# Introduction to Imitation Learning & the Behavior Cloning Algorithm

# Annoucements

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

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

# Recap

## Infinite horizon Discounted MDPs

$$\mathcal{M} = \{S, A, \gamma, r, P, \mu\}$$

Average state distribution: $d^\pi = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^\pi$

# Recap

## Infinite horizon Discounted MDPs

$$\mathcal{M} = \{S, A, \gamma, r, P, \mu\}$$

What if $r$ is unknown

Average state distribution: $d^\pi = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^\pi$

# Recap

## Infinite horizon Discounted MDPs

$$\mathcal{M} = \{S, A, \gamma, r, P, \mu\}$$

What if $r$ is unknown

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

# Recap

## Infinite horizon Discounted MDPs

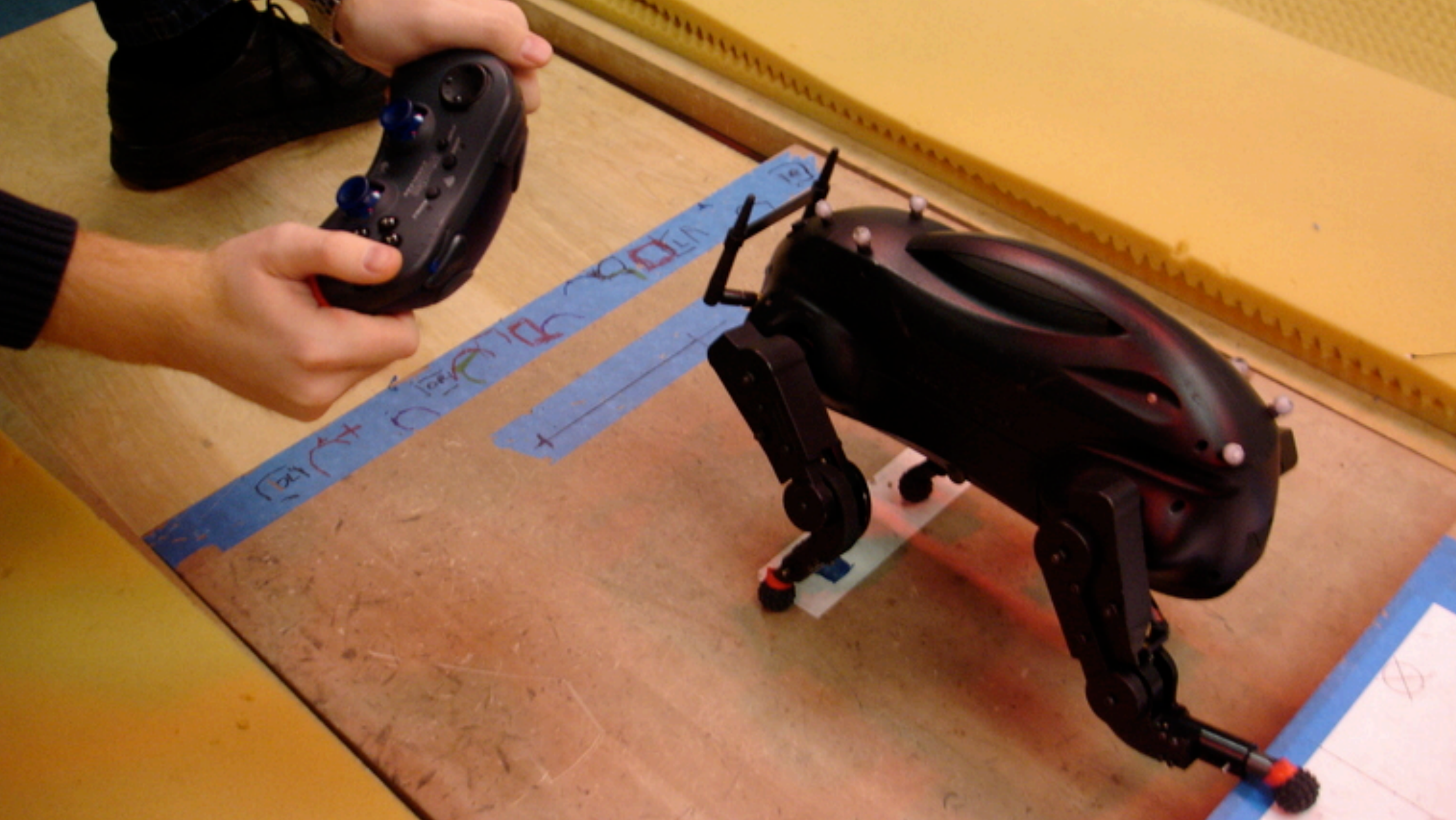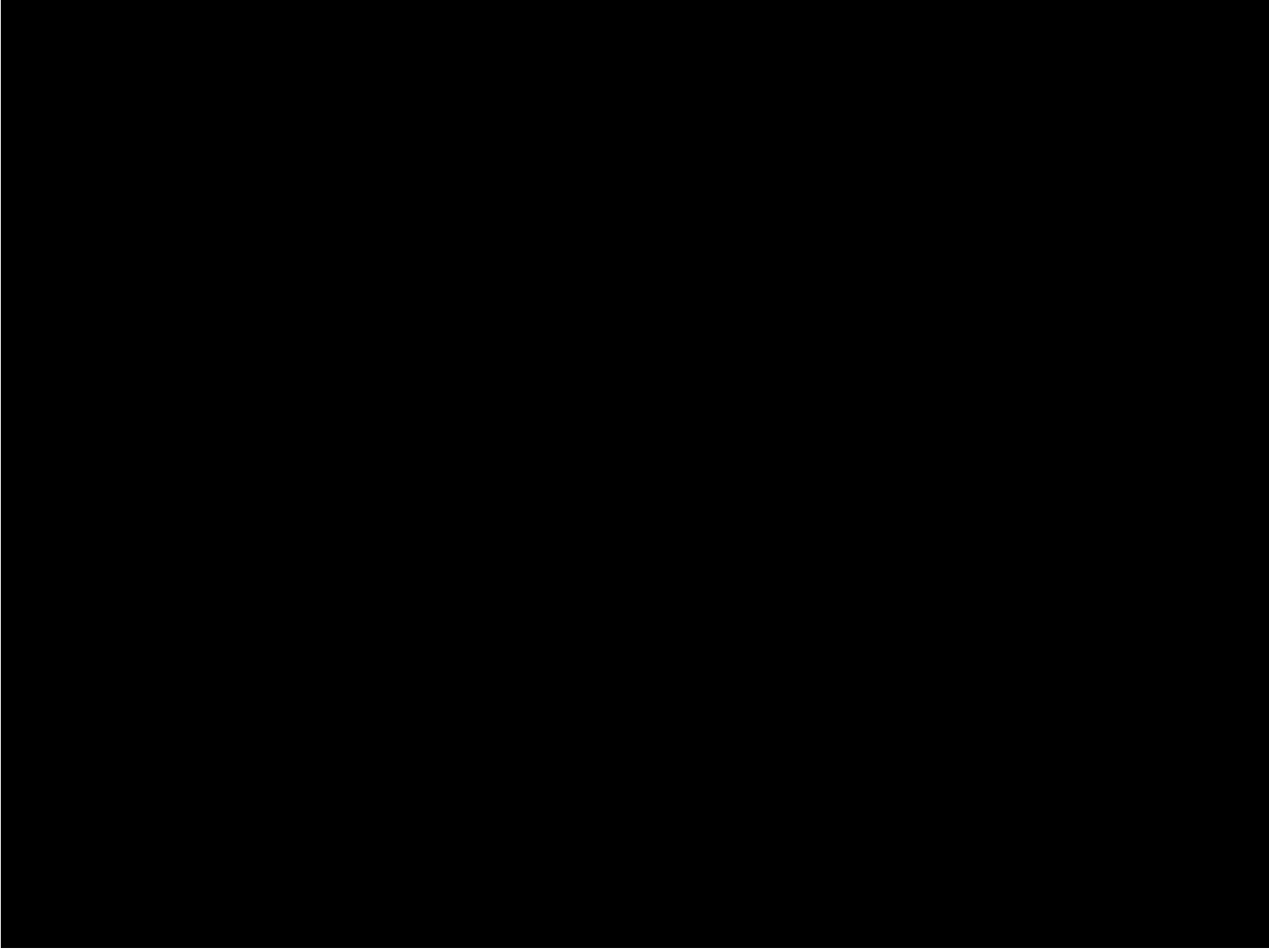$$\mathscr{M} = \{S, A, \gamma, r, P, \mu\}$$

