# Introduction to Imitation Learning & the Behavior Cloning Algorithm

#### Annoucements

2. Releasing the next reading quiz on DPO

3. No class this Wednesday and no office hour this Thursday — traveling to DC for DoD meetings

1. We had a typo in 2.2 of the homework, fixed and updated pdf/latex are posted on ED



#### Recap

#### **Infinite horizon Discounted MDPs**

 $\mathcal{M} = \{S_i\}$ 

Average state distribut

$$\{A, \gamma, r, P, \mu\}$$
 What if *r* is unknown

ution: 
$$d^{\pi} = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^{h} \mathbb{P}_{h}^{\pi}$$

#### Recap

#### **Infinite horizon Discounted MDPs**

#### $\mathcal{M} = \{S\}$

Today: how to learn directly from expert demonstations

$$\{A, \gamma, r, P, \mu\}$$
 What if *r* is unknown

We have covered how to learn a reward from binary preference data...

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

1. Offline Imitation Learning: Behavior Cloning

2. Performance difference lemma and its application to proving BC's bound



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





30x32 Video Input Retina

Figure 1: ALVINN Architecture

# Imitation Learning



# Imitation Learning

#### Expert Demonstrations



- SVM

. . .

- LWR



 Gaussian Process Kernel Estimator • Deep Networks **Random Forests** 

## Maps states to <u>actions</u>

# Learning to Drive by Imitation

## Input:



### Camera Image

## [Pomerleau89, Saxena05, Ross11a] Output:





Steering Angle in [-1, 1]

## Supervised Learning Approach: Behavior Cloning

**Expert Trajectories** 



control (steering direction)

[Widrow64, Pomerleau89]

#### Dataset



## LLMs are trained via BC in their pre-training phase

Take a sentence from the web:

State

Reinforcement learning (RL) is an interdisciplinary area of machine learning and optimal control concerned with how an intelligent agent should take actions in a dynamic environment in order to maximize a reward signal.

#### Action

## LLMs are trained via BC in their pre-training phase

Take a sentence from the web:

State

Reinforcement learning (RL) is an interdisciplinary area of machine learning and optimal control concerned with how an intelligent agent should take actions in a dynamic environment in order to maximize a reward signal.

Action

## LLMs are trained via BC in their pre-training phase

Take a sentence from the web:

State

Reinforcement learning (RL) is an interdisciplinary area of machine learning and optimal control concerned with how an intelligent agent should take actions in a dynamic environment in order to maximize a reward signal.

Forcing LLM to predict the next "action" conditioned on past...

Action

#### Let's formalize the offline IL Setting and the Behavior Cloning algorithm

Discounted infinite horizon

We have a dataset

$$\mathsf{MDP}\,\mathscr{M} = \{S, A, \gamma, r, P, \rho, \pi^{\star}\}$$

Ground truth reward  $r(s, a) \in [0,1]$  is unknown; For simplicity, let's assume expert is a (nearly) optimal policy  $\pi^{\star}$ 

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

Goal: learn a policy from  $\mathscr{D}$  that is as good as the expert  $\pi^{\star}$ 

### Let's formalize the Behavior Cloning algorithm

BC is a Reduction to Supervised Learning:

 $\hat{\pi} = \arg \min$  $\pi \in \mathbb{R}$ 

Many choices of loss functions:

BC Algorithm input: a restricted policy class  $\Pi = \{\pi : S \mapsto \Delta(A)\}$ 

$$\prod_{\Pi} \sum_{i=1}^{M} \ell(\pi, s^{\star}, a^{\star})$$

1. Negative log-likelihood (NLL):  $\ell(\pi, s^{\star}, a^{\star}) = -\ln \pi(a^{\star} | s^{\star})$ 

2. square loss (i.e., regression for continuous action):  $\ell(\pi, s^*, a^*) = \|\pi(s^*) - a^*\|_2^2$ 

#### Analysis

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

Does that imply that  $\hat{\pi}$  is a good policy? What's the performance difference between  $\hat{\pi}$  and  $\pi^*$ ?



Assumption: we are going to assume Supervised Learning succeeded

## $\neq \pi^{\star}(s) \leq \epsilon \in \mathbb{R}^+$





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

1. Offline Imitation Learning: Behavior Cloning

2. Performance difference lemma and its application to proving BC's bound

Given two policies  $\pi : S \mapsto \Delta(A), \pi'$ :



#### **Performance Difference Lemma**

$$S \mapsto \Delta(A), \text{ recall } V^{\pi}(s_0) = \mathbb{E}\left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \,|\, \pi\right]$$

**Performance Difference Lemma (PDL):** 

$$\mathbb{E}_{s \sim d_{s_0}^{\pi}} \left[ \mathbb{E}_{a \sim \pi(\cdot|s)} Q^{\pi'}(s,a) - V^{\pi'}(s) \right]$$
  
$$\mathbb{E}_{s \sim d_{s_0}^{\pi}} \left[ \mathbb{E}_{a \sim \pi(\cdot|s)} A^{\pi'}(s,a) \right]$$

