

Interactive Imitation Learning (continue)

Recap

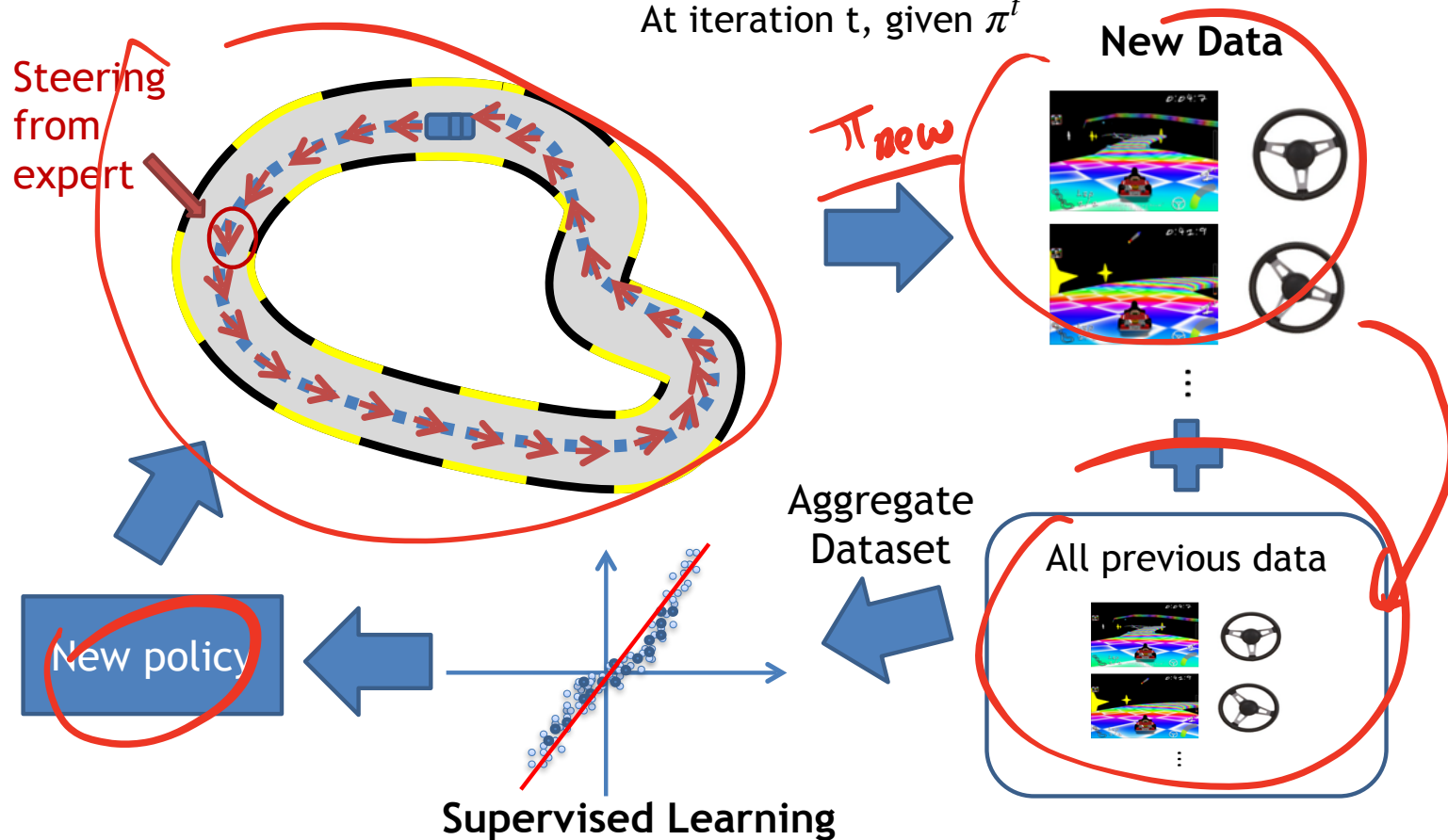
Interactive Imitation Learning Setting

Key assumption:

we can query expert π^\star at any time and any state during training

Dagger Revisit

At iteration t , given π^t



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 $\ell_0, \dots, \ell_{t-1}, \ell_t$, learner updates to a new decision:

$$\text{FTL: } \theta_{t+1} = \min_{\theta \in \Theta} \sum_{i=0}^t \ell_i(\theta) + \lambda R(\theta) \quad ||\theta||_2$$

(Note: The image contains handwritten red annotations: a triangle under FTL, a horizontal line under the summation term, and the handwritten expression ||θ||₂ to the right of the equation.)

Recap on the Follow-the-Regularized Leader Guarantee:

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

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

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

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

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 policy is a binary classifier)

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

DAgger Analysis: A reduction to no-regret online learning

infinite horizon MDP

(assume discrete action space—**in fact let's assume 2 actions**, so policy is a binary classifier)

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

Classification:

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

$$\hat{\pi} = \arg \min_{\pi} \sum_{x, y} \left[\ell(\pi, x, y) \right]$$

loss

$x \in \mathbb{R}^d$

DAgger Analysis: A reduction to no-regret online learning

infinite horizon MDP

(assume discrete action space—in fact let's assume 2 actions, so policy is a binary classifier)

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

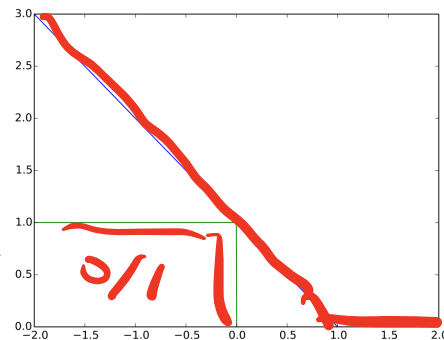
Classification:

SUM

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

$$\hat{\pi} = \arg \min_{\pi} \sum_{x,y} \left[\ell(\pi, x, y) \right]$$

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



DAgger Analysis: A reduction to no-regret online learning



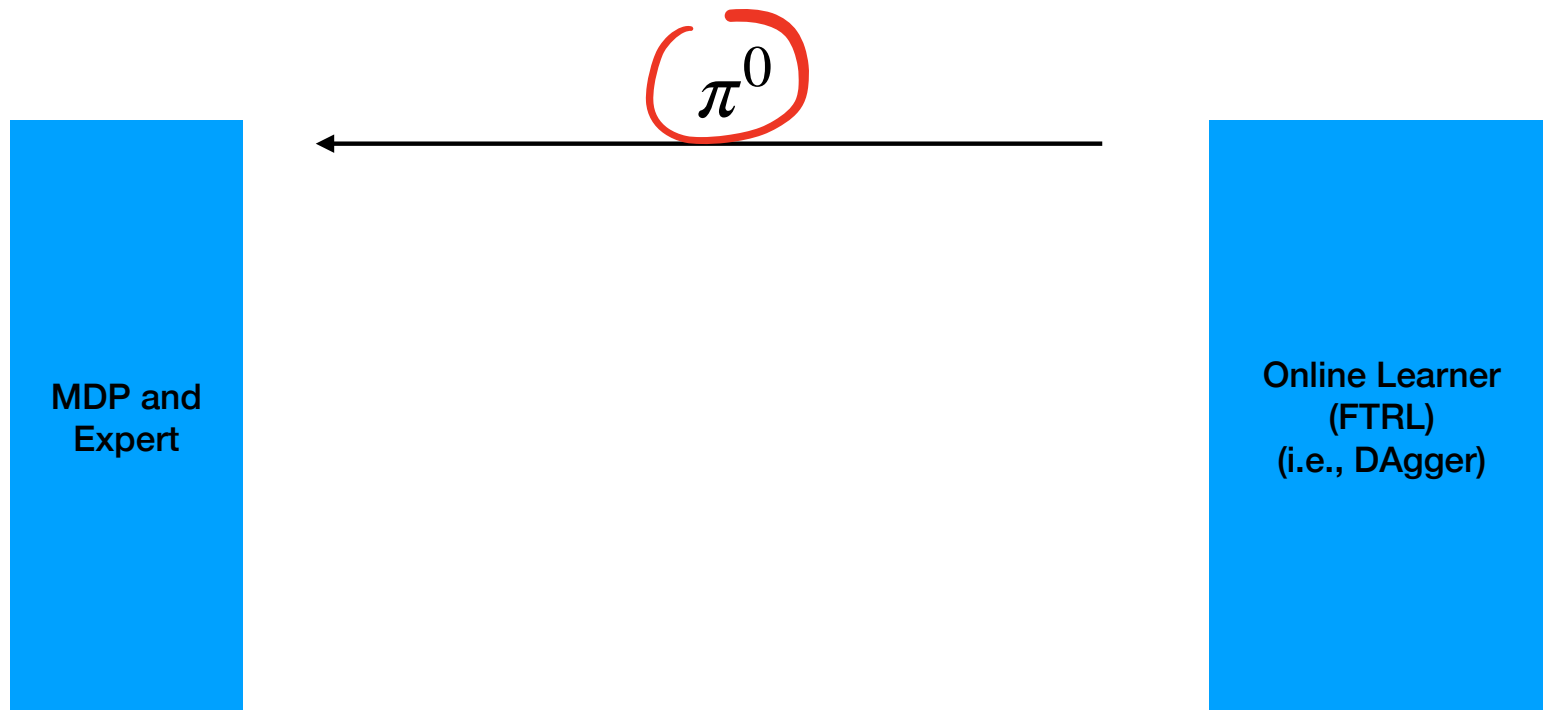
MDP and
Expert

Online Learner
(FTRL)
(i.e., DAgger)

...

Total loss so far:

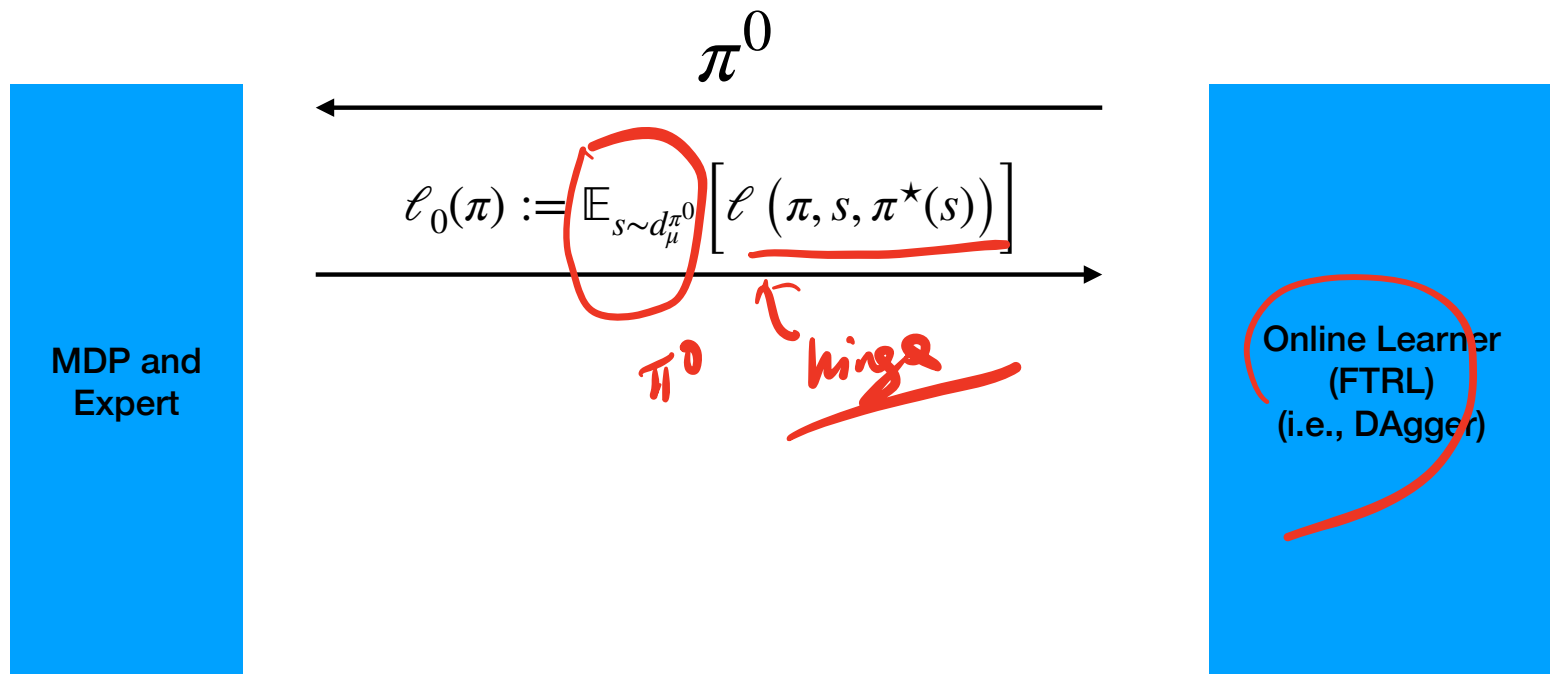
DAgger Analysis: A reduction to no-regret online learning



...

