Imitation Learning as Game-Solving

Gokul Swamy

Outline for Today 1. Why do we need interaction in imitation learning?

2. What else do we need to tell which mistakes matter?

3. How do we learn a policy that recovers from mistakes that matter if we don't know what the reward function is?

Outline for Today 1. Why do we need interaction in imitation learning? **A**: to be able to tell that we've made a mistake that compounds.

2. What else do we need to tell which mistakes matter?

3. How do we learn a policy that recovers from mistakes that matter if we don't know what the reward function is?

The Pitfalls of Behavioral Cloning



-) Gaussian Policies => MsE $\arg\min_{\pi\in\Pi} \mathbb{E}_{\xi\sim\pi_E} \left| \log\left(\prod_h \pi(a_h | s_h)\right) \right| = \arg\min_{\pi\in\Pi} \sum_h \mathbb{E}_{s_h, a_h\sim\pi_{s}} \left[\log \pi(a_h | s_h)\right]$

What went wrong?

Train Time: $\ell_{BC}(\pi)$

Test Time:

$$= \mathbb{E}_{s_h, a_h \sim \pi_E} \left[-\log \pi(a_h | s_h) \right]$$

$$P_{torr}(x) \neq P_{truin}(x)$$

i coudrate shift"

Covariate Shift \Rightarrow Compounding Errors

No offline IL algorithm can tell the difference between π_1 and π_2

$$a) = \mathbf{1}[s = s_1] s_0) = \pi_2(s_0) = [1,0] s_1) = \pi_2(s_1) = [\varepsilon, 1 - \varepsilon] (z) = [1,0] \pi_2(s_2) = [0,1]$$

 $\mathfrak{g}((\pi_1) = \mathfrak{l}_{\mathfrak{g}}(\pi_2)$

$$(\pi_{\mathsf{G}}, \mathsf{r}) - \mathcal{J}(\pi_{\mathsf{I}}, \mathsf{r}) = \mathcal{L} \cdot \mathsf{H}^{2}$$

$$(\pi_{\mathsf{G}}, \mathsf{r}) - \mathcal{J}(\pi_{\mathsf{I}}, \mathsf{r}) = \mathcal{L} \cdot \mathsf{H}^{2}$$

What we talk about when we talk about ε

1. Finite-sample error: limited number of expert demos.

A: Get more data.

2. Optimization error: imperfect search over policy class.

A: Use more compute.

3. Misspecification error: irreducible error from $\pi_F \notin \Pi$.

A: Use an interactive algorithm.

Interaction Generates Samples from the Test Distribution

Outline for Today 1. Why do we need interaction in imitation learning? **A**: to be able to tell that we've made a mistake that compounds.

2. What else do we need to tell which mistakes matter?

3. How do we learn a policy that recovers from mistakes that matter if we don't know what the reward function is?

Outline for Today 1. Why do we need interaction in imitation learning? A: to be able to tell that we've made a mistake that compounds. 2. What else do we need to tell which mistakes matter? A: information about the set of rewards we could be judged on.

3. How do we learn a policy that recovers from mistakes that matter if we don't know what the reward function is?

Not All Mistakes are Made Equal

We need to be able to tell which mistakes cost us performance.

Moments in Imitation Learning

- distance to realist cal - distance to conter of lang - distance from odge of row - speed (speed limity - ausil dittaut to hive on any - désturie tron recreit persons - listance for the goal

Ok... but which $f \in \mathcal{R}$??? I'll even tell you that $r \in \mathcal{R}$.

P Idea: Be good under all *f* ∈ \Re !

Outline for Today

- 1. Why do we need interaction in imitation learning?
- **A**: to be able to tell that we've made a mistake that compounds.
- 2. What else do we need to tell which mistakes matter?
- A: information about the set of rewards we could be judged on.
- 3. How do we learn a policy that recovers from mistakes that matter if we don't know what the reward function is?
- **A**: Find the policy that is the least distinguishable from the expert's under any reward function in the moment set \mathcal{R} .

