

Policy Iteration

Local planning / Online planning

P.I

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

π_0 ML

π_1

π_2

⋮

π_{k+1} greedy w.r.t v^{π_k}

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

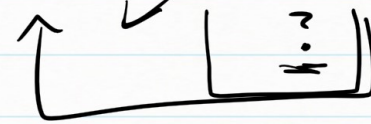
$$\pi_{k+1}(s) = \underset{a}{\operatorname{argmax}} \underbrace{r_a(s) + \gamma \langle P_a(s), v^{\pi_k} \rangle}_{\text{sys. tie resolution}}$$

$\forall s$

sys. tie resolution

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

Proof: Claim: $v^* \geq v^{\pi_k} \geq T^k v^{\pi_0}$



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

\uparrow
 $v^* = T v^*$

Need: $v^{\pi_k} \geq T^k v^{\pi_0}$

$$k-1 \Rightarrow k \quad \begin{matrix} \max \\ \downarrow \end{matrix} \quad \begin{matrix} \pi_{k-1} \\ \downarrow \end{matrix}$$

$$T_{\pi_k} v^{\pi_{k-1}} = T v^{\pi_{k-1}} \geq T_{\pi_{k-1}} v^{\pi_{k-1}}$$

$$T_{\pi_k}^2 v^{\pi_{k-1}} \geq T_{\pi_k} v^{\pi_k} = T v^{\pi_{k-1}} \geq v^{\pi_{k-1}}$$

$$v^{\pi_k} \leftarrow T_{\pi_k}^i v^{\pi_{k-1}} \geq T v^{\pi_{k-1}} \geq T^2 v^{\pi_{k-2}} \quad // \text{Qu. e.}$$

Value-difference identity

$$v^{\pi'} = (I - \gamma P_{\pi'})^{-1} r_{\pi'}$$

π, π' ML

$$v^{\pi'} = r_{\pi'} + \gamma P_{\pi'} v^{\pi'} \quad \checkmark$$

$$v^{\pi'} = T_{\pi'} v^{\pi'}$$

$$v^{\pi'} - v^{\pi} = (I - \gamma P_{\pi'})^{-1} [r_{\pi'} - (I - \gamma P_{\pi'}) v^{\pi}]$$

$$= (I - \gamma P_{\pi'})^{-1} [T_{\pi'} v^{\pi} - v^{\pi}]$$

$g(\pi', \pi)$

[advantage of π' relative to π]

$\pi_0 \downarrow$ v^* \downarrow

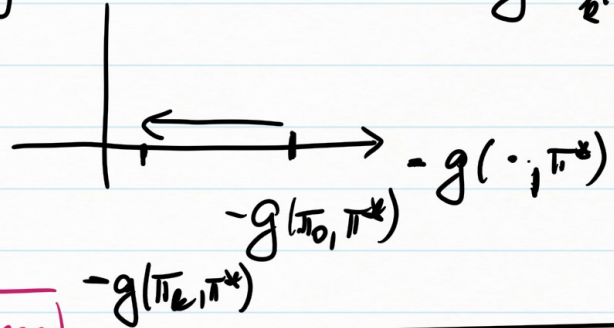
$$v^{\pi'} - v^{\pi} = (I - \gamma P_{\pi'})^{-1} g(\pi', \pi)$$

$\uparrow \pi_k$ $\uparrow \pi^*$

$\pi = \pi^*$ an opt. ML policy

$$\forall \pi': g(\pi', \pi^*) \leq 0$$

Progress measure: $-g(\pi_k, \pi^*)$



Progress

Lemma:

$$\forall \pi_0 \quad v^{\pi_0} \neq v^* \quad \{ \pi_k \} \text{ PI}$$

$$\exists s_0 \in S \text{ s.t. } \forall k \geq k^*(\delta)$$

$$\pi_k(s_0) \neq \pi_0(s_0)$$

Proof:

$$-g(\pi_k, \pi^*) = (I - \gamma P_{\pi_k})(v^* - v^{\pi_k})$$

$$= v^* - v^{\pi_k} - \gamma P_{\pi_k}(v^* - v^{\pi_k})$$

$\underbrace{\qquad}_{\geq 0} \quad \underbrace{\qquad}_{\geq 0}$

$$\leq v^* - v^{\pi_k}$$

$$\|g(\pi_k, \pi^*)\|_{\infty} \leq \|v^* - v^{\pi_k}\|_{\infty} \leq \gamma^k \|v^* - v^{\pi_0}\|_{\infty}$$

