi_\ell} = ?$   /   $T_{\pi_\ell}^i v^{\overline{\pi_{\ell-1}}} \dots$

$$\begin{cases} T_{\overline{\pi_\ell}} v^{\overline{\pi_\ell}} = v^{\overline{\pi_\ell}} \\ r_{\overline{\pi_\ell}} + \gamma P_{\overline{\pi_\ell}} v^{\overline{\pi_\ell}} = v^{\pi_\ell} \\ v^{\pi_\ell} = \left(I - \gamma P_{\overline{\pi_\ell}}\right)^{-1} r_{\pi_k} \end{cases}$$

$S \times S$

---

1.   $\|v^{\pi_k} - v^*\|_\infty \leq \gamma^k \|v^{\pi_0} - v^*\|_\infty$

$\exists s_0 \in S$ ; assuming $\pi_0$ not optimal

2.   After $k^* = \tilde{O}\left(\frac{1}{1-\gamma}\right)$

$\forall k \geq k^*$   $\pi_k(s_0) \neq \pi_0(s_0)$

$$V^* \geq V^{\pi_{k+1}} \geq TV^{\pi_k} \geq T^k V^{\pi_0}$$

$$O(S^3) \lor O(S^2 A)$$

per iteration
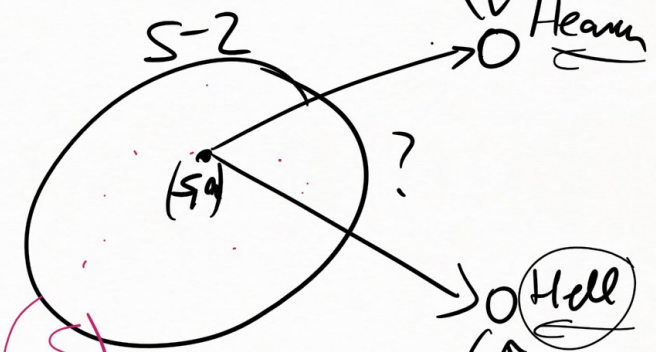
After $\tilde{O}\left(\dfrac{SA}{1-\gamma}\right)$

$\pi_k: \quad V^{\pi_k} = V^*$

No $\log\left(\dfrac{1}{\varepsilon}\right)$

$\Omega(S^2 A)$

$A^S$

$S-2$

1

○ Heaven

$(s,a)$ ?

○ Hell

$\Omega(S)$

○ 0

V.I

$\varnothing$

$(s,a)$

1 2 . . . . . . . . S