

Lecture 23: Interactive Imitation Learning

1) Dataset Aggregation with DAgger

Setting: Discounted Infinite Horizon MDP

$$\mathcal{M} = \{ \mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma \} \quad \begin{array}{l} \text{unknown} \\ \text{possibly unobserved} \end{array}$$

Expert knows optimal policy π^* and we can query the expert at any state during training.

Algorithm: DAgger

Initialize π^0 and dataset $\mathcal{D} = \emptyset$

For $t=0, \dots, T-1$:

1) Generate Dataset with π^t & Query Expert

$$\mathcal{D}^t = \{ s_i, a_i^* \} \text{ where } s_i \sim d_{\pi^t} \text{ and } a_i^* = \pi^*(s_i)$$

2) Data Aggregation: $\mathcal{D} = \mathcal{D} \cup \mathcal{D}^t$

3) Update policy via SL:

$$\pi^{t+1} = \underset{\pi \in \Pi}{\operatorname{argmin}} \sum_{s, a \in \mathcal{D}} \ell(\pi, s, a)$$

Last lecture we used an assumption that SL succeeds to reason about the performance of imitation learning. Because DAgger aggregates data, we need to consider a slightly different learning framework.

2) Online learning

The online learning setting captures the idea of learning from additional data over time. It is iterative with two components:

for $t=0, 1, \dots$

1) Learner chooses θ_t (i.e. from past data)

2) Suffer the risk $R_t(\theta_t) = \mathbb{E}_{z \sim \mathcal{D}_t} [l(\theta_t, z)]$
(expected loss)

We care about average regret

$$\frac{1}{T} R(T) = \frac{1}{T} \left[\sum_{i=0}^{T-1} R_t(\theta_t) - \min_{\theta} \sum_{i=0}^{T-1} R_t(\theta) \right]$$

The baseline for regret is the best learned parameter in hindsight.

ex - supervised learning with \mathcal{D}_t a random sample from \mathcal{D} . This is like ingesting a large dataset one training example at a time (streaming) and we hope that the performance is similar to batch learning.

Why is this different from the SL setting?
 \mathcal{D}_t (and thus R_t) can vary in other ways

Example: in DAGGER, we choose π^t and then suffer $l(\pi^t, s_i, \pi^*(s_i))$ for $s_i \sim \mathcal{D}_t$

In this case \mathcal{D}_t actually depends on θ_t

How should the learner choose θ_t ?

Algorithm: Follow-the-Regularized Leader

For $t=0, 1, \dots, T-1$

$$\theta_t = \min_{\theta} \underbrace{\sum_{k=0}^{t-1} R_k(\theta)}_{\text{data aggregation!}} + \lambda f(\theta) \quad \leftarrow \text{regularizer}$$
$$= \sum_{k=0}^{t-1} \mathbb{E}_{z_k \sim \mathcal{D}_k} [l(\theta, z_k)] = \mathbb{E}_{z_k \sim \mathcal{D}_k} \left[\sum_{k=0}^{t-1} l(\theta, z_k) \right]$$

Theorem (FTL): If losses are convex and regularizer is strongly convex, then even if risks R_t (ie. distributions \mathcal{D}_t) are chosen adversarially,

$$\max_{R_0, \dots, R_{T-1}} \frac{1}{T} \left[\sum_{t=0}^{T-1} R_t(\theta_t) - \min_{\theta} \sum_{t=0}^{T-1} R_t(\theta) \right] = O\left(\frac{1}{\sqrt{T}}\right)$$

3) Analysis of DAGGER

We can view DAGGER as an instance of FTL

Corollary: if $\ell(\pi^*, s, \pi^*(s)) = 0$, then

$$\min_{0 \leq t \leq T-1} \mathbb{E}_{s \sim d_{\pi^*}^t} [\ell(\pi_t^t, s, \pi^*(s))] \leq O(1/\sqrt{T}) = \epsilon_{\text{FTL}}$$

