Interactive Imitation Learning (continue)

Interactive Imitation Learning Setting

Key assumption: we can query expert π^{\star} at any time and any state during training

Recap

DAgger Revisit



Data Aggregation = Follow-the-Regularized-Leader Online Learner

At iteration t, given π^t

New Data

Recap on the Follow-the-Regularized Leader Guarantee:

FTL: $\theta_{t+1} = \min_{\theta \in \Theta} \theta_{\theta}$

Theorem (FTL) (optional): if Θ is convex, and ℓ_t is convex for all t, and $R(\theta)$ is strongly convex, then for regret of FTL, we have: $\frac{1}{T} \left[\sum_{t=0}^{T-1} \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=0}^{T-1} \ell_t(\theta) \right] = O\left(1 - O\left(1 -$

At the end of iteration t, learner has seen $\ell_0, \ldots, \ell_{t-1}, \ell_t$, learner updates to a new decision:

$$\lim_{\theta \to \Theta} \sum_{i=0}^{t} \ell_i(\theta) + \lambda R(\theta)$$

$$\int_{\Theta} \sum_{t=0}^{T-1} \mathscr{C}_{t}(\theta) = O\left(\frac{1}{\sqrt{T}}\right)$$



1. Finish DAgger's Analysis

2. Intro to Maximum Entropy Inverse RL (We have offline demonstrations, but learner can interact with the environments)

Today's Plan

infinite horizon MDP (assume discrete action space — in fact let's assume 2 actions, so policy is a binary classifier)

$$\mathscr{M} = \left\{ S, A, \gamma, r, P, \mu \right\}$$

Classification:

$$\widehat{\pi} = \arg\min_{\pi} \sum_{x,y} \left[a \right]_{x,y}$$





Total loss so far: $\ell_0(\pi^0) + \ell_1(\pi^1) + \ell_2(\pi^2) + ...$

 π^0

 π^2

. . .

$$\frac{d_{\mu}^{\pi^{0}}\left[\ell\left(\pi,s,\pi^{\star}(s)\right)\right]}{\pi^{1}}$$

Online Learner (FTRL) (i.e., DAgger)

After in total T many iterations, we have the following regret for DAgger:

Avg-Regret_T =
$$\frac{1}{T} \left[\sum_{t=0}^{T-1} \ell_t(\pi^t) - \min_{\pi \in \Pi} \sum_{t=0}^{T-1} \ell_t(\pi) \right] \le O\left(\frac{1}{\sqrt{T}}\right)$$

Recall we assume $\pi^* \in \Pi$, we must have:
 $T-1$ $T-1$



$$\min_{t \in \{0...T-1\}} \ell_t(\pi^t) \le \frac{1}{T} \sum_{t=0}^{T-1} \ell_t(\pi^t) \le \epsilon_{reg}$$

$$0 \leq \sum_{t=0} \ell_t(\pi^*) = 0$$

Which implies that:

 $\ell_t(\pi^t) \leq \epsilon_{reg}$

$$\mathscr{C}_{t}\left(\pi^{t}\right) = \mathbb{E}_{s \sim d_{\mu}^{\pi^{t}}}\left[\mathscr{C}\left(\pi, s, \pi^{\star}(s)\right)\right] \leq \epsilon_{reg}$$

Recall BC, we had: $\mathbb{E}_{s \sim d^{\pi^{\star}}}\left[\ell(\hat{\pi}, s, \pi^{\star}(s))\right] \leq \epsilon, \text{ i.e., we matched to } \pi^{\star} \text{ under } \pi^{\star}\text{'s distribution}$

Summary so far: we know that there must exists $t \in \{0, ..., T-1\}$, such that:

Recall the definition of $\ell_t(\pi^t)$

 π^t matches to π^* under its own state distribution!

