# Interactive Imitation Learning (continue)

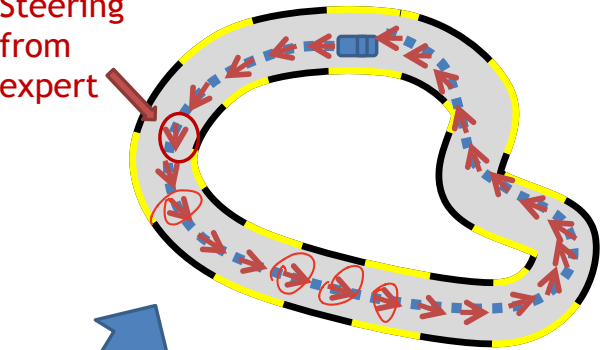# Recap

**Interactive Imitation Learning Setting**

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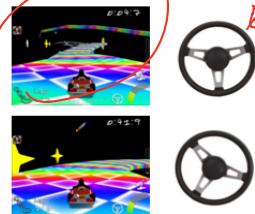we can query expert $\pi^\star$ at any time and any state during training

# DAgger Revisit

At iteration $t$, given $\pi^t$



$d_\mu^{\pi^t}$

$\pi^*(s)$

**New Data**

**All previous data**

Aggregate Dataset

**Supervised Learning**

New policy

$\pi^{t+1}$

Steering from expert

# DAgger Revisit

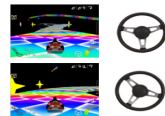At iteration t, given $\pi^t$     **New Data**



Steering from expert

Aggregate Dataset

**+**

All previous data

New policy

**Supervised Learning**

# DAgger Revisit

At iteration t, given $\pi^t$

**New Data**



$$1\left\{\pi(s) \neq \pi^\star(s)\right\}$$

$$\ell_t(\pi) = \mathbb{E}_{s \sim d_\mu^{\pi^t}}\left[c\left(\underset{\Delta}{\pi}, s, \pi^\star(s)\right)\right]$$

**Steering from expert**

Aggregate Dataset

**All previous data**

**New policy**

**Supervised Learning**

# DAgger Revisit

At iteration t, given $\pi^t$

**New Data**

Steering from expert

$$\ell_t(\pi) = \mathbb{E}_{s \sim d_\mu^{\pi^t}}\left[c\left(\pi, s, \pi^\star(s)\right)\right]$$

Aggregate Dataset
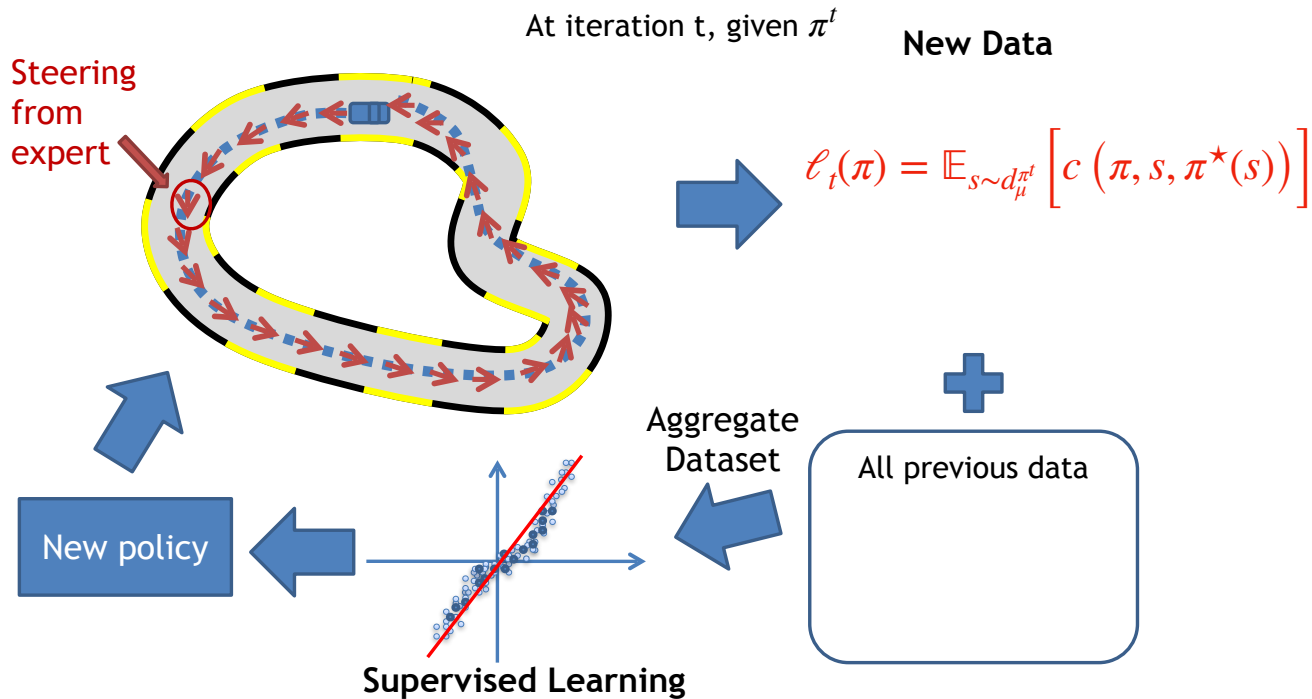
**+**

All previous data

New policy

**Supervised Learning**

# DAgger Revisit

At iteration t, given $\pi^t$

**New Data**

$$\ell_t(\pi) = \mathbb{E}_{s \sim d_\mu^{\pi^t}}\Big[ c\left(\pi, s, \pi^\star(s)\right) \Big]$$

Steering from expert



Aggregate Dataset

**+**

All previous data

$$\sum_{i=0}^{t} \ell_i(\pi)$$

New policy

**Supervised Learning**

# DAgger Revisit

At iteration t, given $\pi^t$

**New Data**

$$\ell_t(\pi) = \mathbb{E}_{s \sim d_\mu^{\pi^t}} \left[ c\left(\pi, s, \pi^\star(s)\right) \right]$$

