

How much worse is policy π over expert π^* ?

$$J(\pi) - J(\pi^*)$$

PERFORMANCE
(TOTAL COST)

WHY IS THIS HARD?

$$J(\pi) - J(\pi^*)$$

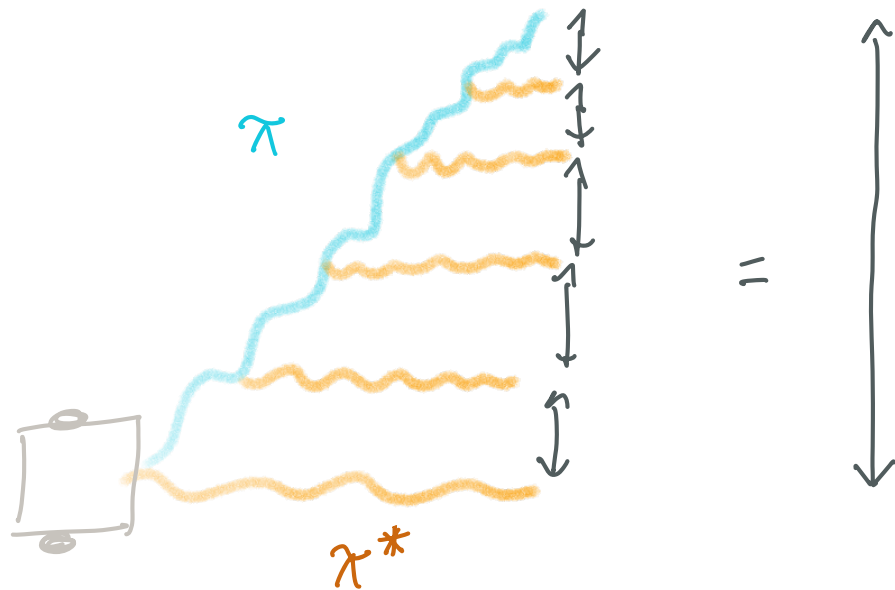
$$= \sum_{t=0}^{T-1} E_{\substack{s_t \sim d_t^\pi \\ a_t \sim \pi}} [c(s_t, a_t)] - \sum_{t=0}^{T-1} E_{\substack{s_t \sim d_t^{\pi^*} \\ a_t \sim \pi^*(s_t)}} [c(s_t, a_t)]$$

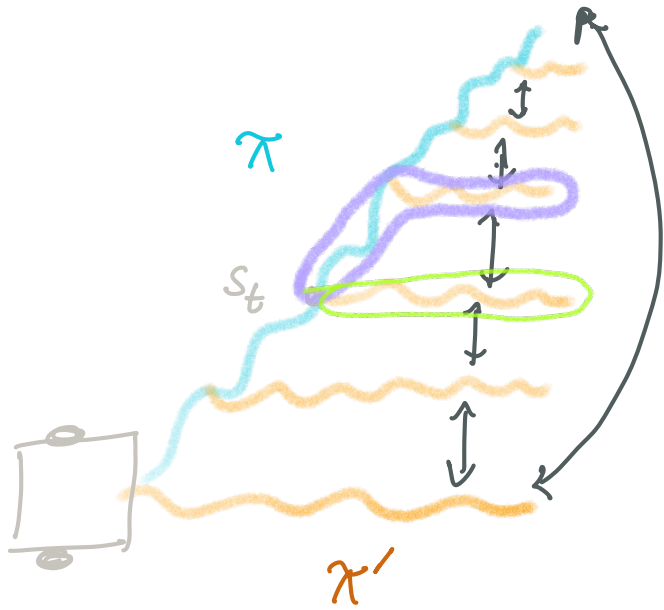
$$= \sum_{t=0}^{T-1} \sum_{s_t, a_t} c(s_t, a_t) \left[P(s_t, a_t | \pi) - P(s_t, a_t | \pi^*) \right]$$

???

CHALLENGE: DIFFERENCE IN DISTRIBUTIONS HARD TO ANALYZE

$$J(\lambda) - J(\lambda^*)$$





PERFORMANCE	DIFFERENCE	LEMMA
	$J(\pi) - J(\pi^*)$	
$= \sum_{t=0}^{T-1} E_{s_t \sim d_t^{\pi}}$	$\left[\underline{Q}(s_t, \pi) - \underline{Q}(s_t, \pi^*) \right]$	
ROLLOUT LEARNER TILL TIMESTEP t	DIFFERENCE IN Q -VALUE OF LEARNER - EXPERT "ADVANTAGE"	$A^{\pi^*}(s_t, \pi)$