Finally, turn things into the performance bound using PDL:

Theorem: There exists a iteration t, such that:

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

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

$$V^{\pi^{\star}} - V^{\pi^{t}} \leq \frac{\max_{s, a} \left| A^{\pi^{\star}}(s, a) \right|}{1 - \gamma} \cdot \epsilon_{reg}$$

 $V^{\pi^{\star}} - V^{\pi^{t}} \le \frac{\max_{s,a} \left| A^{\pi^{\star}}(s,a) \right|}{1 - \nu} \cdot \epsilon_{reg}$

This bound indicates that:

We avoid quadratic error if expert π^{\star} can quickly recover from a mistake

 $\max_{s,a} |A^{\pi^*}(s,a)| \le c \in \mathbb{R}^+$ S, a

i.e., at any state, π^{\star} can quickly recover from your mistake (take action a)



Summary of DAgger

$\mathbb{E}_{s \sim d^{\widehat{\pi}}_{\mu}} \left[\mathbf{1} \{ \widehat{\pi}(s) \neq \pi \right]$

$$V^{\pi^{\star}} - V^{\pi^{t}} \leq$$

DAgger finds a policy $\hat{\pi}$ such that it matches to π^* under $d_{\mu}^{\hat{\pi}}$

$$\star(s)\}] \le \epsilon_{reg} = O(1/\sqrt{T})$$

If expert can quickly recover from a deviation, i.e., $|Q^{\pi^*}(s, a) - V^{\pi^*}(s)|$ is small for all s,

 $\leq O\left(\frac{1}{1-\gamma}\cdot\epsilon_{reg}\right)$



Today's Plan

2. Intro to Maximum Entropy Inverse RL (We have offline demonstrations, but learner can interact with the environments)

Review of the IL settings that we covered so far

1. Offline IL Setting:

We have a dataset

No expert interaction, no real world interaction

2. Interactive IL setting:

$$\mathbf{t} \mathcal{D} = (s_i^{\star}, a_i^{\star})_{i=1}^M \sim d^{\pi^{\star}}$$

- We have access to π^{\star} during training
- Interaction w/ expert and interaction w/ the world (i.e., we can try out our policies)

A new setting (more realistic maybe??)

1. We have an offline dataset $\mathscr{D} = (s_i^{\star}, a_i^{\star})_{i=1}^M \sim d^{\pi^{\star}}$ (e.g., a pre-collected demonstrations)

2. And we can interact with the world (e.g., try out our policy and see what happens)

Hybrid:

Running Example: Human trajectory forecasting



paths and destinations from noisy vision-input

High-level assumptions:

Experts may have some cost function regarding walking in their mind (1) Experts are (approximately) optimizing the cost function (2)

[Kitani, et al, ECCV 12]

Fig. 1. Given a single pedestrian detection, our proposed approach forecasts plausible

Finite horizon MDP ./

We have a dataset

Setting

$$\mathscr{M} = \{S, A, H, c, P, \mu, \pi^{\star}\}$$

(1) Ground truth cost c(s, a) is unknown; (2) assume expert is the optimal policy π^{\star} of the cost c(3) transition P is known

$$\mathbf{t} \mathcal{D} = (s_i^{\star}, a_i^{\star})_{i=1}^M \sim d^{\pi^{\star}}$$

Key Assumption on cost: $c(s, a) = \langle \theta^{\star}, \phi(s, a) \rangle$, linear w.r.t feature $\phi(s, a)$

Running Example: Define feature map

Key Assumption on cost: $c(s, a) = \langle \theta^{\star}, \phi(s, a) \rangle$, linear wrt feature $\phi(s, a)$





sidewalk

State *s*: pixel or a group of neighboring pixels in image)

 $\mathbb{P}(\text{pixels being building})$ $\mathbb{P}(\text{pixels being grass})$ $\phi(s,a) =$ $\mathbb{P}(\text{pixels being sidewalk})$ $\mathbb{P}(\text{pixels being car})$

> Maybe colliding with cars or buildings has **high** cost, but walking on sideway or grass has low cost





Running Example: Human Trajectory Forecasting



State space: grid, action space: 4 actions



We predict that we are more likely to use sidewalk

We will talk about the algorithm (MaxEnt-IRL) behind it next week