Proof: π_t plays the roll of θ_t , & $(s, \pi^*(s))$ is z & \mathcal{D}_t is $d_{\pi^*}^t$

$$\min_{0 \leq t \leq T-1} R_t(\pi_t) \leq \frac{1}{T} \sum_{t=0}^{T-1} R_t(\pi_t) \quad (\text{min} \leq \text{avg})$$

(π^* has 0 loss)

$$= \frac{1}{T} \sum_{t=0}^{T-1} R_t(\pi_t) - R_t(\pi^*)$$

(π^* is minimizer)

$$\leq \frac{1}{T} \left(\sum_{t=0}^{T-1} R_t(\pi_t) - \min_{\pi} \sum_{t=0}^{T-1} R_t(\pi^*) \right)$$

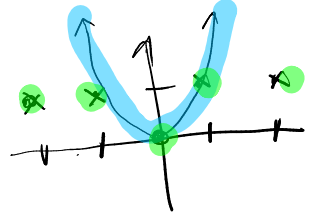
(less than worst-case distributions)
(FTL theorem)

$$\leq \max_{\mathcal{D}_0, \mathcal{D}_T} \frac{1}{T} \left(\sum_{t=0}^{T-1} R_t(\pi_t) - \min_{\pi} \sum_{t=0}^{T-1} R_t(\pi^*) \right) \leq O(1/\sqrt{T}) = \epsilon_{\text{FTL}} \quad \square$$

Notice that this guarantee concerns the performance of π_t on $d_{\pi^*}^t$, i.e. on the state distribution that it induces!

(contrast with supervised ML where we only had guarantees with respect to $d_{\pi^*}^t$!)

Theorem: if $\ell(\pi^t, s, \pi^*(s)) \geq \mathbb{1}\{\pi^t(s) \neq \pi^*(s)\}$,



There exists $0 \leq t \leq T-1$ such that

$$\mathbb{E}_{s \sim \mu_0} \left[V^{\pi^*}(s) - V^{\pi^t}(s) \right] \leq \frac{\max_{s,a} |A^{\pi^*}(s,a)|}{1-\gamma} \cdot \epsilon_{FTL}$$

Proof: We apply PDL in the other direction

$$\mathbb{E}_{s \sim \mu_0} \left[V^{\pi^t}(s) - V^{\pi^*}(s) \right] = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\gamma}^{\pi^t}} \left[A^{\pi^*}(s, \pi^t(a)) \right] \quad (\text{PDL})$$

$$(A^{\pi^*}(s, \pi^t(a)) - A^{\pi^*}(s, \pi^*(a))) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\gamma}^{\pi^t}} \left[A^{\pi^*}(s, \pi^t(a)) - A^{\pi^*}(s, \pi^*(a)) \right]$$

$$\geq \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\gamma}^{\pi^t}} \left[\max_{s,a} |A^{\pi^*}(s,a)| \mathbb{1}\{\pi^t(s) \neq \pi^*(s)\} \right]$$

$$\mathbb{E}_{s \sim \mu_0} \left[V^{\pi^*}(s) - V^{\pi^t}(s) \right] \leq \frac{1}{1-\gamma} \max_{s,a} |A^{\pi^*}(s,a)| \mathbb{E}_{s \sim d_{\gamma}^{\pi^t}} \mathbb{1}\{\pi^t(s) \neq \pi^*(s)\}$$

Assumption on loss

$$\leq \frac{1}{1-\gamma} \max_{s,a} |A^{\pi^*}(s,a)| \underbrace{R_t(\pi^t)}_{\epsilon_{FTL}}$$

How to interpret $\max_{s,a} |A^{\pi^*}(s,a)|$?

Small if expert π^* can quickly recover from mistake.
i.e. if we take action a at state s instead of $\pi^*(s)$, it doesn't impact future rewards too much as long as we follow $\pi^*(s)$ going forward.

($Q^{\pi^*}(s,a)$ is not too much smaller than $V^{\pi^*}(s)$)