Steering from expert

**New policy**

Aggregate Dataset

All previous data

$$\sum_{i=0}^{t} \ell_i(\pi)$$

**Supervised Learning**

# DAgger Revisit

At iteration t, given $\pi^t$

**New Data**

**Steering from expert**



$$\ell_t(\pi) = \mathbb{E}_{s \sim d_\mu^{\pi^t}} \left[ c \left( \pi, s, \pi^\star(s) \right) \right]$$

Aggregate Dataset

All previous data

$$\sum_{i=0}^{t} \ell_i(\pi)$$

**New policy**

$$\pi^{t+1} = \arg \min_\pi \sum_{i=0}^{t} \ell_i(\pi) + \lambda R(\pi)$$

**Supervised Learning**

# DAgger Revisit

At iteration t, given $\pi^t$

**New Data**



**Steering from expert**

$$\ell_t(\pi) = \mathbb{E}_{s \sim d_\mu^{\pi^t}}\left[c\left(\pi, s, \pi^\star(s)\right)\right]$$

Aggregate Dataset

All previous data

$$\sum_{i=0}^{t} \ell_i(\pi)$$

**New policy**

$$\pi^{t+1} = \arg\min_\pi \sum_{i=0}^{t} \ell_i(\pi) + \lambda R(\pi)$$

**Supervised Learning**

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

## Recap on the Follow-the-Regularized Leader Guarantee:

At the end of iteration $t$, learner has seen $\ell_0, \ldots \ell_{t-1}, \ell_t$, learner updates to a new decision:

$$\textbf{FTL: } \theta_{t+1} = \min_{\theta \in \Theta} \sum_{i=0}^{t} \ell_i(\theta) + \lambda R(\theta)$$

# Recap on the Follow-the-Regularized Leader Guarantee:

At the end of iteration $t$, learner has seen $\ell_0, \ldots \ell_{t-1}, \ell_t$, learner updates to a new decision:

$$\textbf{FTL: } \theta_{t+1} = \min_{\theta \in \Theta} \sum_{i=0}^{t} \ell_i(\theta) + \lambda R(\theta)$$

↑ Data Aggregation

**Theorem (FTL) (optional):** if $\Theta$ is convex, and $\ell_t$ is convex for all $t$, and $R(\theta)$ is strongly convex, then for regret of FTL, we have:

$$\frac{1}{T} \left[ \sum_{t=0}^{T-1} \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=0}^{T-1} \ell_t(\theta) \right] = O\left( 1/\sqrt{T} \right)$$

$T \to \infty$

↑ Best in Hindsight

**Today's Plan**

1. Finish DAgger's Analysis

2. Intro to Maximum Entropy Inverse RL
(We have offline demonstrations, but learner can interact with the environments)

# DAgger Analysis: A reduction to no-regret online learning

infinite horizon MDP
(assume discrete action space—**in fact let's assume 2 actions**, so we do binary classification)

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

# DAgger Analysis: A reduction to no-regret online learning

infinite horizon MDP
(assume discrete action space—**in fact let's assume 2 actions**, so we do binary classification)

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

$A = \{-1, +1\}$

**Function approximation:**

Decision set $\Pi := \{\pi : S \mapsto A\}$ (assume $\pi^\star \in \Pi$)

← Realizability

$A \to$ Binary classifier

# DAgger Analysis: A reduction to no-regret online learning

infinite horizon MDP
(assume discrete action space—**in fact let's assume 2 actions**, so we do binary classification)

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

**Function approximation:**

Decision set $\Pi := \{\pi : S \mapsto A\}$ (assume $\pi^\star \in \Pi$)

**Classification algorithm (oracle) $\mathcal{A}$:** ← SVM

Given a binary-class data distribution $\rho$, where $\{x, y\} \sim \rho, y \in \{-1, 1\}$

label

A feature

$$\widehat{\pi} = \mathcal{A}\left(\Pi, \rho\right) := \arg\min_{\pi \in \Pi} \mathbb{E}_{x,y\sim\rho}\left[c\left(\pi, x, y\right)\right]$$

# DAgger Analysis: A reduction to no-regret online learning

infinite horizon MDP
(assume discrete action space—**in fact let's assume 2 actions**, so we do binary classification)

$$\mathcal{M} = \left\{ S, A, \gamma, r, P, \mu \right\}$$
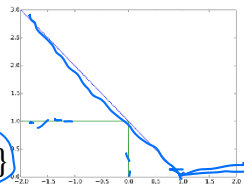
**Function approximation:**

Decision set $\Pi := \{\pi : S \mapsto A\}$ (assume $\pi^\star \in \Pi$)

**Classification algorithm (oracle) $\mathscr{A}$:**

Given a binary-class data distribution $\rho$, where $\{x, y\} \sim \rho, y \in \{-1, 1\}$

$$\widehat{\pi} = \mathscr{A}\left(\Pi, \rho\right) := \arg\min_{\pi \in \Pi} \mathbb{E}_{x,y \sim \rho}\left[c\left(\pi, x, y\right)\right]$$

$$c(\pi, x, y) = \max\{0, 1 - \pi(x) \cdot y\}$$

# DAgger Analysis: A reduction to no-regret online learning

Decision set $\Pi$ (assume $\pi^\star \in \Pi$)

$\Delta$

MDP and
Expert

Online Learner w/ $\mathscr{A}$
(i.e., DAgger)

…

Total loss so far:

# DAgger Analysis: A reduction to no-regret online learning

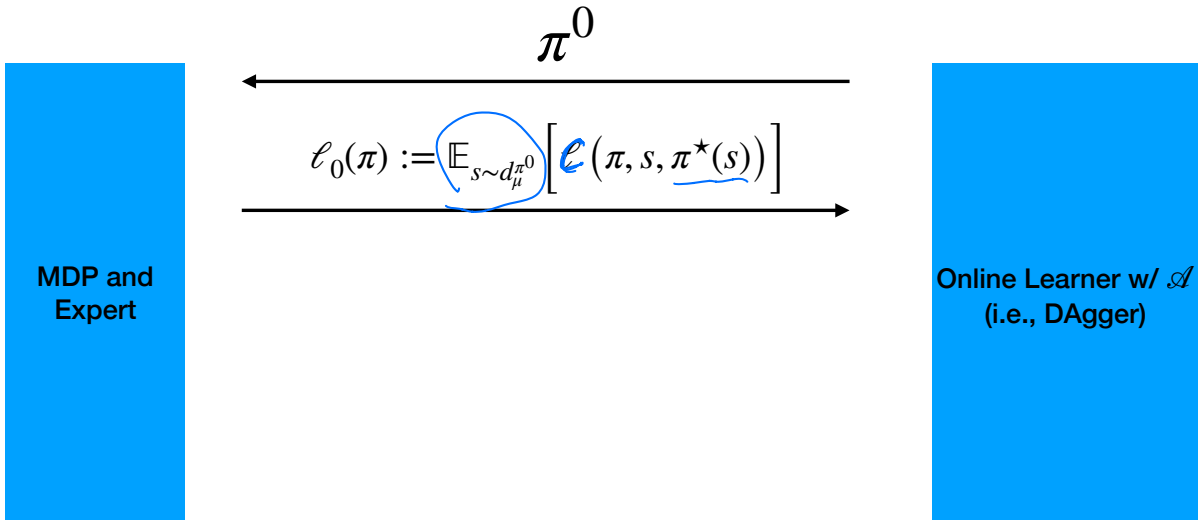Decision set $\Pi$ (assume $\pi^\star \in \Pi$)

$$\pi^0 \in \Pi$$

MDP and Expert

Online Learner w/ $\mathcal{A}$ (i.e., DAgger)

...

Total loss so far:

# DAgger Analysis: A reduction to no-regret online learning

Decision set $\Pi$ (assume $\pi^{\star} \in \Pi$)

$$\pi^0$$

$$\ell_0(\pi) := \mathbb{E}_{s \sim d_\mu^{\pi^0}} \left[ \ell \left( \pi, s, \pi^{\star}(s) \right) \right]$$

MDP and
Expert

Online Learner w/ $\mathcal{A}$
(i.e., DAgger)

...

Total loss so far:

# DAgger Analysis: A reduction to no-regret online learning
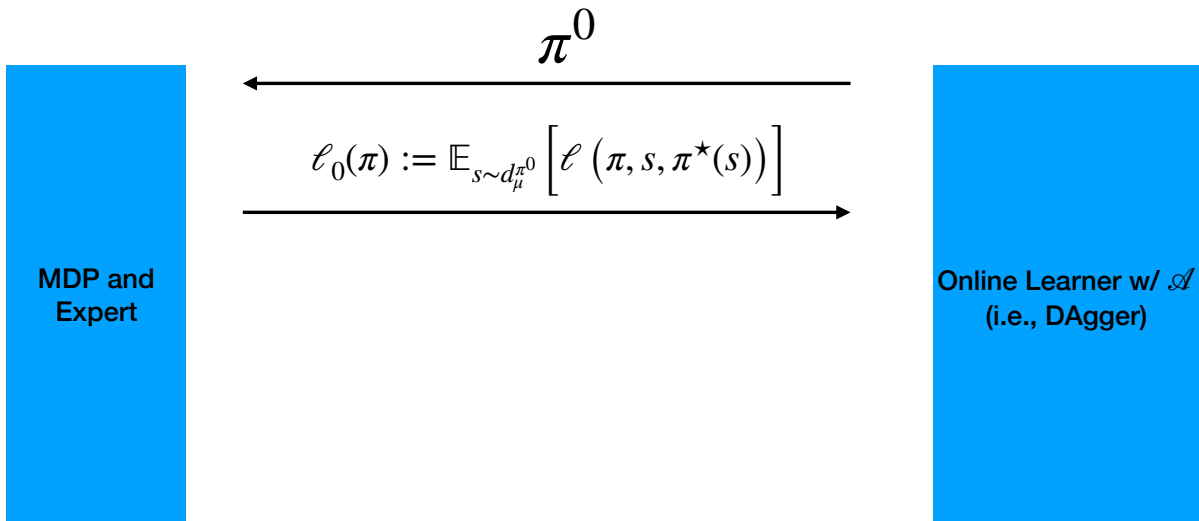
Decision set $\Pi$ (assume $\pi^\star \in \Pi$)

$$\pi^0$$

$$\ell_0(\pi) := \mathbb{E}_{s \sim d_\mu^{\pi^0}} \left[ \ell \left( \pi, s, \pi^\star(s) \right) \right]$$

MDP and
Expert

Online Learner w/ $\mathscr{A}$
(i.e., DAgger)

...

Total loss so far: $\ell_0(\pi^0)$

# DAgger Analysis: A reduction to no-regret online learning

<span style="color:red">Decision set $\Pi$ (assume $\pi^\star \in \Pi$)</span>

$$\pi^0$$

$$\ell_0(\pi) := \mathbb{E}_{s \sim d_\mu^{\pi^0}} \left[ \ell \left( \pi, s, \pi^\star(s) \right) \right]$$

$$\pi^1$$

**MDP and Expert**

**Online Learner w/ $\mathcal{A}$ (i.e., DAgger)**

...

Total loss so far:  $\ell_0(\pi^0)$

# DAgger Analysis: A reduction to no-regret online learning

Decision set $\Pi$ (assume $\pi^\star \in \Pi$)



$$\pi^0$$

$$\ell_0(\pi) := \mathbb{E}_{s \sim d_\mu^{\pi^0}} \left[ \ell \left( \pi, s, \pi^\star(s) \right) \right]$$

**MDP and Expert**

$$\pi^1$$

**Online Learner w/ $\mathscr{A}$ (i.e., DAgger)**

$$\ell_1(\pi) := \mathbb{E}_{s \sim d_\mu^{\pi^1}} \left[ \ell \left( \pi, s, \pi^\star(s) \right) \right]$$

...

