Interactive Imitation Learning

Recap

- **The Behavior Cloning algorithm:**
- Choose regression (for continuous action) or classification loss $\ell(\pi(s), a)$, and perform SL:



What could go wrong? [Pomerleau89,Daume09] Predictions affect future inputs/

observations

Learned Policy



Distribution Shift: Example



Assume SL returned such policy $\widehat{\pi}$

$$\widehat{a}_{0} = \begin{cases} a_{1} \quad \text{w/ prob } 1 - \epsilon/(1 - \gamma) \\ a_{2} \quad \text{w/ prob } \epsilon/(1 - \gamma) \end{cases}, \quad \widehat{\pi}(s_{1}) = a_{2}, \ \widehat{\pi}(s_{2}) = a_{2} \\ \widehat$$

We will have good supervised learning error:

$$\mathbb{E}_{s \sim d_{s_0}^{\pi^*}} \mathbb{E}_{a \sim \hat{\pi}(\cdot|s)} \mathbf{1} \left(a \neq \pi^*(s) \right) = \epsilon$$

But we have quadratic error in performance:

$$V_{s_0}^{\widehat{\pi}} = \frac{\gamma}{1-\gamma} - \frac{\epsilon\gamma}{(1-\gamma)^2} = V_{s_0}^{\pi^*} - \frac{\epsilon\gamma}{(1-\gamma)^2}$$

Issue: once we make a mistake at s_0 , we end up in s_2 which is not in the training data!



An Autonomous Land Vehicle In A Neural Network [Pomerleau, NIPS '88]



"If the network is not presented with sufficient variability in its training exemplars to cover the conditions it is likely to encounter...[it] will perform poorly"

Question for today:

How to mitigate the distribution shift issue?

Solution:

Interactive Imitation Learning Setting

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

(Recall that previously we only had an offline dataset $\mathscr{D} = (s_i^{\star}, a_i^{\star})_{i=1}^M \sim d_{\mu}^{\pi^{\star}}$)

2. Analysis of DAgger: DAgger as online learning

Outline for today:

1. The DAgger (Data Aggregation) Algorithm

Recall the Main Problem from Behavior Cloning:

No training data of "recovery" behavior

Learned Policy



Intuitive solution: Interaction

Use interaction to collect data where learned policy goes



General Idea: Iterative Interactive Approach



Updated Policy

All DAgger slides credit: Drew Bagnell, Stephane Ross, Arun Venktraman



[Ross11a] DAgger: Dataset Aggregation **Oth iteration**





Supervised Learning

DAgger: Dataset Aggregation [Ross11a] 1st iteration

Execute π_1 and Query Expert





[Ross11a] DAgger: Dataset Aggregation 1st iteration

Execute π_1 and Query Expert





New Data







[Ross11a] DAgger: Dataset Aggregation 1st iteration

Execute π_1 and Query Expert





New Data

[Ross11a] DAgger: Dataset Aggregation 1st iteration

Execute π_1 and Query Expert



DAgger: Dataset Aggregation [Ross11a] 2nd iteration

Execute π_2 and Query Expert



17

[Ross11a] DAgger: Dataset Aggregation nth iteration

Execute π_{n-1} and Query Expert



Success!



[Ross AISTATS 2011]

Average Falls/Lap



[Ross AISTATS 2011]

FPS: 24 Attempt: 1 of 1 AgentLinear Selected Actions:

RIGHT



More fun than Video Games...



[Ross ICRA 2013] 22

Forms of the Interactive Experts

Interactive Expert is expensive, especially when the expert is human...

But expert does not have to be human...

Example: high-speed off-road driving [Pan et al, RSS 18, Best System Paper]



Fig. 4: The AutoRally car and the test track.

Goal: learn a racing control policy that maps from data on cheap on-board sensors (raw-pixel imagine) to low-level control (steer and throttle)



Steering + throttle

(a) raw image



Forms of the Interactive Experts

Example: high-speed off-road driving [Pan et al, RSS 18, Best System Paper]

The MPC is the expert in this case!