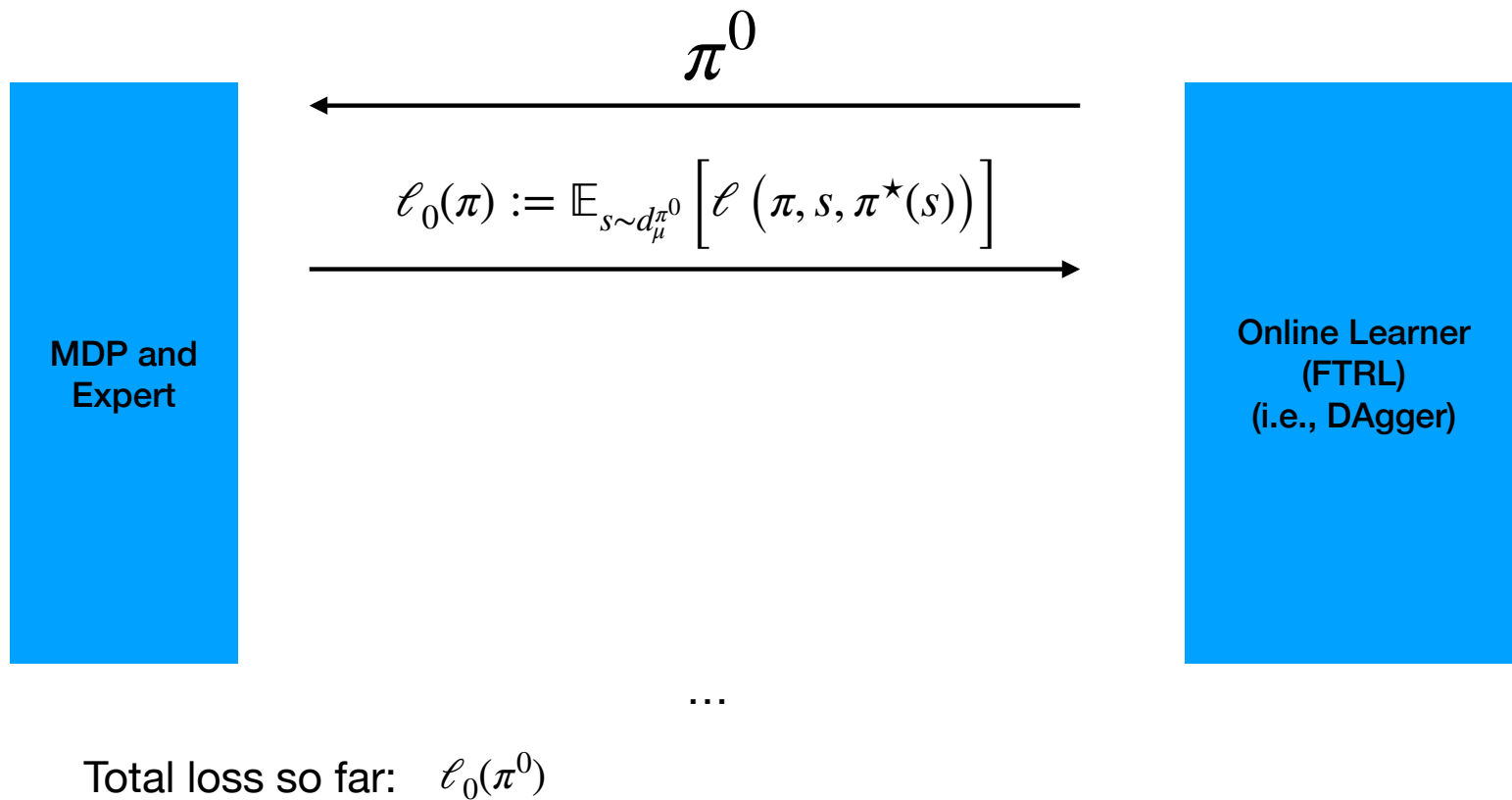
Total loss so far:

DAgger Analysis: A reduction to no-regret online learning

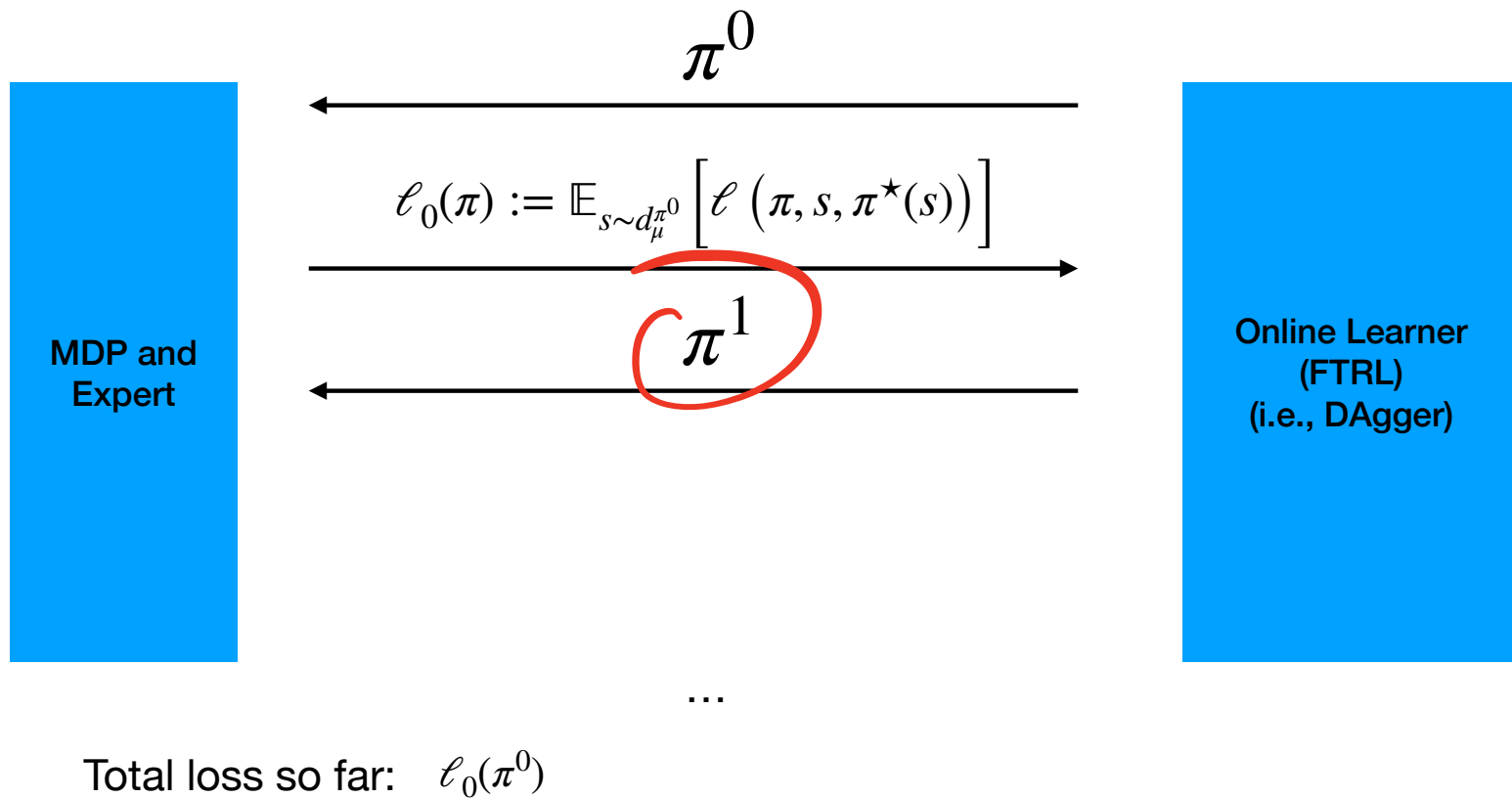


Total loss so far:

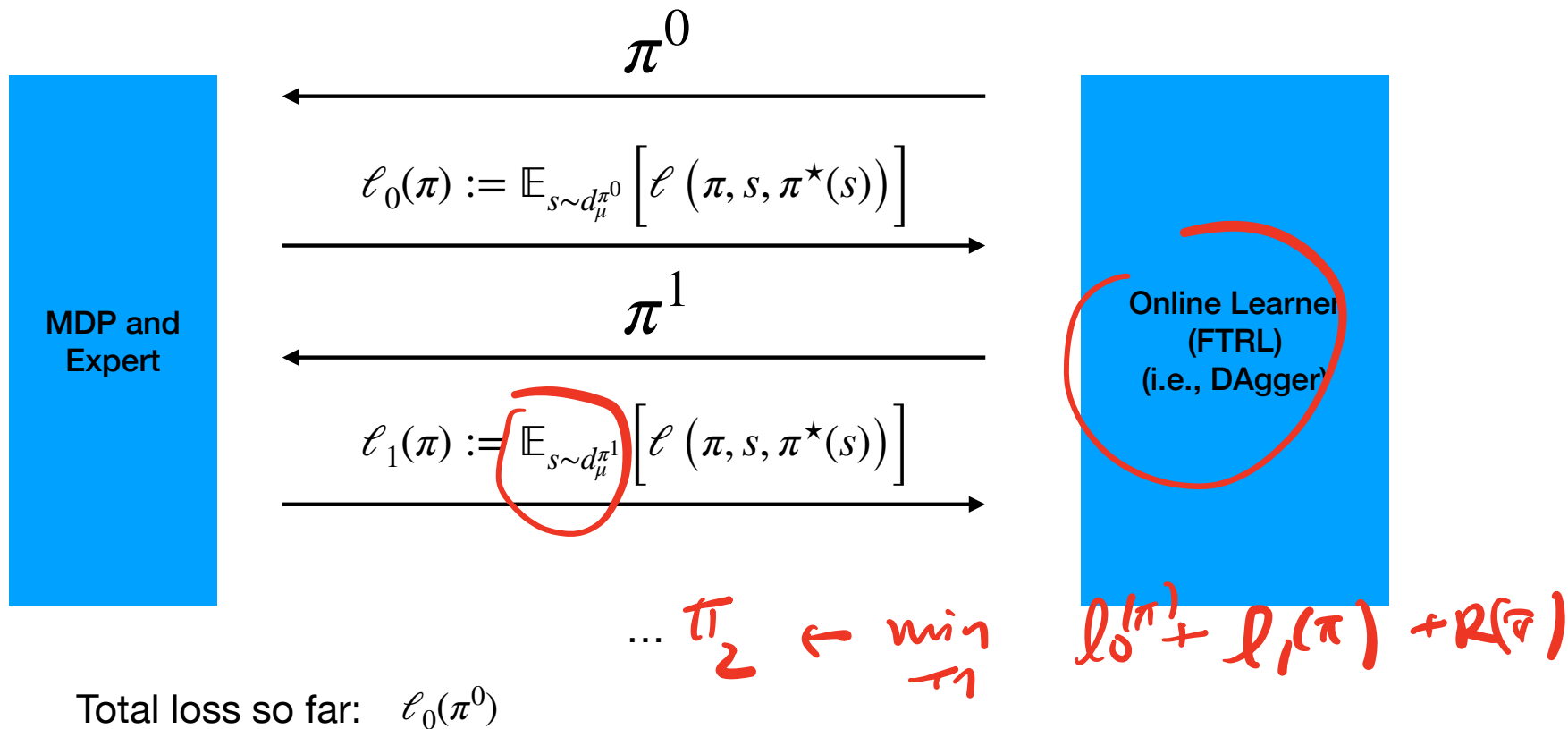
DAgger Analysis: A reduction to no-regret online learning



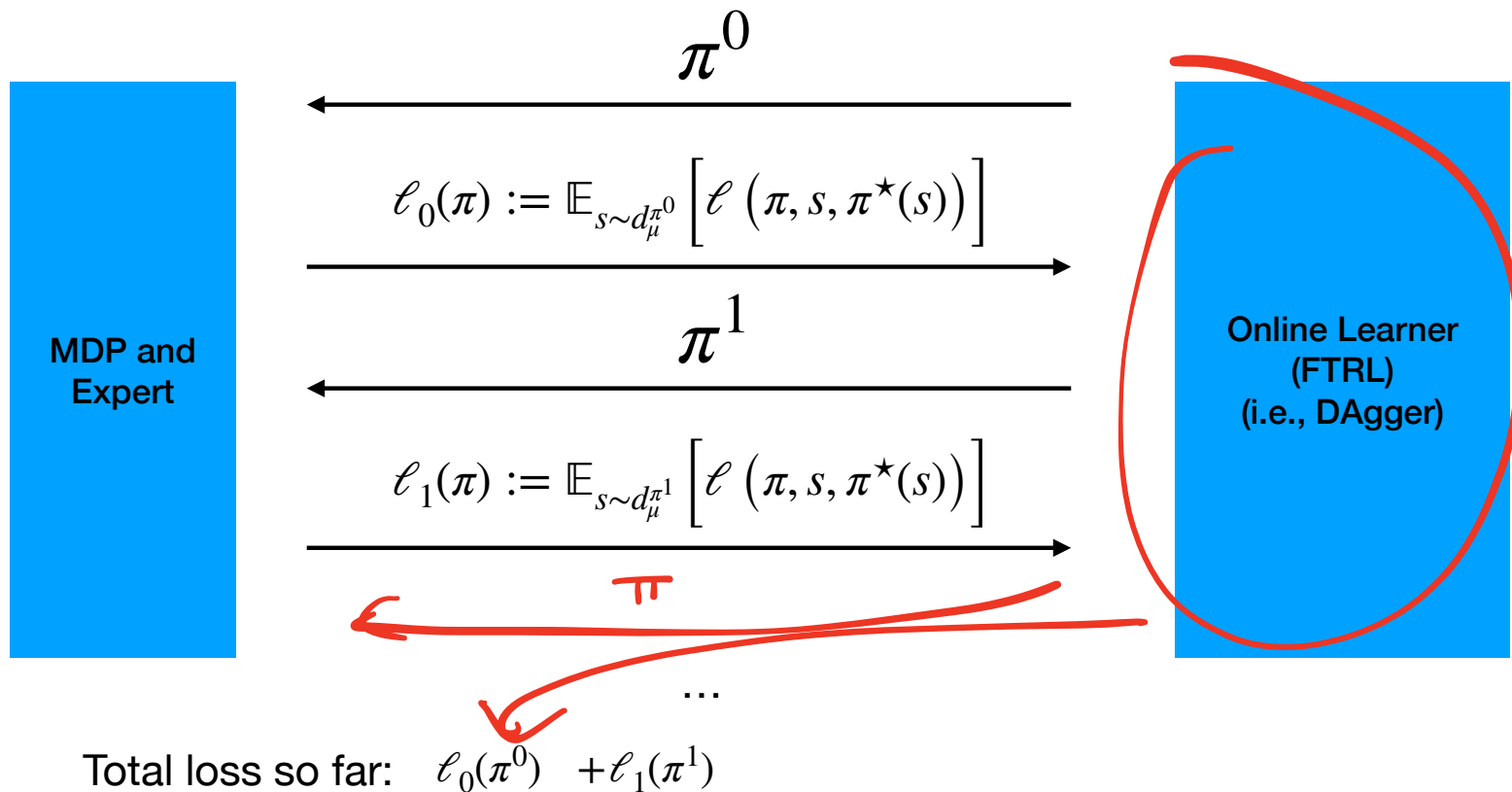
DAgger Analysis: A reduction to no-regret online learning