Total loss so far:  $\ell_0(\pi^0)$

# DAgger Analysis: A reduction to no-regret online learning
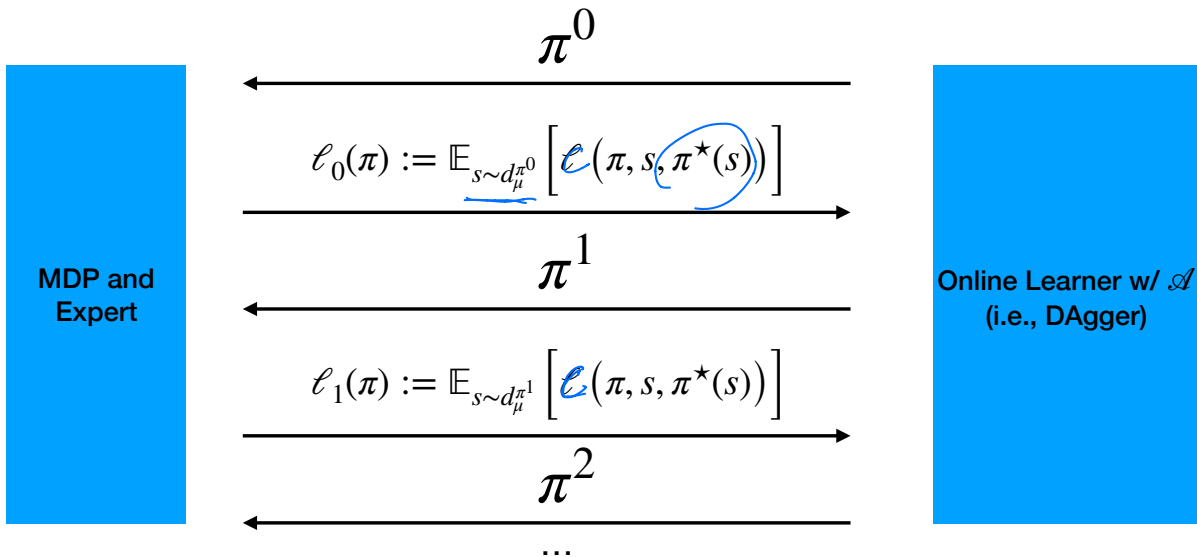
Decision set $\Pi$ (assume $\pi^\star \in \Pi$)

$$\pi^0$$

MDP and Expert

$$\ell_0(\pi) := \mathbb{E}_{s \sim d_\mu^{\pi^0}} \left[ \ell \left( \pi, s, \pi^\star(s) \right) \right]$$

$$\pi^1$$

$$\ell_1(\pi) := \mathbb{E}_{s \sim d_\mu^{\pi^1}} \left[ \ell \left( \pi, s, \pi^\star(s) \right) \right]$$

Online Learner w/ $\mathscr{A}$ (i.e., DAgger)

...

Total loss so far: $\ell_0(\pi^0)$ $\left( +\ell_1(\pi^1) \right)$

# DAgger Analysis: A reduction to no-regret online learning

Decision set $\Pi$ (assume $\pi^\star \in \Pi$)

$\pi^0$

$\ell_0(\pi) := \mathbb{E}_{s \sim d_\mu^{\pi^0}} \left[ \ell \left( \pi, s, \pi^\star(s) \right) \right]$

MDP and Expert

$\pi^1$

$\ell_1(\pi) := \mathbb{E}_{s \sim d_\mu^{\pi^1}} \left[ \ell \left( \pi, s, \pi^\star(s) \right) \right]$

Online Learner w/ $\mathscr{A}$ (i.e., DAgger)

$\pi^2$

...

Total loss so far:   $\ell_0(\pi^0)$  $+\ell_1(\pi^1)$

# DAgger Analysis: A reduction to no-regret online learning

Decision set $\Pi$ (assume $\pi^\star \in \Pi$)

*Interactive setting*

$$\pi^0$$

MDP and Expert

Online Learner w/ $\mathcal{A}$ (i.e., DAgger)

$$\ell_0(\pi) := \mathbb{E}_{s \sim d_\mu^{\pi^0}} \left[ \ell\left(\pi, s, \pi^\star(s)\right) \right]$$

$$\pi^1$$

$$\ell_1(\pi) := \mathbb{E}_{s \sim d_\mu^{\pi^1}} \left[ \ell\left(\pi, s, \pi^\star(s)\right) \right]$$

$$\pi^2$$

...

Total loss so far: $\ell_0(\pi^0) \; + \ell_1(\pi^1) \; + \ell_2(\pi^2) + \ldots$

# DAgger Analysis: A reduction to no-regret online learning

After in total $T$ many iterations, we have the following regret for DAgger:

$$\text{Avg-Regret}_T = \frac{1}{T}\left[\sum_{t=0}^{T-1}\ell_t(\pi^t) - \min_{\pi\in\Pi}\sum_{t=0}^{T-1}\ell_t(\pi)\right] \leq O\left(\frac{1}{\sqrt{T}}\right)$$

Total-loss
DAgger suffered

$\epsilon_{reg}$

Goal: Turn $\epsilon_{reg}$

to $V^{\pi} - V^{\pi^*}$

# DAgger Analysis: A reduction to no-regret online learning

After in total $T$ many iterations, we have the following regret for DAgger:

$$\text{Avg-Regret}_T = \frac{1}{T}\left[\sum_{t=0}^{T-1}\ell_t(\pi^t) - \min_{\pi\in\Pi}\sum_{t=0}^{T-1}\ell_t(\pi)\right] \le O\left(\underbrace{\frac{1}{\sqrt{T}}}_{\epsilon_{reg}}\right)$$

Recall we assume $\pi^\star \in \Pi$, we must have:

$$\min_{\pi\in\Pi}\sum_{t=0}^{T-1}\ell_t(\pi) \le \sum_{t=0}^{T-1}\ell_t(\pi^\star) = 0$$

$$1\left\{\pi(s) \ne \pi^*(s)\right\}$$

$$\ell_t(\pi) = \mathbb{E}_{s\sim d_M^{\pi^t}}\left[c(\pi, s, \pi^*(s))\right]$$

$$\ell_t(\pi^*) = 0$$

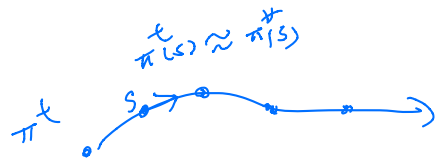# DAgger Analysis: A reduction to no-regret online learning

$$\frac{1}{T}\sum \ell_t(\pi^t) \leq \min_{\pi \in \Pi} \sum \ell_t(\pi) + \varepsilon_{reg} \leq 0 + \varepsilon_{reg}$$