$$\|g(\pi_k, \pi_*)\|_\infty \leq \gamma^k \|v^* - v^{\pi_0}\|_\infty =$$

$$(I - \gamma P_{\pi_0})^{-1} g(\pi_0, \pi_*) = v^{\pi_0} - v^*$$

$$\rightarrow = \gamma^k \|(I - \gamma P_{\pi_0})^{-1} g(\pi_0, \pi_*)\|_\infty$$

$$\leq \gamma^k \|(I - \gamma P_{\pi_0})^{-1}\|_\infty \|g(\pi_0, \pi_*)\|_\infty$$

$$\|\sum_{i=0}^{\infty} \gamma^i P_{\pi_0}^i\|_\infty \leq \sum_{i=0}^{\infty} \gamma^i = \frac{1}{1-\gamma}$$

$$\leq \frac{\gamma^k}{1-\gamma} \|g(\pi_0, \pi_*)\|_\infty = \frac{\gamma^k}{1-\gamma} (-g(\pi_0, \pi_*)(s_0)) > 0$$

✓ g expand

$$(I - A)^{-1} = \sum_{i=0}^{\infty} A^i$$

$$\frac{1}{1-a} = \sum_{i=0}^{\infty} a^i$$

$$k \geq k^*(\gamma) \rightarrow$$

$$\frac{\gamma^k}{1-\gamma} < 1$$

$$\Rightarrow [v^* - T_{\pi_0} v^*](s_0) > [v^* - T_{\pi_k} v^*](s_0)$$

$$-g(\pi_0, \pi_*)(s_0) = \|g(\pi_0, \pi_*)\|_\infty > 0$$

$\pi_0 \neq \pi^*$

Claim: s_0 suits the ball!

$$\frac{\gamma^k}{1-\gamma} (-g(\pi_0, \pi_*)(s_0)) \geq \|g(\pi_k, \pi_*)\|_\infty$$

$$\geq -g(\pi_k, \pi_*)(s_0)$$

$$\frac{\gamma^k}{1-\gamma} [v^* - T_{\pi_0} v^*](s_0)$$

$$\geq [v^* - T_{\pi_k} v^*](s_0)$$

$$k \geq k^*(\gamma)$$

$$(T_{\pi_k} v^*)(s_0) > (T_{\pi_0} v^*)(s_0)$$

$$r_{\pi_k(s_0)} + \gamma < P_{\pi_k(s_0), v^*}$$

$$> r_{\pi_0(s_0)} + \gamma < P_{\pi_0(s_0), v^*}$$

$$\Rightarrow \pi_k(s_0) \neq \pi_0(s_0)$$

$$k^*(\gamma) \quad \frac{\gamma^k}{1-\gamma} < 1$$

$$k^*(\gamma) := \lceil H_{\gamma, 1} \rceil = \left\lceil \frac{\log\left(\frac{1}{1-\gamma}\right)}{\log\left(\frac{1}{\gamma}\right)} \right\rceil = \left\lceil \frac{\log\left(\frac{1}{1-\gamma}\right)}{1-\gamma} \right\rceil$$