min π

$$\sum_{t=0}^{T-1} E_{s_t \sim d_t^\pi}$$

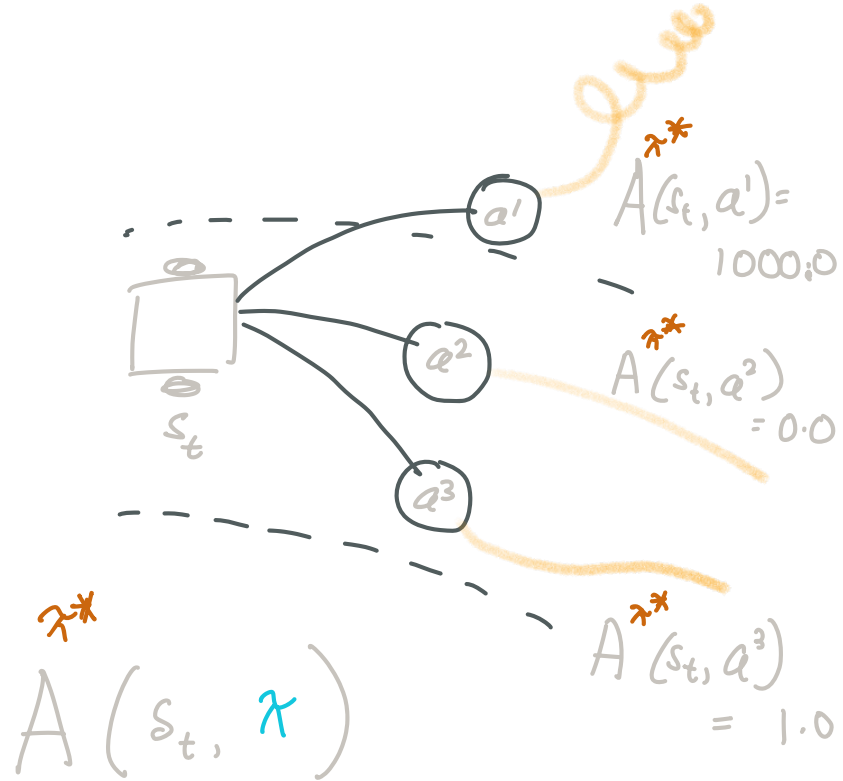
DISTRIBUTION

INTUITION

$$A(s_t, \pi)$$

Loss

STATES VISITED
By π

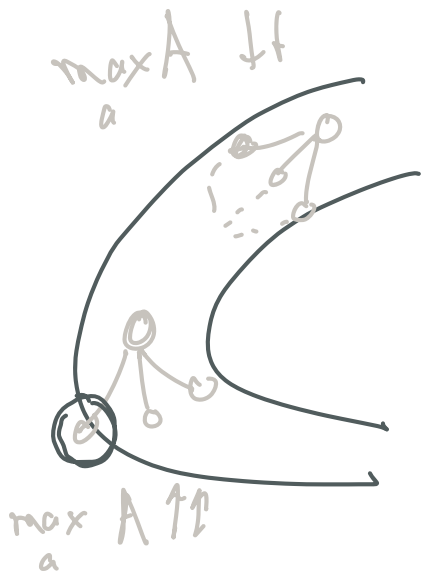


$$= \begin{bmatrix} \pi(a^1) \\ \pi(a^2) \\ \pi(a^3) \end{bmatrix} \cdot \begin{bmatrix} A(a^1) \\ A(a^2) \\ A(a^3) \end{bmatrix} \begin{matrix} 1000.0 \\ 0 \\ 1.0 \end{matrix}$$

COST SENSITIVE CLASSIFICATION

T-1 π^*

$$J(\pi) - J(\pi^*) = \sum_{t=0}^{\infty} E_{s_t \sim d_{\pi}^t} A(s_t, \pi)$$



$$\leq \sum_{t=0}^{T-1} E_{s_t \sim d_{\pi}^t} \left[\max_a A(s_t, a) \right] \cdot \mathbb{1}(\pi(s_t) \neq \pi^*(s_t))$$

$$\leq \underbrace{\max_{s, a} A(s, a)}_M \sum_{t=0}^T E_{s_t \sim d_{\pi}^t} \mathbb{1}(\pi(s_t) \neq \pi^*(s_t))$$

$$\leq O(\epsilon T \cdot M)$$

APPLICATIONS

①

FUNDAMENTALS: POLICY ITERATION

$\pi^+ \equiv \arg \max_{\pi} Q^{\pi^-}(s, \pi) \quad \forall s \Rightarrow$ DOES THIS MONOTONICALLY IMPROVE?

$$J(\pi^+) - J(\pi^-) = \sum_{t=0}^{T-1} \mathbb{E}_{s_t, d_t^{\pi^+}} A^{\pi^-}(s_t, \pi^+) \geq 0 \rightarrow \text{YES!}$$

② ALGORITHMS IN IMITATION LEARNING / REINFORCEMENT LEARNING

$$\max_{\pi} J(\pi) - J(\pi^*) \equiv \max_{\pi} \sum_{t=0}^{T-1} \mathbb{E}_{s_t, d_t^{\pi}} A^{\pi^*}(s_t, \pi)$$

AGGREGATE (ROSS & BACHNILL '14)

ROLL IN LEARNER

③ OTHER LEMMAS: SIMULATION LEMMA

$$J_M(\bar{\pi}) - J_{M'}(\bar{\pi})$$

" PERFORMANCE IN
MODEL M "

" PERFORM IN
MODEL M' "

PROOF

$$J(\pi) - J(\pi') = E_{s_0} Q^{(\pi)}(s_0, \pi) - E_{s_0} Q^{(\pi')}(s_0, \pi')$$

$$= E_{s_0} Q^{(\pi)}(s_0, \pi) - E_{s_0} Q^{(\pi)}(s_0, \pi) + E_{s_0} Q^{(\pi)}(s_0, \pi) - E_{s_0} Q^{(\pi')}(s_0, \pi')$$

$$= \cancel{E_{s_0, a_0}^{(\pi)} r(s_0, a_0)} - \cancel{E_{s_0, a_0}^{(\pi)} r(s_0, a_0)} + E_{s_0} A^{(\pi')}(s_0, \pi)$$

$$+ \underbrace{E_{s_1, d_1}^{(\pi)} Q(s_1, \pi) - E_{s_1, d_1}^{(\pi)} Q(s_1, \pi')}_{\text{RECURSIVELY EXPAND THIS}}$$

$$= \sum_{t=0}^{T-1} E_{s_t, d_t}^{(\pi)} A^{(\pi')}(s_t, \pi)$$