After in total $T$ many iterations, we have the following regret for DAgger:

$$\text{Avg-Regret}_T = \frac{1}{T}\left[\sum_{t=0}^{T-1} \ell_t(\pi^t) - \min_{\pi \in \Pi} \sum_{t=0}^{T-1} \ell_t(\pi)\right] \leq O\left(\frac{1}{\sqrt{T}}\right)$$

$$\underbrace{\phantom{xxxxxxx}}_{\epsilon_{reg}}$$

$$\{\pi, \pi^b\}$$

$$\min_{\pi + (\pi, \pi^b)} \ell(\hat{\pi}) \leq \ell(\pi^b)$$

Recall we assume $\pi^\star \in \Pi$, we must have:

$$\frac{1}{T}\sum_{t=0}^{T-1} \ell_t(\pi^t)$$

$$\min_{\pi \in \Pi} \sum_{t=0}^{T-1} \ell_t(\pi) \leq \sum_{t=0}^{T-1} \ell_t(\pi^\star) = 0$$

$$\leq 0 + \varepsilon_{reg}$$

$$\min \leq \text{Avg}$$

Which implies that:

$$\min_{t \in \{0 \dots T-1\}} \ell_t(\pi^t) \leq \frac{1}{T}\sum_{t=0}^{T-1} \ell_t(\pi^t) \leq \epsilon_{reg}$$

# DAgger Analysis: A reduction to no-regret online learning

Summary so far: we know that there must exists $t \in \{0, \ldots, T-1\}$, such that:

$$\ell_t\left(\pi^t\right) \le \epsilon_{reg} \qquad \left( \min_{t \in \{0, \ldots, T-1\}} \ell_t(\pi^t) \le \epsilon_{reg} \right)$$
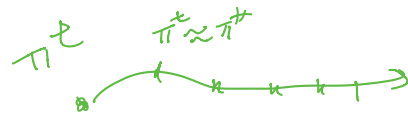
# DAgger Analysis: A reduction to no-regret online learning

Summary so far: we know that there must exists $t \in \{0, \ldots, T-1\}$, such that:

$$\ell_t\left(\pi^t\right) \leq \epsilon_{reg}$$

Recall the definition of $\ell_t(\pi^t)$

$$\ell_t\left(\pi^t\right) = \mathbb{E}_{s \sim d_\mu^{\pi^t}}\left[c\left(\pi^t, s, \pi^\star(s)\right)\right] \leq \epsilon_{reg}$$

# DAgger Analysis: A reduction to no-regret online learning

Summary so far: we know that there must exists $t \in \{0,\ldots,T-1\}$, such that:

$$\ell_t \left( \pi^t \right) \le \epsilon_{reg}$$

$\pi^t(s) \approx \pi^*(s)$

Recall the definition of $\ell_t(\pi^t)$

$$\ell_t \left( \pi^t \right) = \mathbb{E}_{s \sim d_\mu^{\pi^t}} \left[ c \left( \pi^t, s, \pi^\star(s) \right) \right] \le \epsilon_{reg}$$

$\pi^t$ matches to $\pi^\star$ under its own state distribution!

# DAgger Analysis: A reduction to no-regret online learning

Summary so far: we know that there must exists $t \in \{0, \ldots, T-1\}$, such that:

$$\ell_t \left( \pi^t \right) \leq \epsilon_{reg}$$

Recall the definition of $\ell_t(\pi^t)$

$$\ell_t \left( \pi^t \right) = \mathbb{E}_{s \sim d_\mu^{\pi^t}} \left[ c \left( \pi, s, \pi^\star(s) \right) \right] \leq \epsilon_{reg}$$

$\pi^t$ matches to $\pi^\star$ under its own state distribution!

Behavior- cloning

Recall BC, we had:

$\mathbb{E}_{s \sim d^{\pi^\star}} \left[ c(\widehat{\pi}, s, \pi^\star(s)) \right] \leq \epsilon$, i.e., we matched to $\pi^\star$ under $\pi^\star$'s distribution
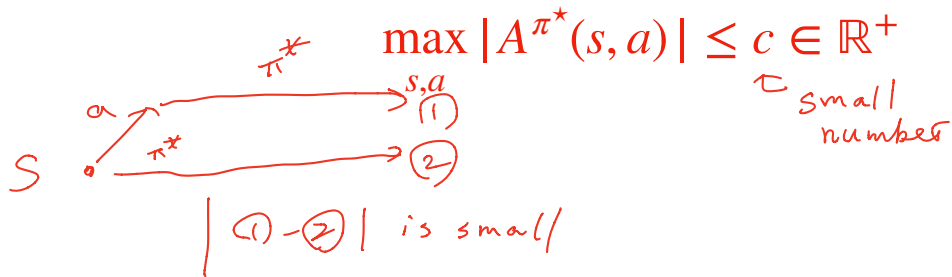
**Finally, turn things into the performance bound using PDL:**

Theorem: There exists a iteration $t$, such that:

$$V^{\pi^\star} - V^{\pi^t} \leq \frac{\max_{s,a} \left| A^{\pi^\star}(s,a) \right|}{1 - \gamma} \cdot \epsilon_{reg}$$

**This bound indicates that:**

**Finally, turn things into the performance bound using PDL:**

Theorem: There exists a iteration $t$, such that:

$$V^{\pi^\star} - V^{\pi^t} \leq \frac{\max_{s,a} \left| A^{\pi^\star}(s,a) \right|}{1 - \gamma} \cdot \epsilon_{reg}$$