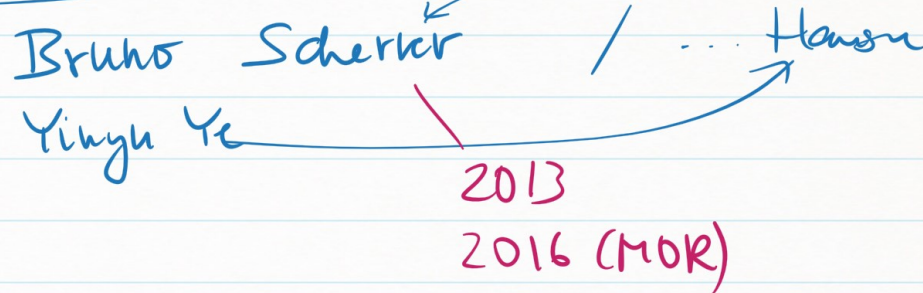
s^A	1	...	A
1			x
...		x	
...			
S			

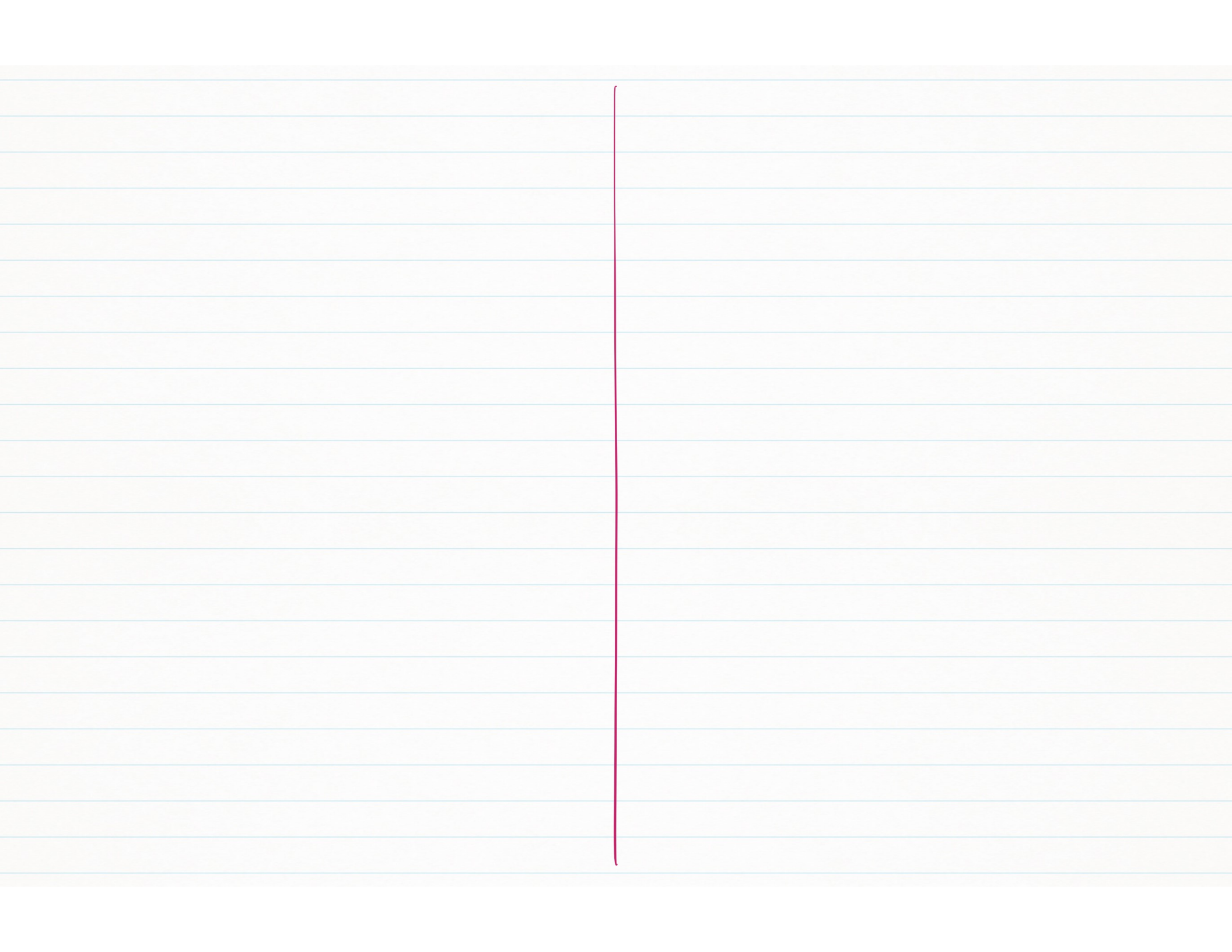
k^*

$SA - S$
Total # SA pairs

Thm: $k \geq \lceil \frac{SA-S}{1-\gamma} \log\left(\frac{1}{1-\gamma}\right) \rceil$

$$\pi_k = \pi^*$$





1