Their Setup: At Training, we have expensive sensors for accurate state estimation and we have computation resources for MPC (i.e., high-frequency replanning)

Distribution Shift: Example



Assume SL returned such policy $\widehat{\pi}$

$$\widehat{\pi}(s_0) = \begin{cases} a_1 & \text{w/ prob } 1 - \epsilon/(1 - \gamma) \\ a_2 & \text{w/ prob } \epsilon/(1 - \gamma) \end{cases}, \quad \widehat{\pi}(s_1) = a_2, \, \widehat{\pi}(s_2) = a_2 \end{cases}$$

We will have good supervised learning error:

$$\mathbb{E}_{s \sim d_{s_0}^{\pi^*}} \mathbb{E}_{a \sim \widehat{\pi}(\cdot|s)} \mathbf{1} \left(a \neq \pi^*(s) \right) = \epsilon$$

But we have quadratic error in performance:

$$V_{s_0}^{\hat{\pi}} = V_{s_0}^{\pi^*} - \frac{\epsilon\gamma}{(1-\gamma)^2}$$



Distribution Shift: Example



Assume SL returned such policy $\widehat{\pi}$ $\widehat{\pi}(s_0) = \begin{cases} a_1 & \text{w/ prob } 1 - \epsilon/(1 - \gamma) \\ a_2 & \text{w/ prob } \epsilon/(1 - \gamma) \end{cases}, \quad \widehat{\pi}(s_1) = a_2, \ \widehat{\pi}(s_2) = a_2 \end{cases}$

Why DAgger can fix this problem?





Outline for today:

2. Analysis of DAgger: DAgger as online learning

Learner picks a decision θ_0





convex Decision set Θ



[Vovk92,Warmuth94,Freund97,Zinkevich03,Kalai05,Hazan06,Kakade08]

Online Learning

Adversary picks a loss $\mathscr{C}_0: \Theta \to \mathbb{R}$

Learner picks a new decision θ_1

Adversary picks a loss $\ell_1: \Theta \to \mathbb{R}$

Adversary



$$\sum_{t=0}^{-1} \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=0}^{T-1} \ell_t(\theta)$$

. . .



Example: online linear regression

- 1. Learner first picks $\theta_t \in \text{Ball} \subset \mathbb{R}^d$
- 2. Adversary **then** picks $x_t \in \mathcal{X} \subset \mathbb{R}^d, y_t \in [a, b]$
 - 3. Learner suffers loss $\ell_t(\theta_t) = (\theta_t^T x_t y_t)^2$
- Learner has to make decision θ_t based on history up to t 1, while adversary could pick (x_t, y_t) even after seeing θ_t
 - Adversary seems too powerful...

- Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?
 - **Every iteration** *t* :

Example: online linear regression

BUT, a very intuitive algorithm actually achieves no-regret property:

1. Learner first picks θ_t that minimizes the aggregated loss

$$\theta_t = \arg\min_{\theta \in \mathsf{Ball}} \sum_{i=0}^{t-1} \left(\theta^\top x_i - y_i\right)^2 + \lambda \|\theta\|_2^2$$

This is called Follow-the-Regularized-Leader (FTRL), and it achieves no-regret property:

$$\sum_{i=0}^{T-1} \ell_i(\theta_i) - \min_{\theta \in \mathsf{Ball}} \sum_{i=0}^{T-1} \ell_i(\theta) = O\left(\sqrt{T}\right)$$

Every iteration *t* :

Generally, Follow-the-Regularized-Leader is no-regret

At time step t, learner has seen $\ell_0, \ldots \ell_{t-1}$, which new decision she could pick?

FTRL: $\theta_t = \min_{\theta \in \Theta} \theta_{\theta}$

Informal Theorem (FTRL): when things are convex, FTRL is no-regret, i.e., $\frac{1}{T} \left[\sum_{t=0}^{T-1} \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=0}^{T-1} \ell_t(\theta) \right] = O\left(1/\sqrt{T}\right)$

$$\inf_{\mathbf{\Theta}} \sum_{i=0}^{t-1} \ell_i(\theta) + \lambda R(\theta)$$

DAgger Revisit



Summary for Today

1. The DAgger algorithm

Initialize π^0 , and dataset 2

For
$$t = 0 \rightarrow T - 1$$
:

$$\mathcal{D} = \mathcal{O}$$

- 1. W/ π^t , generate dataset $\mathscr{D}^t = \{s_i, a_i^{\star}\}, s_i \sim d_{\mu}^{\pi^t}, a_i^{\star} = \pi^{\star}(s_i)$ 2. Data aggregation: $\mathscr{D} = \mathscr{D} + \mathscr{D}^t$ 3. Update policy via Supervised-Learning: $\pi^{t+1} = SL(\mathscr{D})$

2. We can see that DAgger is essentially an online-learning algorithm (FTRL)