What if $r$ is unknown

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

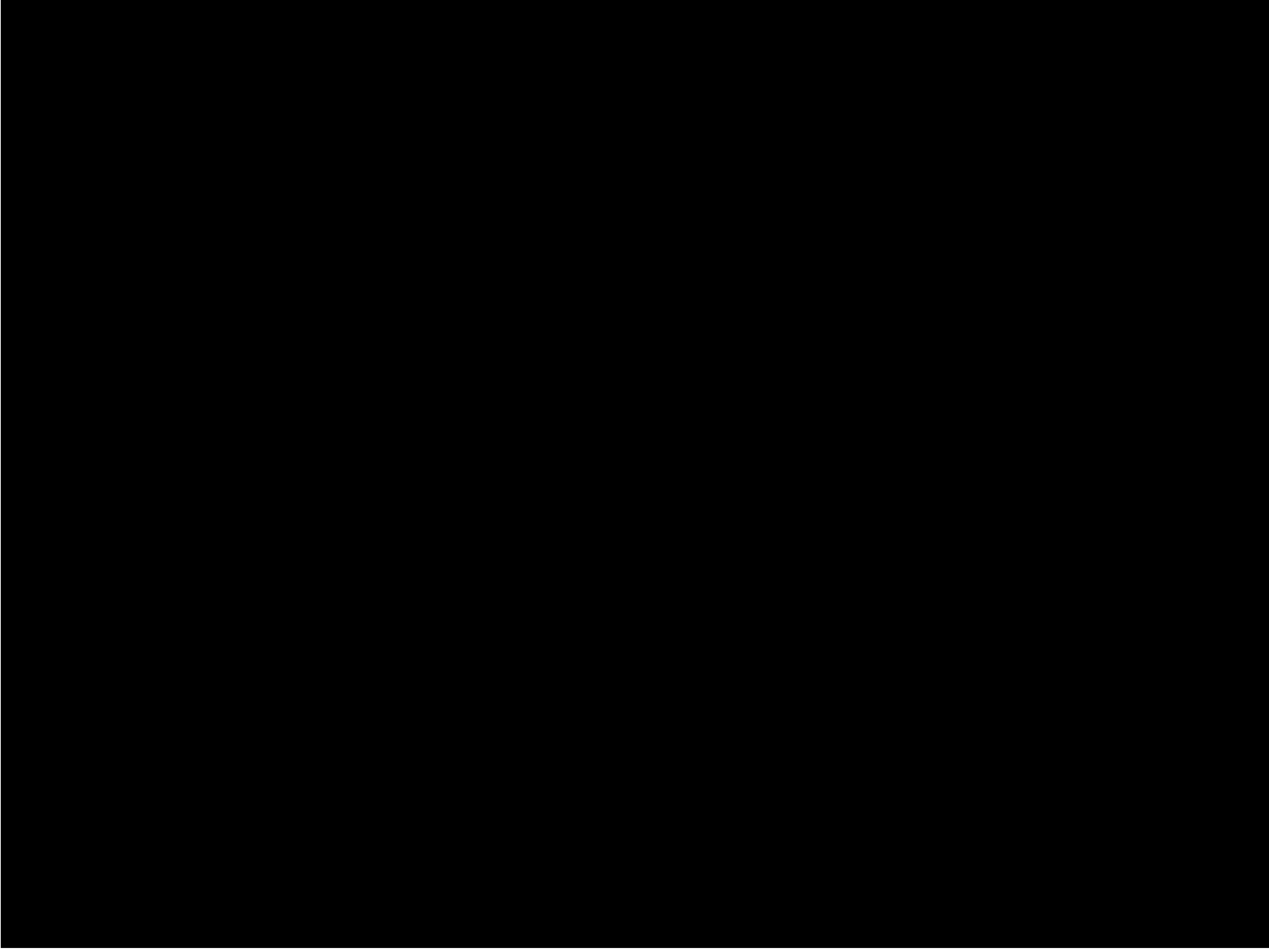Today: how to learn directly from expert demonstations

