# **Interactive Imitation Learning**

#### Recap

#### The Behavior Cloning algorithm:

Choose regression (for continuous action) or classification loss  $\ell(\pi(s), a)$ , and perform SL:

$$\widehat{\pi} = \min_{\pi \in \Pi} \sum_{i=1}^{M} \ell(\pi(s_i^{\star}), a_i^{\star}) \qquad (s^{\star}, a^{\star})$$
ERM
$$s^{\star} - d^{\star}$$

$$s^{\star} - d^{\star}$$

$$c^{\star} - \tau^{\star}(\cdot)(s^{\star})$$







{s. a.) {s. a.)









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







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:

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

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"

### 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** 

#### **Solution:**

#### **Interactive Imitation Learning Setting**

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

#### **Solution:**

#### **Interactive Imitation Learning Setting**

#### Key assumption: we can query expert $\pi^*$ at any time and any state during training

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

#### **Outline for today:**

1. The DAgger (Data Aggregation) Algorithm

2. Analysis of DAgger: DAgger as online learning

### Recall the Main Problem from Behavior Cloning:



### Intuitive solution: Interaction

Use interaction to collect data where learned policy goes

### General Idea: Iterative Interactive Approach



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

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







### DAgger: Dataset Aggregation [Ross11a] 1st iteration

Execute  $\pi_1$  and Query Expert

**New Data** 



States from the learned policy

### DAgger: Dataset Aggregation [Ross11a] 1st iteration

Execute  $\pi_1$  and Query Expert



**New Data** 



### DAgger: Dataset Aggregation [Ross11a] 1st iteration

Execute  $\pi_1$  and Query Expert

New Data



### DAgger: Dataset Aggregation [Ross11a] 2nd iteration

Execute  $\pi_2$  and Query Expert

**New Data** 

17



### DAgger: Dataset Aggregation [Ross11a] n<sup>th</sup> iteration

Execute  $\pi_{n-1}$  and Query Expert



# Success!



# Success!



# Success!



[Ross AISTATS 2011]

### Average Falls/Lap



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

RIGHT



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

RIGHT


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

RIGHT



## More fun than Video Games...



## More fun than Video Games...



## More fun than Video Games...



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

But expert does not have to be human...

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.

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]



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)

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

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)



(a) raw image

Steering + throttle

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

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

Their Setup:

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

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

Their Setup:

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

The MPC is the expert in this case!

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

Their Setup:

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

The MPC is the expert in this case!



### **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:

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

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

### **Distribution Shift: Example**



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

Why DAgger can fix this problem?

**Outline for today:** 

1. The DAgger (Data Aggregation) Algorithm

2. Analysis of DAgger: DAgger as online learning

. . .

## Learner



**Adversary** 





. . .

convex Decision set  $\Theta$ 





convex Decision set  $\Theta$ 

Learner picks a decision  $\theta_0$ 

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

. . .

 $l_0 \in \Theta_0$ =  $(\varphi_0^T \chi_0 - \gamma_0)^T$ 

**Adversary** 



## Learner



convex Decision set  $\boldsymbol{\Theta}$ 

Learner picks a decision  $\theta_0$ 

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

Learner picks a new decision  $\theta_1$ 

Adversary picks a loss 
$${\mathscr C}_1:\Theta o {\mathbb R}$$

Adversary





Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?

Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?

**Every iteration** *t* :

Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?

**Every iteration** *t* :

1. Learner first picks  $\theta_t \in \text{Ball} \subset \mathbb{R}^d$ 

Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?

**Every iteration** *t* :

1. Learner first picks  $\theta_t \in \text{Ball} \subset \mathbb{R}^d$ 2. Adversary **then** picks  $x_t \in \mathcal{X} \subset \mathbb{R}^d, y_t \in [a, b]$ 

Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?

**Every iteration** *t* :

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  $\ell_t(\theta_t) = (\theta_t^{\mathsf{T}} x_t - y_t)^2$ 

Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?

**Every iteration** *t* :

1. Learner first picks  $\theta_t \in \text{Ball Q} \mathbb{R}^d$ 2. Adversary **then** picks  $x_t \in \mathcal{X} \subset \mathbb{R}^d, y_t \in [a, b]$ 

3. Learner suffers loss  $\ell_t(\theta_t) = (\theta_t^{\mathsf{T}} x_t - y_t)^2$ 

Learner has to make decision  $\theta_t$  based on history up to t - 1, while adversary could pick  $(x_t, y_t)$  even after seeing  $\theta_t$ 

Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?

**Every iteration** *t* :

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 
$$\mathscr{C}_t(\theta_t) = (\theta_t^{\top} x_t - y_t)^2$$

Learner has to make decision  $\theta_t$  based on history up to t - 1, while adversary could pick  $(x_t, y_t)$  even after seeing  $\theta_t$ 

Adversary seems too powerful...

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

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



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

#### **Every iteration** *t* :

1. Learner first picks  $\theta_t$  that minimizes the aggregated loss

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

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

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

#### **Every iteration** *t* :

1. Learner first picks  $\theta_t$  that minimizes the aggregated loss

$$\theta_t = \arg\min_{\theta \in \mathsf{Ba}} \left[ \sum_{i=0}^{t-1} \theta^{\mathsf{T}} x_i - y_i \right]^2 + \lambda \|\theta\|_2^2$$

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

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

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



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

At time step *t*, learner has seen  $\ell_0, \dots \ell_{t-1}$ , which new decision she could pick? **FTRL:**  $\theta_t = \min_{\theta \in \Theta} \sum_{i=0}^{t-1} \ell_i(\theta) + \lambda R(\theta)$ 

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
















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

#### **Summary for Today**

#### 1. The DAgger algorithm

Initialize  $\pi^0$ , and dataset  $\mathcal{D} = \emptyset$ For  $t = 0 \rightarrow T - 1$ : 1. W/  $\pi^t$ , generate dataset  $\mathcal{D}^t = \{s_i, a_i^{\star}\}, s_i \sim d_{\mu}^{\pi^t}, a_i^{\star} = \pi^{\star}(s_i)$ 2. Data aggregation:  $\mathcal{D} = \mathcal{D} + \mathcal{D}^t$ 3. Update policy via Supervised-Learning:  $\pi^{t+1} = SL(\mathcal{D})$ 

#### **Summary for Today**

1. The DAgger algorithm

Initialize  $\pi^0$ , and dataset  $\mathscr{D} = \mathscr{O}$ For  $t = 0 \to T - 1$ : 1. W/  $\pi^t$ , 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)