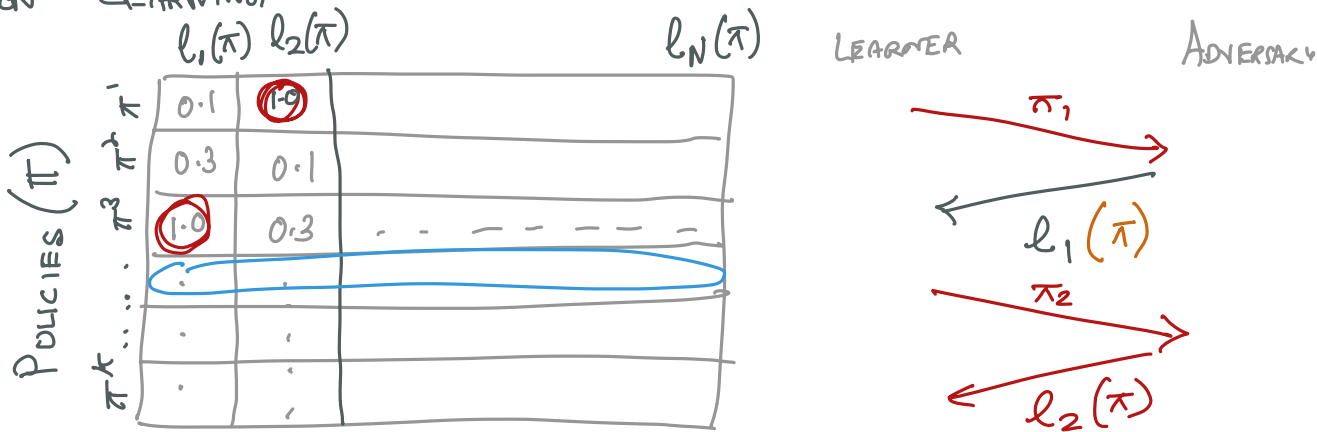


ANALYZING DROCFR

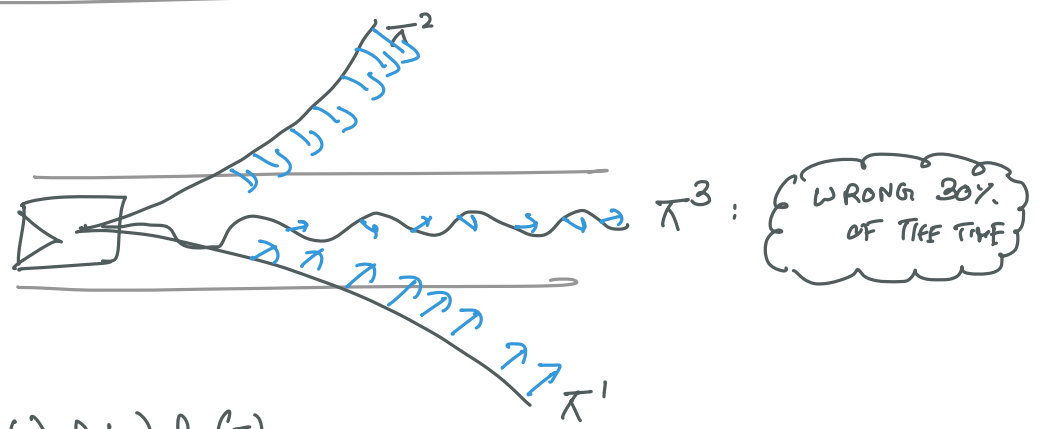
① IMITATION LEARNING CAN BE VIEWED AS INTERACTIVE ONLINE LEARNING



GOAL: MINIMIZE AVG REGRET

$$\text{AVG REGRET}(\pi_{1:N}) = \frac{1}{N} \left(\underbrace{\sum_{i=1}^N l_i(\pi)}_{\text{Loss of learner}} - \min_{\pi \in \Pi} \sum_{i=1}^N l_i(\pi) \right)$$

NO REGRET ALGORITHM := $\lim_{N \rightarrow \infty} \text{AVG REGRET}(\pi_{1:N}) \rightarrow 0$



| | | | | |
|---------|------------|------------|------------|-----|
| | $l_1(\pi)$ | $l_2(\pi)$ | $l_3(\pi)$ | |
| π^1 | 1.0 | 0.1 | 0.7 | 0.7 |
| π^2 | 0.1 | 1.0 | 0.7 | 0.7 |
| π^3 | 0.3 | 0.3 | 0.3 | 0.3 |

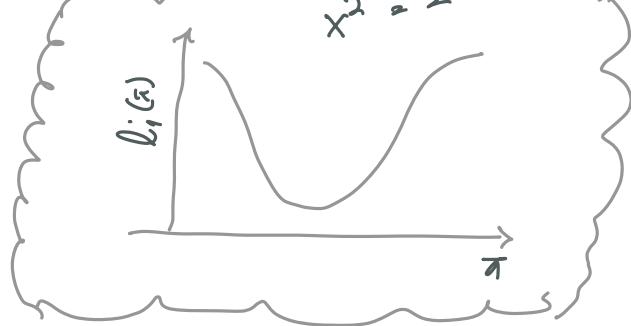
② FOLLOW THE LEADER

$$\pi^i = \underset{\pi \in \Pi}{\text{arg min}} \sum_{j=0}^{i-1} l_j(\pi)$$

$x = 0$

If $l_j(\pi)$ is strongly convex ∞

then FTL is no-regret



③ DAGGER IS FTL.

$$\begin{aligned} \pi_{i+1} &= \text{TRACON}(\mathcal{D}) \\ &= \underset{\pi \in \Pi}{\text{argmin}} \sum_{j=0}^i l_j(\pi) \end{aligned}$$

④ DAGGER RETURNS AT LEAST ONE POLICY π_i

THAT DOES WELL ON ITS OWN INDUCED DISTRIBUTION

$$\pi_i \text{ s.t. } \left[\sum_{t=0}^{T-1} \mathbb{E}_{S_t \sim d_t^{\pi_i}} l(S_t, \pi_i) \right] \leq \underline{O(\dots T)}$$

ASSUMPTION: [RICH POLICY CLASS] \mathcal{D}

$$l(\pi) = \sum_{t=0}^{T-1} \mathbb{E}_{S_t \sim d_t} l(S_t, \pi)$$

FOR ANY LOSS FUNCTION $l(\pi)$, $\exists \pi \in \Pi$ THAT IS GOOD

$$\min_{\pi \in \Pi} l(\pi) \leq O(\epsilon T H)$$

\hookrightarrow recoverability of MDP

PROOF:

LOOK AT THE BEST POLICY THAT DAGGER RETURNS.

$$\min_{i=1, \dots, N} l_i(\pi_i)$$

$$\leq \frac{1}{N} \sum_{i=1}^N l_i(\pi_i)$$

$$\leq \frac{1}{N} \left(\sum_{i=1}^N l_i(\pi) - \min_{\pi \in \Pi} \sum_{i=1}^N l_i(\pi) \right) + \min_{\pi \in \Pi} \frac{1}{N} \sum_{i=1}^N l_i(\pi)$$

~

$$\text{Avg REG} (\pi_{1:n})$$

As $N \rightarrow \infty$, $\text{Avg REG} \rightarrow 0$

$$\frac{\log N}{N}$$

By our Assumption

$$\leq O(\epsilon H T)$$