# 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]



Road Intensity Feedback Unit

45 Direction Output Units

29 Hidden Units

30x32 Video Input Retina
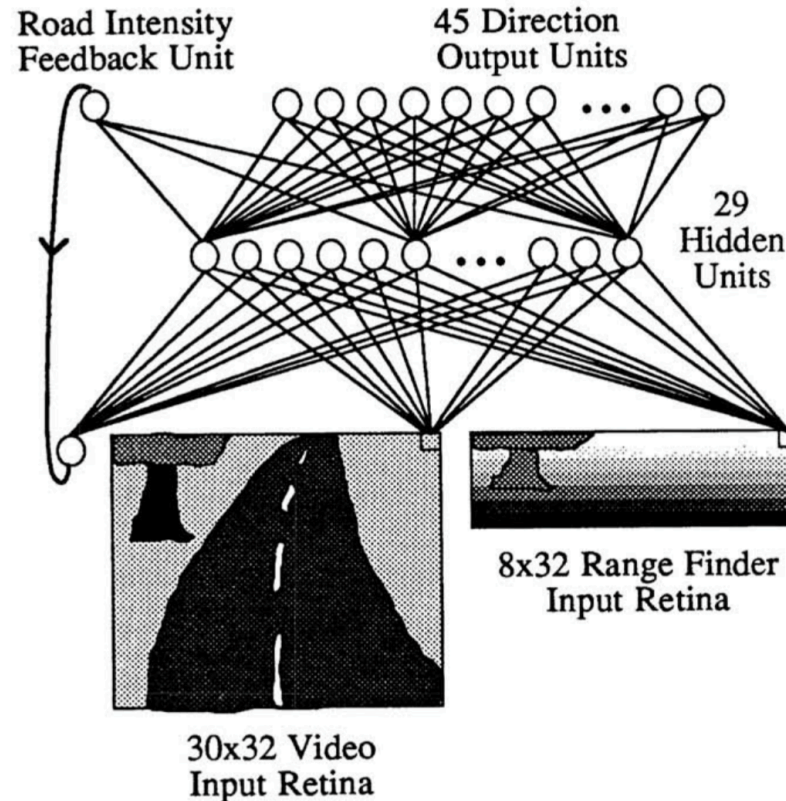
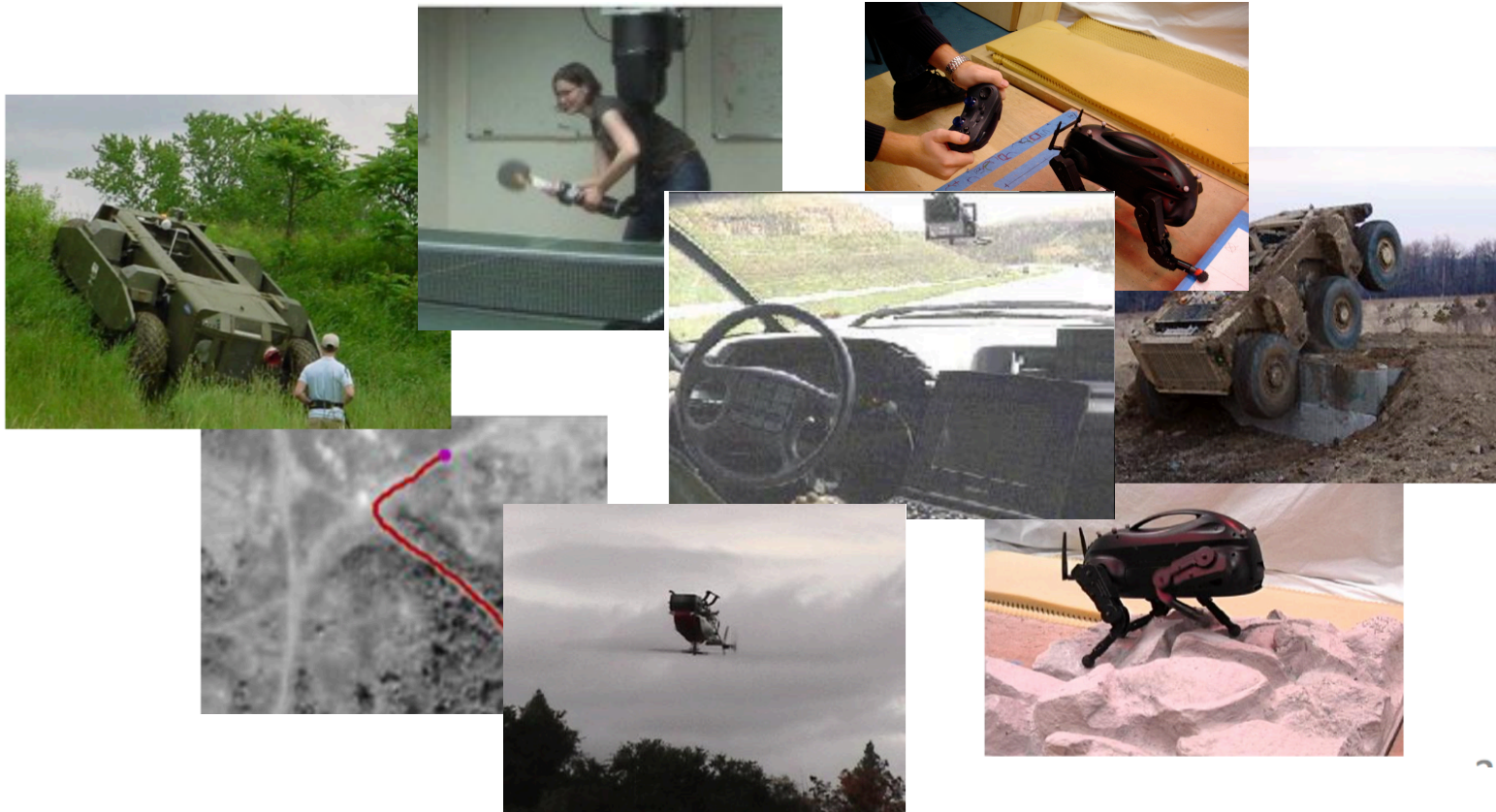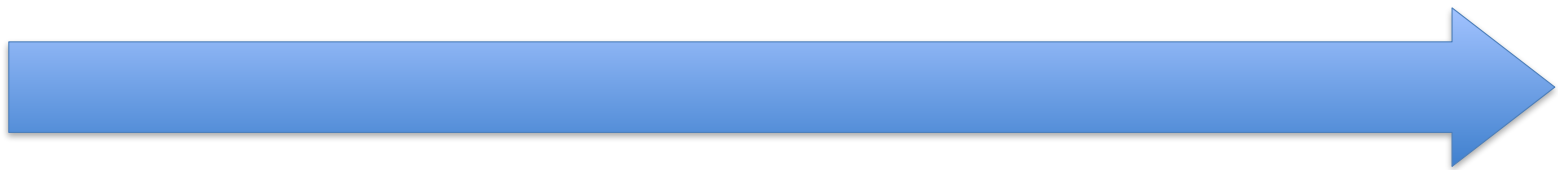8x32 Range Finder Input Retina