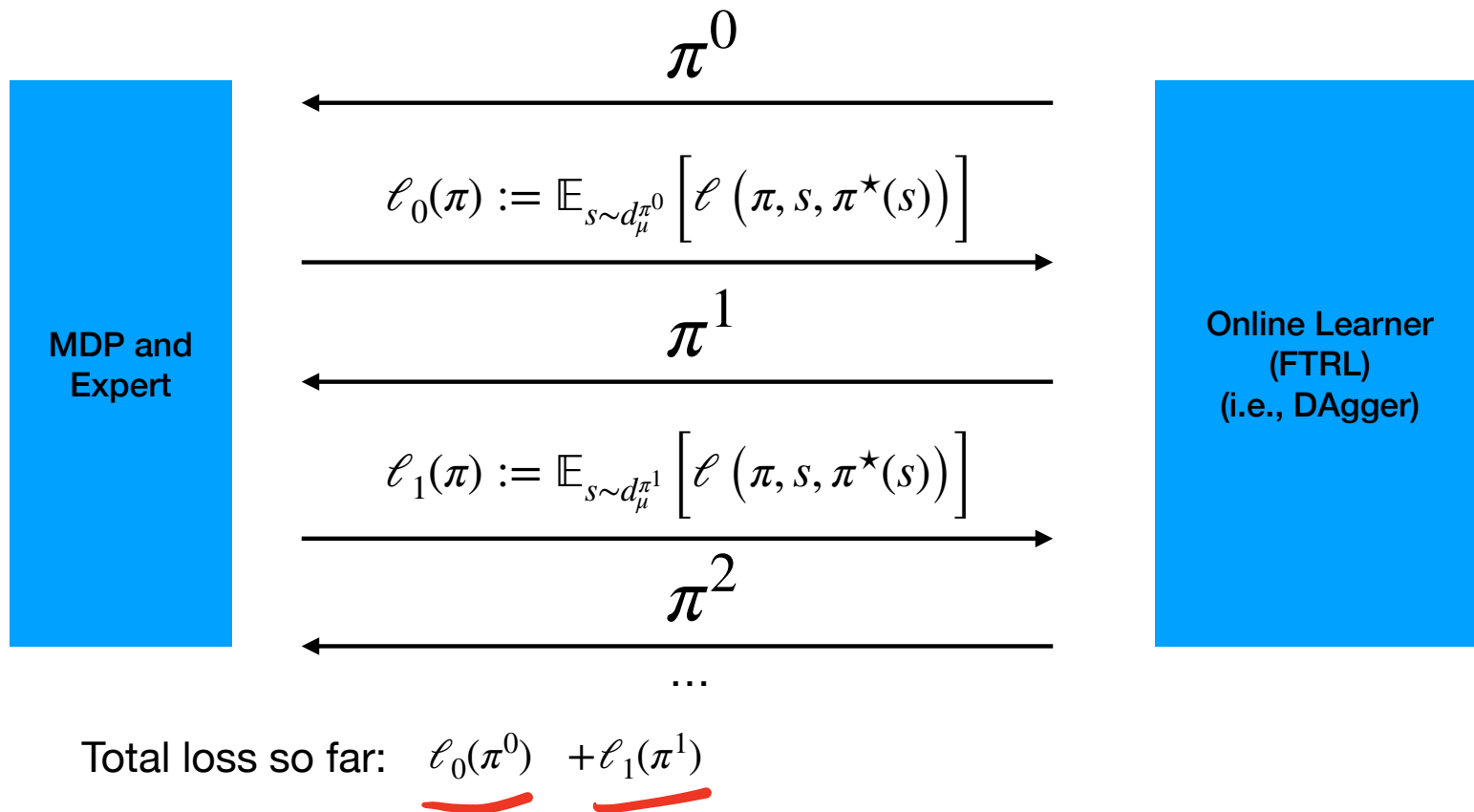
DAgger Analysis: A reduction to no-regret online learning



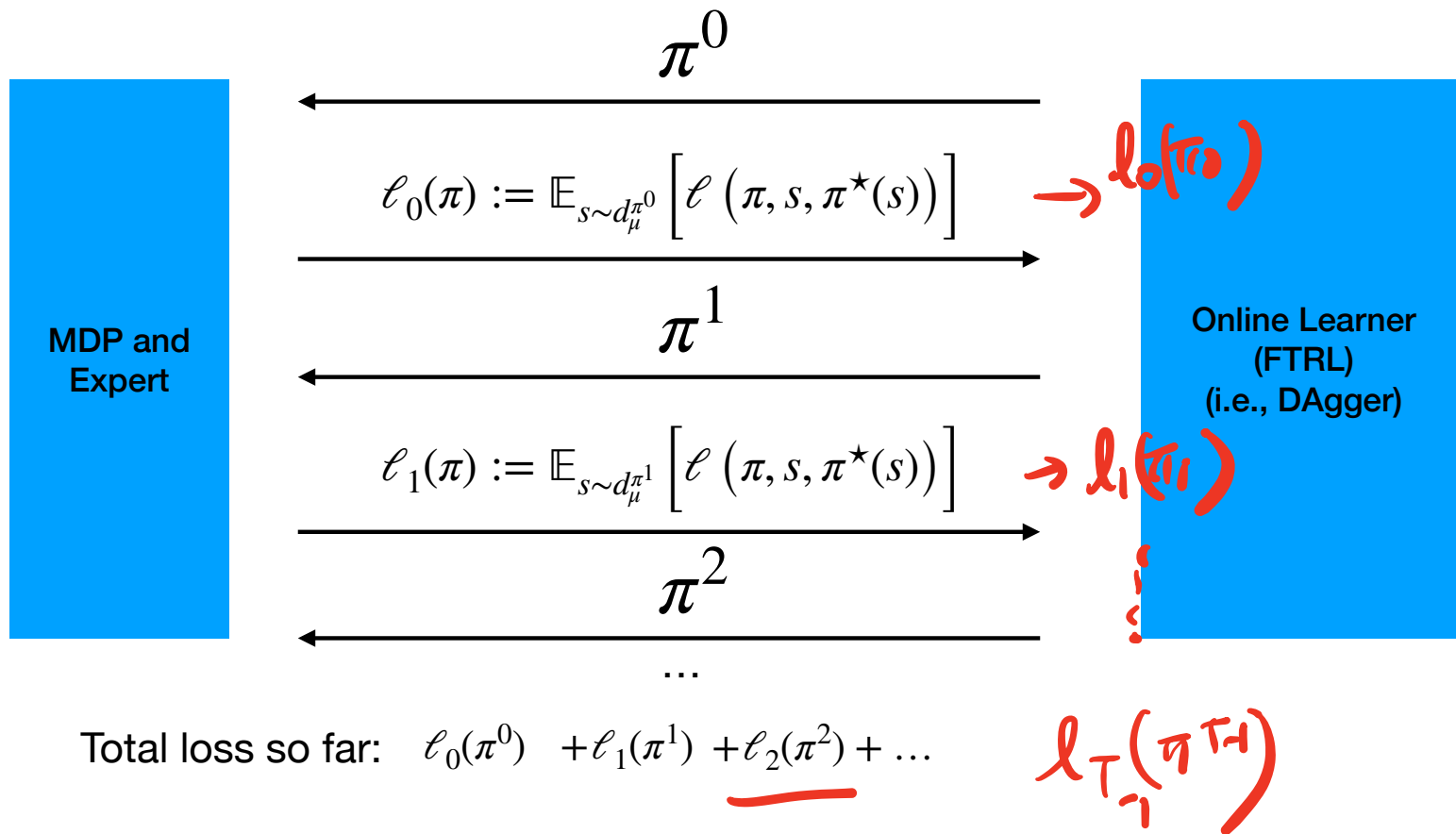
DAgger Analysis: A reduction to no-regret online learning