#### **PDL Explanation**

 $V^{\pi}(s_0) - V^{\pi'}(s_0) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{s_0}^{\pi}} \left[ \mathbb{E}_{a \sim \pi(\cdot|s)} A^{\pi'}(s, a) \right]$ 

# $V^{\pi}(s_0) - V^{\pi'}(s_0) = \frac{1}{1 - \gamma}$

Proof of Sketch (see reading material for detailed steps)

$$\begin{split} V^{\pi}(s_{0}) &- V^{\pi'}(s_{0}) \\ &= V^{\pi}(s_{0}) - \mathbb{E}_{a_{0} \sim \pi(\cdot|s_{0})} \left[ r(s_{0}, a_{0}) + \gamma \mathbb{E}_{s' \sim P(s_{0}, a_{0})} V^{\pi'}(s') \right] + \mathbb{E}_{a_{0} \sim \pi(\cdot|s_{0})} \left[ r(s_{0}, a_{0}) + \gamma \mathbb{E}_{s' \sim P(s_{0}, a_{0})} V^{\pi'}(s') \right] - V^{\pi'}(s_{0}) \\ &= \gamma \mathbb{E}_{a_{0} \sim \pi(\cdot|s_{0})} \mathbb{E}_{s_{1} \sim P(s_{0}, a_{0})} \left[ V^{\pi}(s_{1}) - V^{\pi'}(s_{1}) \right] + \mathbb{E}_{a_{0} \sim \pi(\cdot|s_{0})} \left[ r(s_{0}, a_{0}) + \gamma \mathbb{E}_{s' \sim P(s_{0}, a_{0})} V^{\pi'}(s') \right] - V^{\pi'}(s_{0}) \\ &= \gamma \mathbb{E}_{a_{0} \sim \pi(\cdot|s_{0})} \mathbb{E}_{s_{1} \sim P(s_{0}, a_{0})} \left[ V^{\pi}(s_{1}) - V^{\pi'}(s_{1}) \right] + \mathbb{E}_{a_{0} \sim \pi(\cdot|s_{0})} \left[ Q^{\pi'}(s_{0}, a_{0}) - V^{\pi'}(s_{0}) \right] \\ &= \gamma \mathbb{E}_{a_{0} \sim \pi(\cdot|s_{0})} \mathbb{E}_{s_{1} \sim P(s_{0}, a_{0})} \left[ V^{\pi}(s_{1}) - V^{\pi'}(s_{1}) \right] + \mathbb{E}_{a_{0} \sim \pi(\cdot|s_{0})} \left[ A^{\pi'}(s_{0}, a_{0}) \right] \end{split}$$

#### **PDL Proof**

$$-\mathbb{E}_{s \sim d_{s_0}^{\pi}}\left[\mathbb{E}_{a \sim \pi(\cdot|s)}A^{\pi'}(s,a)\right]$$

#### An Application of PDL in Policy Iteration

Recall that  $\pi^{t+1}(x)$ 

Show monotonic improvement using PDL:

$$V^{\pi^{t+1}}(s_0) - V^{\pi^t}(s_0) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{s_0}^{\pi^{t+1}}} \mathbb{E}_{a \sim \pi^{t+1}(\cdot|s)} A^{\pi^t}(s, a)$$





$$s) = \arg\max_{a} A^{\pi^{t}}(s, a)$$

$$A^{\pi^t}(s,\pi^{t+1}(s))$$

$$A^{\pi^t}(s,\pi^t(s)) = 0$$

#### Analysis of BC

Theorem [BC Performance] With probability at least  $1 - \delta$ , BC returns a policy  $\hat{\pi}$ :  $V^{\pi^*} - V^{\hat{\pi}} \leq \frac{2}{(1 - \gamma)^2} \epsilon$ 

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

$$= \mathbb{E}_{s \sim d^{\pi^{\star}}} A^{\widehat{\pi}}(s, \pi^{\star}(s)) - \mathbb{E}_{s \sim d^{\pi^{\star}}} A^{\widehat{\pi}}(s, \widehat{\pi}(s))$$

$$\leq \mathbb{E}_{s \sim d^{\pi^{\star}}} \frac{2}{1 - \gamma} \mathbf{1} \left\{ \widehat{\pi}(s) \neq \pi^{\star}(s) \right\}$$

$$\leq \frac{2}{1-\gamma}\epsilon$$

The quadratic amplification is annoying; Related to pre-trained LLM hallucination; Will see how to fix it next lecture

PDL: how to capture the performance difference between two policies (as important as Simulation lemma)

#### Summary

BC: simple algorithm that directly learns from human demonstrations; used in robotics and NLP