Figure 1: ALVINN Architecture

# Imitation Learning

# Imitation Learning

# Imitation Learning

# Imitation Learning



Expert Demonstrations

# Imitation Learning



Expert Demonstrations

Machine Learning Algorithm

- SVM
- Gaussian Process
- Kernel Estimator
- Deep Networks
- Random Forests
- LWR
- …

# Imitation Learning



Expert Demonstrations → Machine Learning Algorithm → Policy $\pi$

- SVM
- Gaussian Process
- Kernel Estimator
- Deep Networks
- Random Forests
- LWR
- …

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

# Learning to Drive by Imitation

[Pomerleau89, Saxena05, Ross11a]

## Input:



Camera Image

## Output:

Policy



Steering Angle
in [-1, 1]

11

# Supervised Learning Approach: Behavior Cloning



Expert Trajectories

Dataset

# Supervised Learning Approach: Behavior Cloning



Expert Trajectories

Dataset

$X$ ⋮ $Y$

# Supervised Learning Approach: Behavior Cloning

Expert Trajectories

Dataset

$X$ ⋮ $Y$

$(x_i, y_i)$

$M$

Supervised Learning

12

# Supervised Learning Approach: Behavior Cloning



Expert Trajectories

Dataset

$X$ ⋮ $Y$

Learned Policy $\pi$

*Mapping from state (image) to control (steering direction)*

$M$ $(x_i, y_i)$

Supervised Learning

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

Take a sentence from the web:

State                                Action

*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.*

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

Take a sentence from the web:

<span style="color:green">State</span>           <span style="color:red">Action</span>

*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.*

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

Take a sentence from the web:

State                                                              Action

*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…*

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

Discounted infinite horizon MDP $\mathcal{M} = \{S, A, \gamma, r, P, \rho, \pi^\star\}$

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

Discounted infinite horizon MDP $\mathcal{M} = \{S, A, \gamma, r, P, \rho, \pi^\star\}$

Ground truth reward $r(s, a) \in [0,1]$ is unknown;

For simplicity, let's assume expert is a (nearly) optimal policy $\pi^\star$

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

Discounted infinite horizon MDP $\mathcal{M} = \{S, A, \gamma, r, P, \rho, \pi^\star\}$

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

We have a dataset $\mathscr{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d^{\pi^\star}$

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

Discounted infinite horizon MDP $\mathcal{M} = \{S, A, \gamma, r, P, \rho, \pi^\star\}$

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

We have a dataset $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^{M} \sim d^{\pi^\star}$

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

# Let's formalize the Behavior Cloning algorithm

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

BC is a Reduction to Supervised Learning:

# Let's formalize the Behavior Cloning algorithm

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

BC is a Reduction to Supervised Learning:

$$\widehat{\pi} = \arg\min_{\pi \in \Pi} \sum_{i=1}^{M} \ell\left(\pi, s^{\star}, a^{\star}\right)$$

# Let's formalize the Behavior Cloning algorithm

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

BC is a Reduction to Supervised Learning:

$$\widehat{\pi} = \arg\min_{\pi \in \Pi} \sum_{i=1}^{M} \ell\left(\pi, s^\star, a^\star\right)$$

Many choices of loss functions:

# Let's formalize the Behavior Cloning algorithm

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

BC is a Reduction to Supervised Learning:

$$\widehat{\pi} = \arg\min_{\pi \in \Pi} \sum_{i=1}^{M} \ell\left(\pi, s^{\star}, a^{\star}\right)$$

Many choices of loss functions:

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

# Let's formalize the Behavior Cloning algorithm

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

BC is a Reduction to Supervised Learning:

$$\hat{\pi} = \arg\min_{\pi \in \Pi} \sum_{i=1}^{M} \ell\left(\pi, s^\star, a^\star\right)$$