Inverse RL as Game-Solving

$\max_{\pi \in \Pi} \min_{f \in \mathscr{R}} J(\pi, f) - J(\pi_E, f)$

where $J(\pi, f) \triangleq \mathbb{E}_{\xi \sim \pi} \left| \sum_{k \in \pi} f(s_h, a_h) \right|.$ h

Approx. Equilibria of IRL Game **Lemma**: Assume $\hat{\pi}$ is an ε -approximate equilibria for the IRL game and for simplicity assume $\pi_F \in \Pi$. Then, $J(\pi_F, r) - J(\pi, r) \leq \mathcal{O}(\epsilon H)$

f∈ℛ

Inverse RL as Game-Solving

- 1. Inverse RL lets avoid compounding errors without needing access to extra expert interaction. -> $\forall \mu \gamma \gamma q \mu \gamma \gamma$
- 2. Inverse RL reduces the search space of policies to just those that are on the Pareto frontier. _> statistic back*
- 3. Inverse RL isn't merely picking a reward that makes the expert look optimal it is fundamentally game-theoretic.

Take-aways from Today

- 1. Why do we need interaction in imitation learning?
- **A**: to be able to tell that we've made a mistake that compounds.
- 2. What else do we need to tell which mistakes matter?
- A: information about the set of rewards we could be judged on.
- 3. How do we learn a policy that recovers from mistakes that matter if we don't know what the reward function is?

under any reward function in the moment set \mathcal{R} .

A: Find the policy that is the least distinguishable from the expert's

If .: MaxEnt Inverse RL min – $\mathbb{H}(\pi)$

s.t. $\forall f \in \mathcal{R}, \quad \mathbb{E}_{\xi \sim \pi} \left[\sum_{h}^{H} f(s_{h}, a_{h}) \right] = \mathbb{E}_{\xi \sim \pi_{E}} \left[\sum_{h}^{H} f(s_{h}, a_{h}) \right]$ $\max_{\lambda \in \mathbb{R}^{|\mathcal{R}|}} \min_{\pi} \mathbb{E}_{\xi \sim \pi} \left[\sum_{h}^{H} \log \pi(a_{h} | s_{h}) \right] + \sum_{f \in \mathcal{Q}} \lambda^{f} (J(\pi, f) - J(\pi_{E}, f))$ $f \in \mathcal{R}$

For some fixed λ_t , we can write the best-response over π as:

 $r_t(s_h, a_h) \triangleq \log \pi(a_h | s_h) + \sum \lambda_t^f f(s_h, a_h)$ $f \in \mathcal{R}$

Let us proceed by backwards-in-time induction over h: Base Case (h = H):

Inductive Step $(h \in [0, H - 1])$:

This is a single-step, action-level maximum entropy problem!

$V_t^{\star}(s_H) \triangleq 0$

 $\pi_t^{\star}(\cdot \mid s_h) = \min_{p \in \Delta(\mathscr{A})} \mathbb{E}_p \left[\log p(a) + \sum_{f \in \mathscr{R}} \lambda_t^f f(s_h, a) + \mathbb{E}_{T(s_h, a)}[V_t^{\star}(s_{h+1})] \right]$

Recall that MaxEnt problems of the form :

f∈ℛ

Have solutions of the form :

 $p^{\star}(x) = \frac{\exp(m(x))}{\sum_{x' \in \mathcal{X}} \exp(m(x'))}$

Here, *m* is just:

$m(a) = \sum_{t} \lambda_{t}^{f} f(s_{h}, a) + \mathbb{E}_{T(s_{h}, a)} [V_{t}^{\star}(s_{h+1})]$

 $\pi_t^{\star}(a_h | s_h) = \frac{\exp\left(\sum_{f \in \mathscr{R}} \lambda_t^f\right)}{\sum_{a \in \mathscr{A}} \exp\left(\sum_{f \in \mathscr{A}} \sum_{f \in \mathscr{A}} \lambda_f^f\right)}$

 $V_t^{\star}(s_h) = \mathbb{E}_{a_h \sim \pi_t^{\star}(s_h)}[\log \pi_t^{\star}(a_h \mid s_h)]$

Can solve for π_t^{\star} via "soft" policy / value iteration! Closely connected to Natural Policy Gradient and Hedge!

$$\frac{f}{t}f(s_h, a_h) + \mathbb{E}_{T(s_h, a_h)}[V_t^{\star}(s_{h+1})]\Big)$$
$$= \Re \lambda_t^f f(s_h, a) + \mathbb{E}_{T(s_h, a)}[V_t^{\star}(s_{h+1})]\Big)$$

$$(S_h) + \lambda_t^f f(s_h, a_h) + \mathbb{E}_{T(s_h, a_h)}[V_t^{\star}(s_{h+1})]$$