**This bound indicates that:**

We **avoid quadratic error** if expert $\pi^\star$ can quickly recover from a mistake

$$\max_{s,a} |A^{\pi^\star}(s,a)| \leq c \in \mathbb{R}^+$$

small number



$s$    $a$    $\pi^\star$    $\xrightarrow{s,a}$ (1)

$\pi^\star$ (2)

$|$ (1) - (2) $|$ is small

# Finally, turn things into the performance bound using PDL:

Theorem: There exists a iteration $t$, such that:

$$V^{\pi^\star} - V^{\pi^t} \le \frac{\max_{s,a} \left| A^{\pi^\star}(s,a) \right|}{1-\gamma} \cdot \epsilon_{reg} \quad \le \quad \frac{c}{1-\gamma} \, \epsilon_{Reg}$$



Recover

$$Q^{\pi^\star}(s_2, a_2) - Q^{\pi^\star}(s_2, a_1)$$

$$\sim 1$$

**This bound indicates that:**

We **avoid quadratic error** if expert $\pi^\star$ can quickly recover from a mistake

$$\max_{s,a} |A^{\pi^\star}(s,a)| \le c \in \mathbb{R}^+$$

i.e., at $s$, taking $a$ then following $\pi^\star$ is almost as good as following $\pi^\star$ directly

# Finally, turn things into the performance bound using PDL:

Theorem: There exists a iteration $t$, such that:

$$V^{\pi^\star} - V^{\pi^t} \leq \frac{\max_{s,a} \left| A^{\pi^\star}(s, a) \right|}{1 - \gamma} \cdot \epsilon_{reg}$$

PDL

$$V^{\pi^t} - V^{\pi^\star} = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^\pi_\mu} \left[ A^{\pi^\star}(s, \pi^t(s)) \right]$$

**This bound indicates that:**

We **avoid quadratic error** if expert $\pi^\star$ can quickly recover from a mistake

$$\max_{s,a} \left| A^{\pi^\star}(s, a) \right| \leq c \in \mathbb{R}^+$$

i.e., at $s$, taking $a$ then following $\pi^\star$ is almost as good as following $\pi^\star$ directly

# Finally, turn things into the performance bound using PDL:

Theorem: There exists a iteration $t$, such that:

$$V^{\pi^\star} - V^{\pi^t} \leq \frac{\max_{s,a}\left|A^{\pi^\star}(s,a)\right|}{1-\gamma} \cdot \epsilon_{reg}$$

$f:$

$\left| f(x) - f(y) \right|$

$\leq \max_{z}\left| f(z) \right| \mathbb{1}(x \neq y)$

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

$A^{\pi^\star}(s, \pi^\star(s))$

$= \frac{1}{1-\gamma}\mathbb{E}_{s \sim d^\pi}\left[A^{\pi^\star}(s, \pi^t(s)) - A^{\pi^\star}(s, \pi^\star(s))\right]$

$= Q^{\pi^\star}(s, \pi^\star(s))$

$\underbrace{\qquad}_{=0}$

$- V^{\pi^\star}(s)$

$= 0$

$\left| A^{\pi^\star}(s, \pi^t(s)) - A^{\pi^\star}(s, \pi^\star(s)) \right|$

$\leq \max_{s \cdot a}\left| A^{\pi^\star}(s,a) \right| \mathbb{1}\left\{ \pi^t(s) \neq \pi^\star(s) \right\}$

**This bound indicates that:**

We **avoid quadratic error** if expert $\pi^\star$ can quickly recover from a mistake

$$\max_{s,a} |A^{\pi^\star}(s,a)| \leq c \in \mathbb{R}^+$$

i.e., at $s$, taking $a$ then following $\pi^\star$ is almost as good as following $\pi^\star$ directly

# Finally, turn things into the performance bound using PDL:

Theorem: There exists a iteration $t$, such that:

$$V^{\pi^\star} - V^{\pi^t} \leq \frac{\max_{s,a}\left|A^{\pi^\star}(s,a)\right|}{1-\gamma} \cdot \epsilon_{reg}$$

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

$$= \frac{1}{1-\gamma}\mathbb{E}_{s\sim d^\pi}\left[A^{\pi^\star}(s, \pi^t(s)) - A^{\pi^\star}(s, \pi^\star(s))\right]$$

$$\geq \frac{-1}{1-\gamma}\mathbb{E}_{s\sim d^{\pi^t}_\mu}\left[\max_{s,a}\left|A^{\pi^\star}(s,a)\right|\mathbf{1}\{\pi^t(s) \neq \pi^\star(s)\}\right]$$

$\geq - \max_{s,a}\left|A^{\pi^\star}(s,a)\right| \mathbf{1}\left(\pi^t(s) \neq \pi^b(s)\right)$

from Dagger: $\ell_t(\pi^t) \leq \varepsilon_{Reg}$

$\Rightarrow \mathbb{E}_{s\sim d^{\pi^t}_\mu} \mathbf{1}\left[\pi^t(s) \neq \pi^b(s)\right] \leq \varepsilon_{Reg}$

**This bound indicates that:**

We **avoid quadratic error** if expert $\pi^\star$ can quickly recover from a mistake

$$\max_{s,a}\left|A^{\pi^\star}(s,a)\right| \leq c \in \mathbb{R}^+$$

i.e., at $s$, taking $a$ then following $\pi^\star$ is almost as good as following $\pi^\star$ directly

# Finally, turn things into the performance bound using PDL:

Theorem: There exists a iteration $t$, such that:

$$V^{\pi^\star} - V^{\pi^t} \leq \frac{\max_{s,a} \left| A^{\pi^\star}(s,a) \right|}{1 - \gamma} \cdot \epsilon_{reg}$$

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

$$= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^\pi} \left[ A^{\pi^\star}(s, \pi^t(s)) - A^{\pi^\star}(s, \pi^\star(s)) \right]$$

$$\geq \frac{-1}{1-\gamma} \mathbb{E}_{s \sim d^t_\mu} \left[ \max_{s,a} \left| A^{\pi^\star}(s,a) \right| \mathbf{1}\{\pi^t(s) \neq \pi^\star(s)\} \right] \geq -\frac{\epsilon_{reg}}{1-\gamma} \cdot \max_{sa} \left\lceil \overset{\pi^t}{A^{(s,a)}} \right\rceil$$

$$V^{\pi^\star} - V^{\pi^t} \leq \frac{\max_{s,a} \left| A^{\pi^\star}(s,a) \right|}{1-\gamma} \cdot \epsilon_{reg} \quad \checkmark$$