Many choices of loss functions:

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): $\ell(\pi, s^\star, a^\star) = \|\pi(s^\star) - a^\star\|_2^2$

$$\widehat{\pi} = \arg\min_{\pi \in \Pi} \sum_{i=1}^{M} \ell\left(\pi, s^{\star}, a^{\star}\right)$$

**Analysis**

Assumption: we are going to assume Supervised Learning succeeded

$$\hat{\pi} = \arg\min_{\pi \in \Pi} \sum_{i=1}^{M} \ell\left(\pi, s^\star, a^\star\right)$$

# Analysis

Assumption: we are going to assume Supervised Learning succeeded

$$\mathbb{E}_{s \sim d^{\pi^\star}} \mathbf{1}\left[\hat{\pi}(s) \neq \pi^\star(s)\right] \leq \epsilon \in \mathbb{R}^+$$

$d^{\hat{\pi}^\star} =$

$$\widehat{\pi} = \arg\min_{\pi \in \Pi} \sum_{i=1}^{M} \ell\left(\pi, s^{\star}, a^{\star}\right)$$

# Analysis

Assumption: we are going to assume Supervised Learning succeeded

$$\mathbb{E}_{s \sim d^{\pi^{\star}}} \mathbf{1}\left[\widehat{\pi}(s) \neq \pi^{\star}(s)\right] \leq \epsilon \in \mathbb{R}^{+}$$



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

# Outline for today:

1. Offline Imitation Learning: Behavior Cloning

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

# Performance Difference Lemma

Given two policies $\pi : S \mapsto \Delta(A)$, $\pi' : S \mapsto \Delta(A)$, recall $V^\pi(s_0) = \mathbb{E}\left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid \pi\right]$
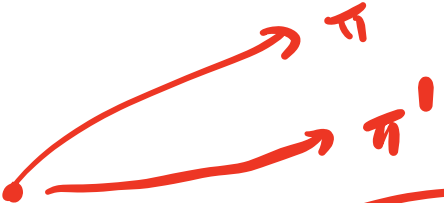
$$V^\pi - V^{\pi'}$$

# Performance Difference Lemma

Given two policies $\pi : S \mapsto \Delta(A)$, $\pi' : S \mapsto \Delta(A)$, recall $V^\pi(s_0) = \mathbb{E}\left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid \pi\right]$

$S_0 \leftarrow$ Intial state

**Performance Difference Lemma (PDL):**

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

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

# PDL Explanation

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

# PDL Proof

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

Proof of Sketch (see reading material for detailed steps)



$\theta^{\pi'}(s_0, a_0) - V^{\pi'}(s_0)$
$= A^{\pi'}(s_0, a_0)$

# PDL Proof

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

Proof of Sketch (see reading material for detailed steps)

$V^{\pi}(s_0) - V^{\pi'}(s_0)$

$= V^{\pi}(s_0) - \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[ r(s_0, a_0) + \gamma \mathbb{E}_{s' \sim P(s_0, a_0)} V^{\pi'}(s') \right] + \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[ r(s_0, a_0) + \gamma \mathbb{E}_{s' \sim P(s_0, a_0)} V^{\pi'}(s') \right] - V^{\pi'}(s_0)$
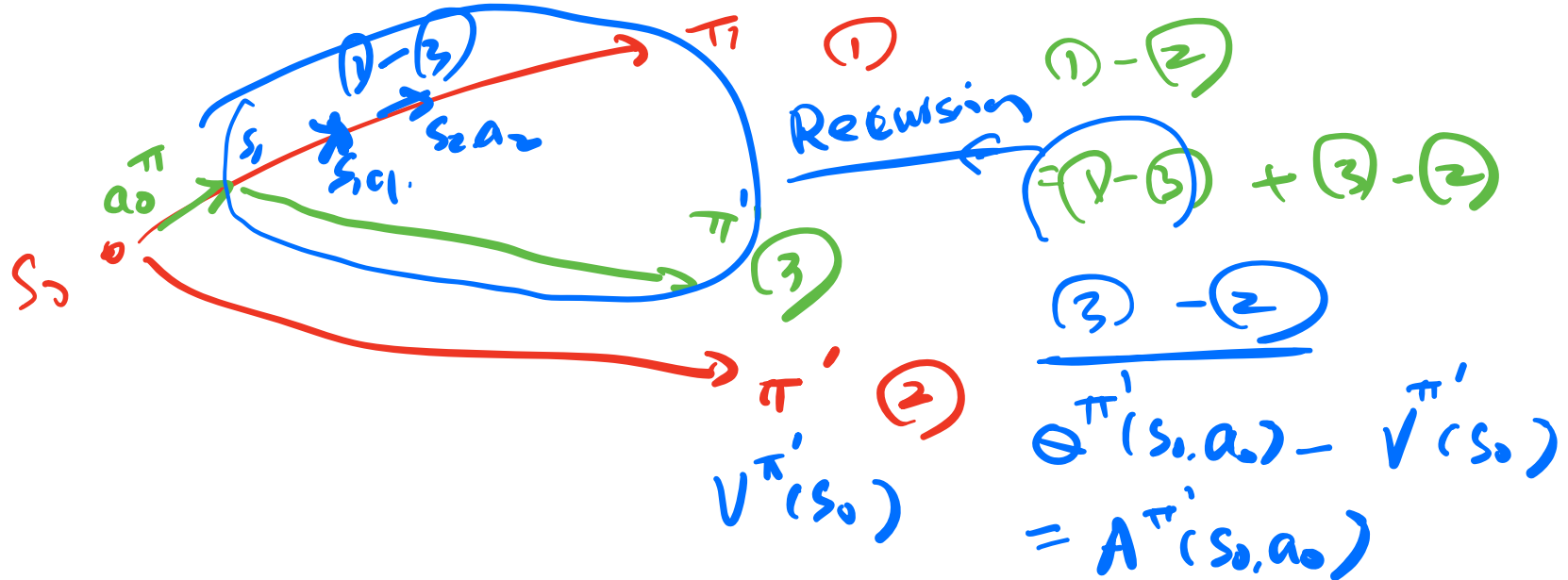
# PDL Proof