DAgger Analysis: A reduction to no-regret online learning



DAgger Analysis: A reduction to no-regret online learning



DAgger Analysis: A reduction to no-regret online learning

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

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

DAgger Analysis: A reduction to no-regret online learning

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

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

≈ 0

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

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

Δ

$$\ell(\pi(s), \pi^*(s))$$

$$\mathbb{1}(\pi(s) \neq \pi^*(s))$$

DAgger Analysis: A reduction to no-regret online learning

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

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

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

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

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 exist $t \in \{0, \dots, T-1\}$, such that:

$$\frac{1}{T} \sum_{t=1}^T \ell_x(\pi_t) \leq \epsilon_{\text{reg}} \rightarrow \underset{\Delta}{\ell_t(\pi^t)} \leq \epsilon_{\text{reg}}$$

DAgger Analysis: A reduction to no-regret online learning

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

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

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

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

DAgger Analysis: A reduction to no-regret online learning

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

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

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

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


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

Dagger Analysis: A reduction to no-regret online learning

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

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

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

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


π^t matches to π^* under its own state distribution!



Recall BC, we had:

$$\mathbb{E}_{s \sim d^{\pi^*}} \left[\ell(\hat{\pi}, s, \pi^*(s)) \right] \leq \epsilon, \text{ i.e., we matched to } \pi^* \text{ under } \pi^* \text{'s distribution}$$


Finally, turn things into the performance bound using PDL:

Theorem: There exists a iteration t , such that:

$$A^{\pi^t} = Q - V^{\pi^t}$$

$$\underline{V^{\pi^*} - V^{\pi^t} \leq \frac{\max_{s,a} |A^{\pi^*}(s,a)|}{(1-\gamma)} \epsilon_{reg}}$$

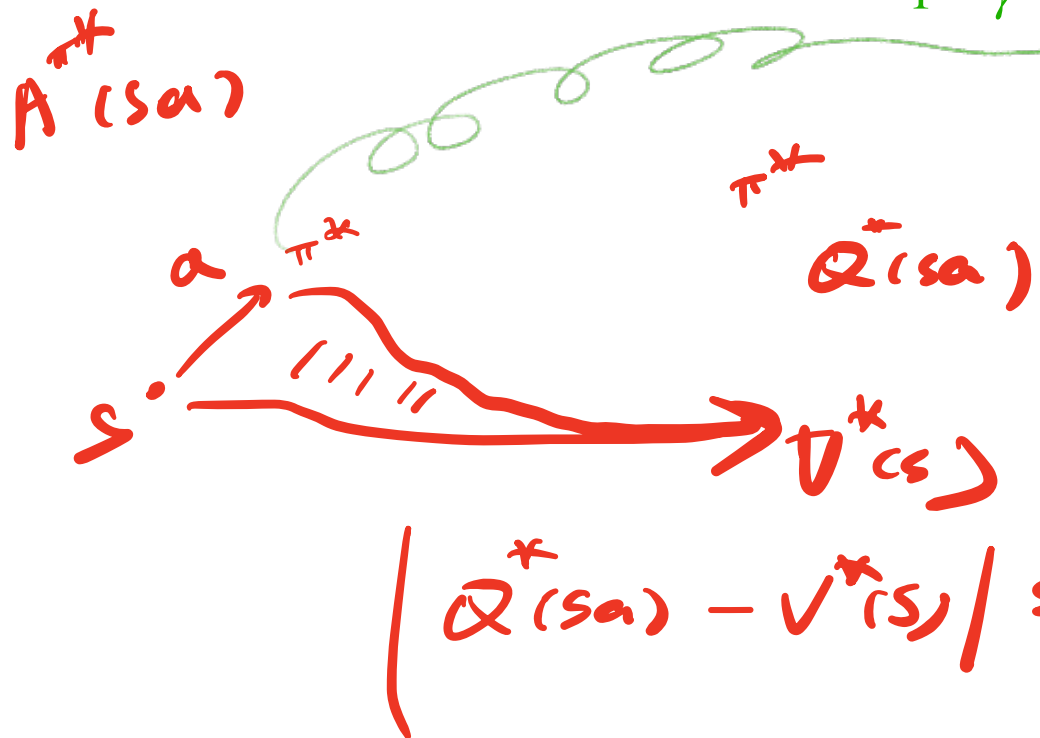
max absolute value
of Adv
of π^*

This bound indicates that:

Finally, turn things into the performance bound using PDL:

Theorem: There exists a iteration t , such that:

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



This bound indicates that:

We **avoid quadratic error** if expert π^* can quickly recover from a mistake

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

$$c \leq \frac{1}{1-\gamma}$$

$$|Q^*(sa) - V^*(s)| \leq c$$

Finally, turn things into the performance bound using PDL:

Theorem: There exists a iteration t , such that:

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

This bound indicates that:

We **avoid quadratic error** if expert π^* can quickly recover from a mistake

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

i.e., at any state, π^* can quickly recover from your mistake (take action a)

Finally, turn things into the performance bound using PDL:

Theorem: There exists a iteration t , such that:

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

PDL

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

$$A^{\pi^*}(s, \pi^*(s)) \neq 0$$

This bound indicates that:

We **avoid quadratic error** if expert π^* can quickly recover from a mistake

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

i.e., at any state, π^* can quickly recover from your mistake (take action a)

Finally, turn things into the performance bound using PDL:

Theorem: There exists a iteration t , such that:

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

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

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

// 0

Case 1, $\pi^t(s) = \pi^*(s)$

Case $\pi^t(s) \neq \pi^*(s)$

$$A^{\pi^*}(s, \pi^t(s)) - \underbrace{A^{\pi^*}(s, \pi^*(s))}_{=0} \geq -\max_{s,a} |A^{\pi^*}(s, a)|$$

This bound indicates that:

We **avoid quadratic error** if expert π^* can quickly recover from a mistake

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

i.e., at any state, π^* can quickly recover from your mistake (take action a)

Finally, turn things into the performance bound using PDL:

Theorem: There exists a iteration t , such that:

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

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

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

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

$$\mathbb{E}_{s \sim d^{\pi^t}} \mathbf{1}\{\pi^t(s) \neq \pi^*(s)\} \leq \epsilon_{reg}$$

This bound indicates that:

We **avoid quadratic error** if expert π^* can quickly recover from a mistake

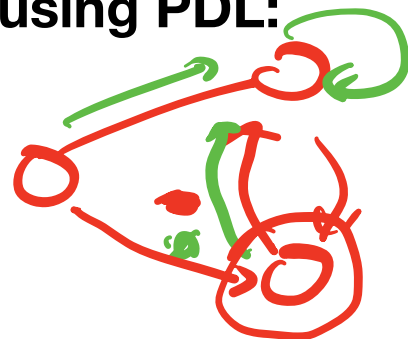
$$\max_{s,a} |A^{\pi^*}(s, a)| \leq c \in \mathbb{R}^+$$

i.e., at any state, π^* can quickly recover from your mistake (take action a)

Finally, turn things into the performance bound using PDL:

Theorem: There exists a iteration t , such that:

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



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

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

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

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

FTRL

This bound indicates that:

We **avoid quadratic error** if expert π^* can quickly recover from a mistake

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

i.e., at any state, π^* can quickly recover from your mistake (take action a)

Summary of DAgger

Summary of DAgger

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

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

Summary of DAgger

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

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

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

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

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

No expert interaction, no real world interaction

2. Interactive IL setting:

We have access to π^\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 $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d^{\pi^\star}$ (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]

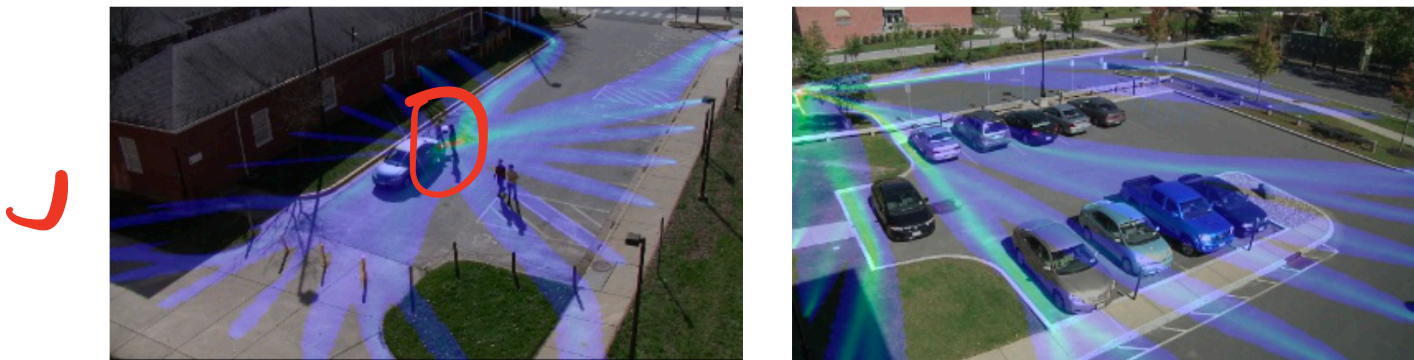


Fig. 1. Given a single pedestrian detection, our proposed approach forecasts plausible paths and destinations from noisy vision-input

Running Example: Human trajectory forecasting

[Kitani, et al, ECCV 12]

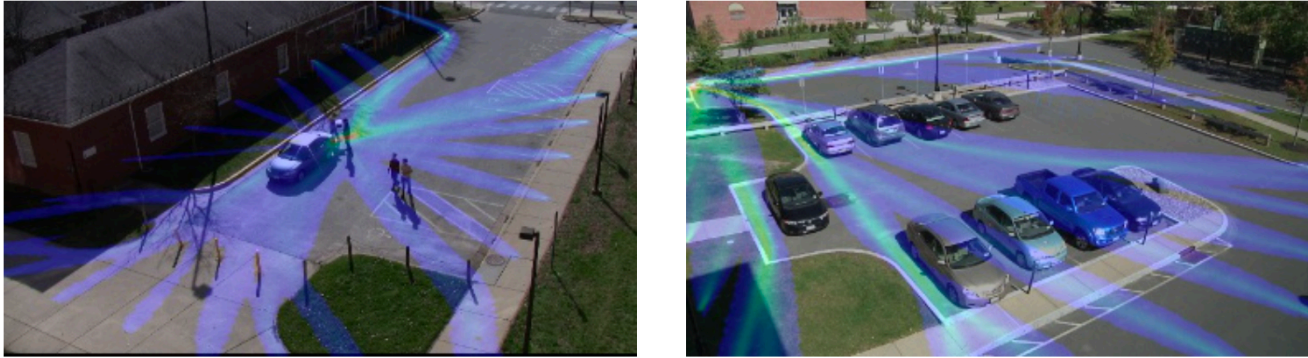


Fig. 1. Given a single pedestrian detection, our proposed approach forecasts plausible paths and destinations from noisy vision-input

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 π^\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 π^\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 π^\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:

$$c(s, a) = \langle \theta^*, \phi(s, a) \rangle, \text{ linear wrt feature } \phi(s, a)$$

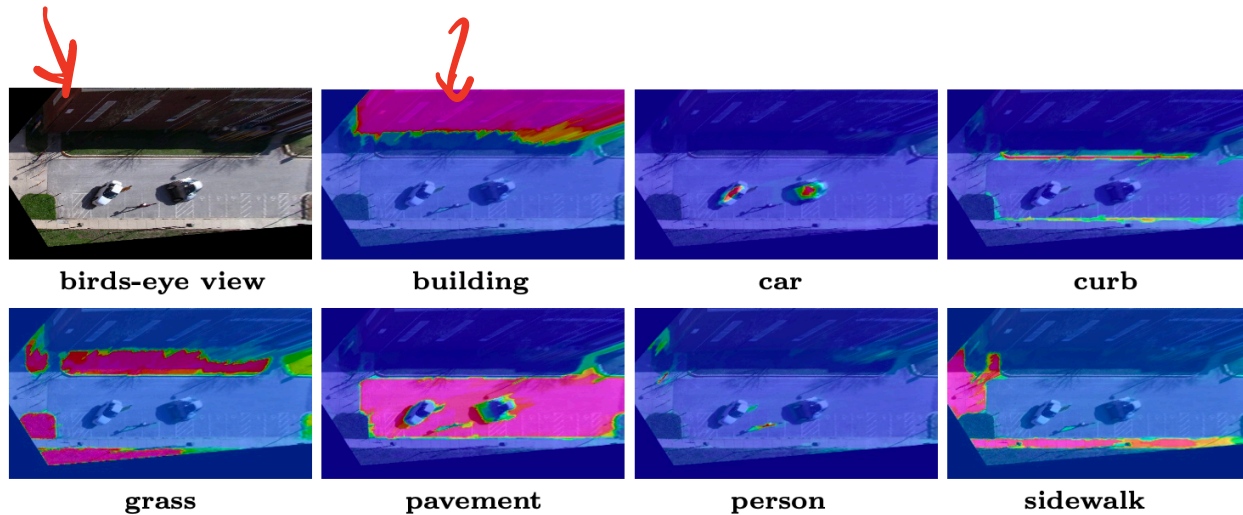


Fig. 4. Classifier feature response maps. Top left is the original image.