**This bound indicates that:**

We **avoid quadratic error** if expert $\pi^\star$ can quickly recover from a mistake

$$\max_{s,a} \left| A^{\pi^\star}(s,a) \right| \leq c \in \mathbb{R}^+$$

i.e., at $s$, taking $a$ then following $\pi^\star$ is almost as good as following $\pi^\star$ directly

# Summary of DAgger

# Summary of DAgger

DAgger finds a policy $\widehat{\pi}$ such that it matches to $\pi^\star$ under $d_\mu^{\widehat{\pi}}$

$$\mathbb{E}_{s \sim d_\mu^{\widehat{\pi}}} \left[ \mathbf{1}\{ \widehat{\pi}(s) \neq \pi^\star(s) \} \right] \leq \epsilon_{reg} = O(1/\sqrt{T})$$

No-Regret
Argument
(FTRL)

# Summary of DAgger

DAgger finds a policy $\widehat{\pi}$ such that it matches to $\pi^\star$ under $d_\mu^{\widehat{\pi}}$

$$\mathbb{E}_{s \sim d_\mu^{\widehat{\pi}}}\left[\mathbf{1}\{\widehat{\pi}(s) \neq \pi^\star(s)\}\right] \leq \epsilon_{reg} = O(1/\sqrt{T})$$

If expert herself can quickly recover from a deviation, i.e., $|Q^{\pi^\star}(s,a) - V^{\pi^\star}(s)|$ is small for all $s$,

$$V^{\pi^\star} - V^{\pi^t} \leq O\left(\frac{1}{1-\gamma} \cdot \epsilon_{reg}\right)$$

$$\overbrace{A}$$
$$A^{\pi^\star}(s,a)$$

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

$$\mathbb{E}_{s \sim d_\mu^{\pi^\star}}\left[A^{\widehat{\pi}}(s, \pi^\star(s))\right]$$

# Today's Plan

✅ 1. Finish DAgger's Analysis

2. Intro to Maximum Entropy Inverse RL
(We have offline demonstrations, but learner can interact with the environments)

# Review of the IL settings that we covered so far

## 1. Offline IL Setting:

We have a dataset $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d^{\pi^\star}$
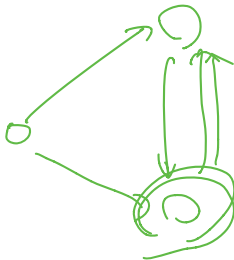
No expert interaction, no real world interaction

# Review of the IL settings that we covered so far

## 1. Offline IL Setting:

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

No expert interaction, no real world interaction

## 2. Interactive IL setting:

We have access to $\pi^\star$ during training

Interaction w/ expert and interaction w/ the world (i.e., we can try out our policies)

# A new setting (more realistic maybe??)

**Hybrid:**

1. We have an offline dataset $\mathscr{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d^{\pi^\star}$ (e.g., a pre-collected demonstrations)

2. And we can interact with the world (e.g., try out our policy and see what happens)
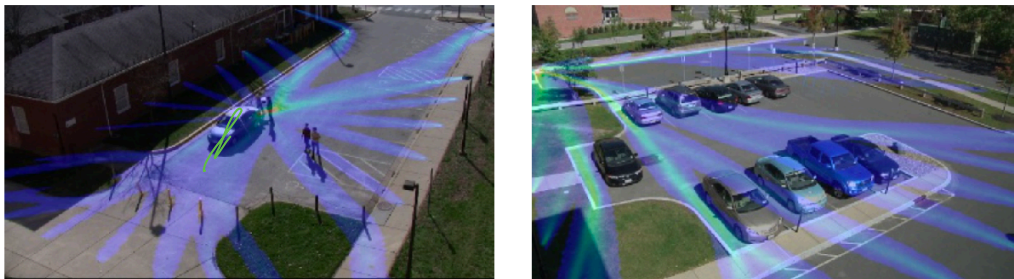
# Running Example: Human trajectory forecasting

[Kitani, et al, ECCV 12]



**Fig. 1.** Given a single pedestrian detection, our proposed approach forecasts plausible paths and destinations from noisy vision-input

# Running Example: Human trajectory forecasting

[Kitani, et al, ECCV 12]



**Fig. 1.** Given a single pedestrian detection, our proposed approach forecasts plausible paths and destinations from noisy vision-input

High-level assumptions:
(1) Experts may have some cost function regarding walking in their mind
(2) Experts are (approximately) optimizing the cost function

# Setting

Finite horizon MDP $\mathcal{M} = \{S, A, H, c, P, \mu, \pi^\star\}$

# Setting

Finite horizon MDP $\mathcal{M} = \{S, A, H, c, P, \mu, \pi^\star\}$

(1) Ground truth cost $c(s, a)$ is unknown;

(2) assume expert is the optimal policy $\pi^\star$ of the cost $c$

(3) **transition P is known**

## Setting

Finite horizon MDP $\mathcal{M} = \{S, A, H, c, P, \mu, \pi^\star\}$

(1) Ground truth cost $c(s, a)$ is unknown;

(2) assume expert is the optimal policy $\pi^\star$ of the cost $c$

(3) **transition P is known**

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

# Setting

Finite horizon MDP $\mathcal{M} = \{S, A, H, c, P, \mu, \pi^\star\}$

(1) Ground truth cost $c(s, a)$ is unknown;

(2) assume expert is the optimal policy $\pi^\star$ of the cost $c$

(3) **transition P is known**

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

**Key Assumption on cost:**

$c(s, a) = \langle \theta^\star, \phi(s, a) \rangle$, **linear w.r.t feature** $\phi(s, a)$