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

Proof of Sketch (see reading material for detailed steps)

$$V^{\pi}(s_0) - V^{\pi'}(s_0)$$

$$= V^{\pi}(s_0) - \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[ r(s_0, a_0) + \gamma \mathbb{E}_{s' \sim P(s_0, a_0)} V^{\pi'}(s') \right] + \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[ r(s_0, a_0) + \gamma \mathbb{E}_{s' \sim P(s_0, a_0)} V^{\pi'}(s') \right] - V^{\pi'}(s_0)$$

$$= \gamma \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \mathbb{E}_{s_1 \sim P(s_0, a_0)} \left[ V^{\pi}(s_1) - V^{\pi'}(s_1) \right] + \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[ r(s_0, a_0) + \gamma \mathbb{E}_{s' \sim P(s_0, a_0)} V^{\pi'}(s') \right] - V^{\pi'}(s_0)$$

# PDL Proof

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

Proof of Sketch (see reading material for detailed steps)

$$V^{\pi}(s_0) - V^{\pi'}(s_0)$$

$$= V^{\pi}(s_0) - \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[ r(s_0, a_0) + \gamma \mathbb{E}_{s' \sim P(s_0,a_0)} V^{\pi'}(s') \right] + \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[ r(s_0, a_0) + \gamma \mathbb{E}_{s' \sim P(s_0,a_0)} V^{\pi'}(s') \right] - V^{\pi'}(s_0)$$

$$= \gamma \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \mathbb{E}_{s_1 \sim P(s_0,a_0)} \left[ V^{\pi}(s_1) - V^{\pi'}(s_1) \right] + \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[ r(s_0, a_0) + \gamma \mathbb{E}_{s' \sim P(s_0,a_0)} V^{\pi'}(s') \right] - V^{\pi'}(s_0)$$

$$= \gamma \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \mathbb{E}_{s_1 \sim P(s_0,a_0)} \left[ V^{\pi}(s_1) - V^{\pi'}(s_1) \right] + \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[ Q^{\pi'}(s_0, a_0) - V^{\pi'}(s_0) \right]$$

# PDL Proof

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

Proof of Sketch (see reading material for detailed steps)

$V^{\pi}(s_0) - V^{\pi'}(s_0)$

$= V^{\pi}(s_0) - \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[ r(s_0, a_0) + \gamma \mathbb{E}_{s' \sim P(s_0, a_0)} V^{\pi'}(s') \right] + \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[ r(s_0, a_0) + \gamma \mathbb{E}_{s' \sim P(s_0, a_0)} V^{\pi'}(s') \right] - V^{\pi'}(s_0)$

$= \gamma \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \mathbb{E}_{s_1 \sim P(s_0, a_0)} \left[ V^{\pi}(s_1) - V^{\pi'}(s_1) \right] + \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[ r(s_0, a_0) + \gamma \mathbb{E}_{s' \sim P(s_0, a_0)} V^{\pi'}(s') \right] - V^{\pi'}(s_0)$

$= \gamma \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \mathbb{E}_{s_1 \sim P(s_0, a_0)} \left[ V^{\pi}(s_1) - V^{\pi'}(s_1) \right] + \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[ Q^{\pi'}(s_0, a_0) - V^{\pi'}(s_0) \right]$

$= \gamma \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \mathbb{E}_{s_1 \sim P(s_0, a_0)} \left[ V^{\pi}(s_1) - V^{\pi'}(s_1) \right] + \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[ A^{\pi'}(s_0, a_0) \right]$

# An Application of PDL in Policy Iteration

Recall that $\pi^{t+1}(s) = \arg\max_a A^{\pi^t}(s,a)$

Show monotonic improvement using PDL:

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

$$V^{\pi^{t+1}} - V^{\pi^t} \geq 0$$

$$= \frac{1}{1-\gamma} \mathop{E}_{s \sim d^{\pi^{t+1}}} \mathop{E}_{a \sim \pi^{t+1}} A^{\pi^t}(s,a)$$

$$= \frac{1}{1-\gamma} \mathop{E}_{s \sim d^{\pi^{t+1}}} \left( A^{\pi^t}(s, \pi^{t+1}(s)) \right) \geq 0$$

$\pi^{t+1}$ is deterministic

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

# An Application of PDL in Policy Iteration

Recall that $\pi^{t+1}(s) = \arg\max_a A^{\pi^t}(s, a)$

Show monotonic improvement using PDL:

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

# An Application of PDL in Policy Iteration

Recall that $\pi^{t+1}(s) = \arg\max_a A^{\pi^t}(s, a)$

Show monotonic improvement using PDL:

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

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

# An Application of PDL in Policy Iteration

Recall that $\pi^{t+1}(s) = \arg\max_{a} A^{\pi^t}(s, a)$

Show monotonic improvement using PDL:

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

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

$$\geq \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{s_0}^{\pi^{t+1}}} A^{\pi^t}(s, \pi^t(s))$$

# An Application of PDL in Policy Iteration

Recall that $\pi^{t+1}(s) = \arg\max_{a} A^{\pi^t}(s, a)$

Show monotonic improvement using PDL:

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

$$= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi^{t+1}}_{s_0}} A^{\pi^t}(s, \pi^{t+1}(s))$$

$$\geq \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi^{t+1}}_{s_0}} A^{\pi^t}(s, \pi^t(s)) = 0$$

# Analysis of BC

$$\epsilon \geq \mathbb{E}_{s \sim d^{\pi^\star}} \mathbb{1}\{\hat{\pi}(s) \neq \pi^\star(s)\}$$
$$0$$

Theorem [BC Performance] With probability at least $1 - \delta$, BC returns a policy $\hat{\pi}$:

$$V^{\pi^\star} - V^{\hat{\pi}} \leq \frac{2}{(1 - \gamma)^2} \epsilon$$

# Analysis of BC

Theorem [BC Performance] With probability at least $1 - \delta$, BC returns a policy $\widehat{\pi}$:

$$V^{\pi^\star} - V^{\widehat{\pi}} \leq \frac{2}{(1 - \gamma)^2}\epsilon$$

PDL

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

0

# Analysis of BC

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

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

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

$$\underbrace{\qquad}_{a} \qquad \underbrace{\qquad}_{b}$$

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

$$\approx 0$$

$$\text{①} \quad \pi^\star(s) = \widehat{\pi}(s)$$

$$a = b \approx 0$$

$$\text{②} \quad \pi^\star(s) \neq \widehat{\pi}(s)$$

$$\left| A^{\widehat{\pi}}(s, \pi^\star(s)) - A^{\widehat{\pi}}(s, \widehat{\pi}(s)) \right| \leq \frac{2}{1-\gamma} \quad \left( \because \begin{array}{c} A^{\pi}(s,a) \\ \in \left\{ -\frac{1}{1-\gamma}, \frac{1}{1-\gamma} \right\} \end{array} \right)$$

# Analysis of BC

Theorem [BC Performance] With probability at least $1 - \delta$, BC returns a policy $\hat{\pi}$:

$$V^{\pi^\star} - V^{\hat{\pi}} \leq \frac{2}{(1-\gamma)^2}\epsilon$$

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

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

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

$$\left| A^\pi(s,a) - A^\pi(s,a') \right| \qquad \epsilon$$

# Analysis of BC

Theorem [BC Performance] With probability at least $1 - \delta$, BC returns a policy $\widehat{\pi}$:

$$V^{\pi^\star} - V^{\widehat{\pi}} \leq \frac{2}{(1-\gamma)^2}\epsilon$$

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

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

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

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

# Analysis of BC

Theorem [BC Performance] With probability at least $1 - \delta$, BC returns a policy $\widehat{\pi}$:

$$V^{\pi^\star} - V^{\widehat{\pi}} \leq \frac{2}{(1 - \gamma)^2}\epsilon$$

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

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

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

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

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

# Summary

$\{s, a\} \sim \pi^*$

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

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