# Running Example: Define feature map

**Key Assumption on cost:**
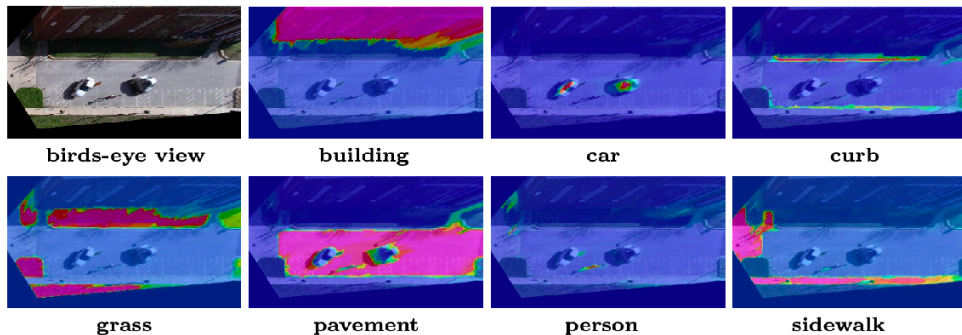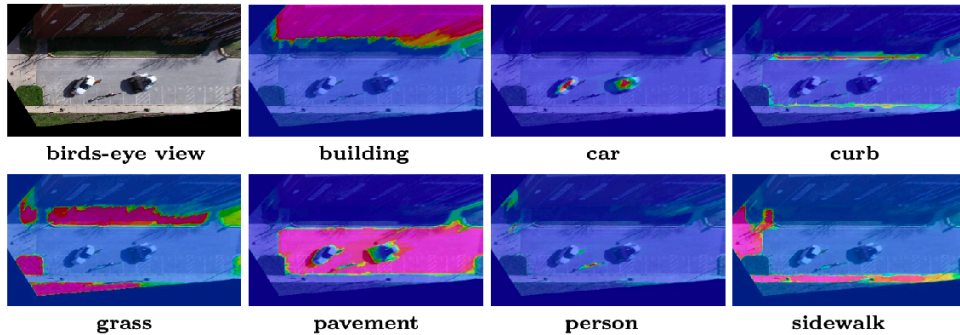$c(s, a) = \langle \theta^\star, \phi(s, a) \rangle$**, linear wrt feature** $\phi(s, a)$



**Fig. 4.** Classifier feature response maps. Top left is the original image.

# Running Example: Define feature map

**Key Assumption on cost:**
$c(s,a) = \langle \theta^\star, \phi(s,a) \rangle$, **linear wrt feature** $\phi(s,a)$

*pixel*

State $s$: pixel or a group of neighboring pixels in image)



**Fig. 4.** Classifier feature response maps. Top left is the original image.

# Running Example: Define feature map

**Key Assumption on cost:**
$c(s, a) = \langle \theta^\star, \phi(s, a) \rangle$**, linear wrt feature** $\phi(s, a)$

State $s$: pixel or a group of neighboring pixels in image)



birds-eye view    building    car    curb

grass    pavement    person    sidewalk

$$\phi(s, a) = \begin{bmatrix} \mathbb{P}(\text{pixels being building}) \\ \mathbb{P}(\text{pixels being grass}) \\ \mathbb{P}(\text{pixels being sidewalk}) \\ \mathbb{P}(\text{pixels being car}) \\ \cdots \end{bmatrix}$$

pixel

**Fig. 4.** Classifier feature response maps. Top left is the original image.

# Running Example: Define feature map

**Key Assumption on cost:**
$c(s,a) = \langle \theta^\star, \phi(s,a) \rangle$, **linear wrt feature** $\phi(s,a)$



**Fig. 4.** Classifier feature response maps. Top left is the original image.

State $s$: pixel or a group of neighboring pixels in image)

$$\phi(s,a) = \begin{bmatrix} \mathbb{P}(\text{pixels being building}) \\ \mathbb{P}(\text{pixels being grass}) \\ \mathbb{P}(\text{pixels being sidewalk}) \\ \mathbb{P}(\text{pixels being car}) \\ \cdots \end{bmatrix}$$

Maybe colliding with cars or buildings has **high** cost, but walking on sideway or grass has **low** cost
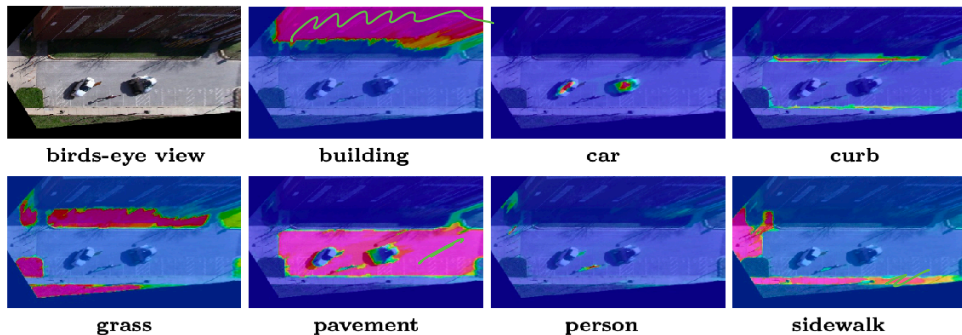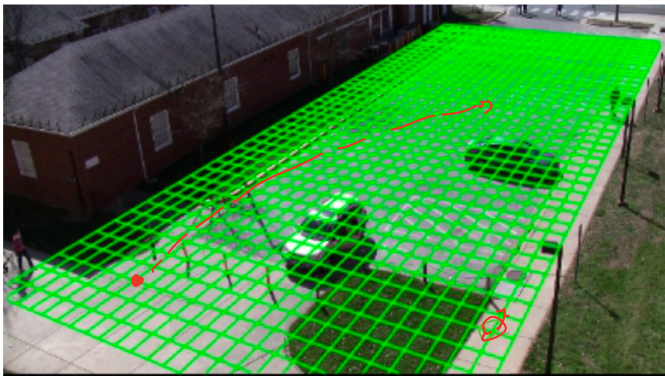
# Running Example: Human Trajectory Forecasting



State space: grid,
action space: 4 actions

We predict that we are more likely to use sidewalk

**We will talk about the algorithm (MaxEnt-IRL) behind it next week**