Running Example: Define feature map

Key Assumption on cost:

$$c(s, a) = \langle \theta^*, \phi(s, a) \rangle, \text{ linear wrt feature } \phi(s, a)$$

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

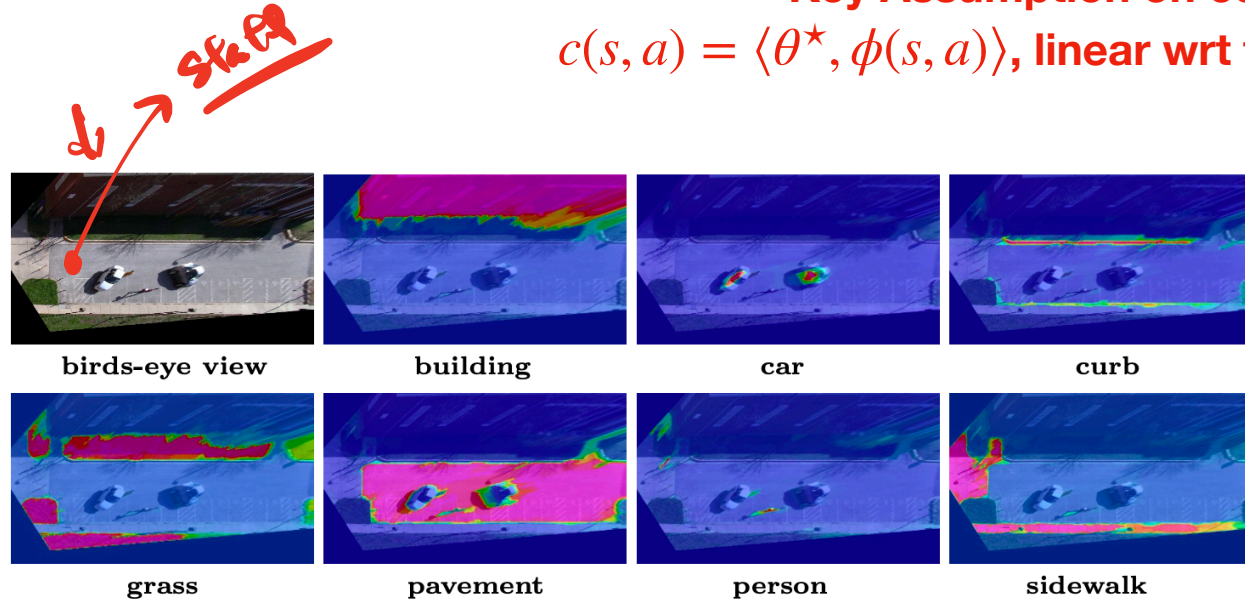


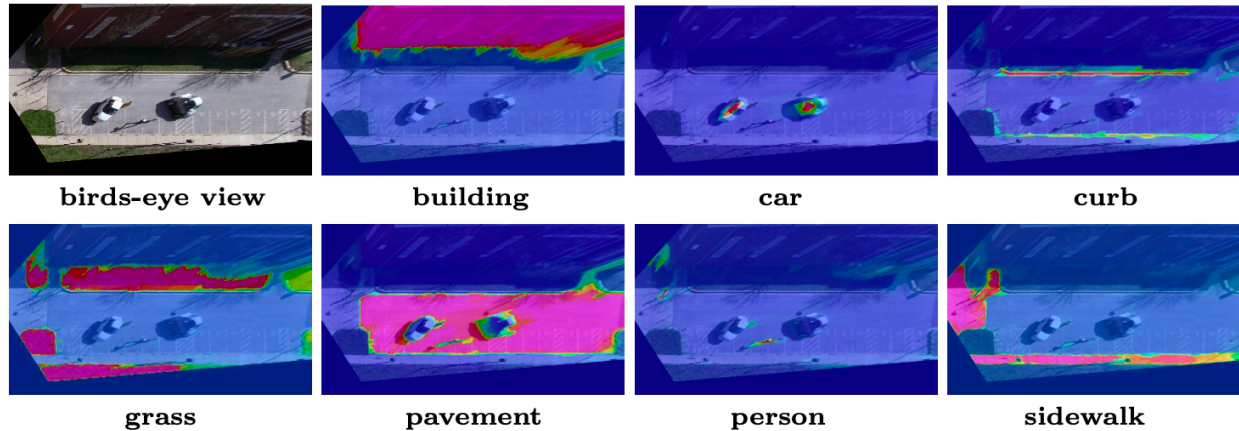
Fig. 4. Classifier feature response maps. Top left is the original image.

Running Example: Define feature map

Key Assumption on cost:

$$c(s, a) = \langle \theta^*, \phi(s, a) \rangle, \text{ linear wrt feature } \phi(s, a)$$

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



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

$$\theta^* \phi(s, a)$$

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

Running Example: Define feature map

Key Assumption on cost:

$$c(s, a) = \langle \theta^*, \phi(s, a) \rangle, \text{ linear wrt feature } \phi(s, a)$$

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

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

Maybe colliding with cars or buildings has **high** cost, but walking on sidewalk or grass has **low** cost

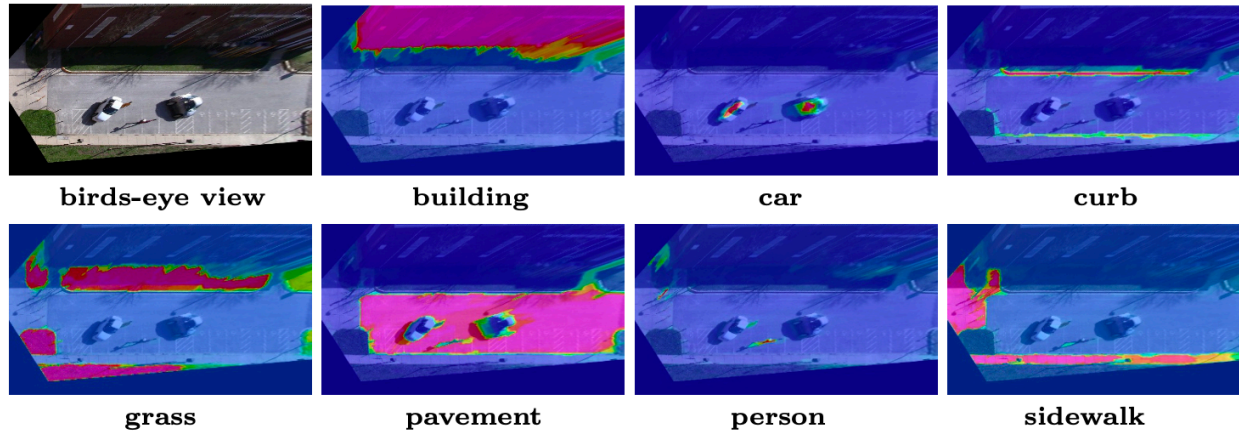
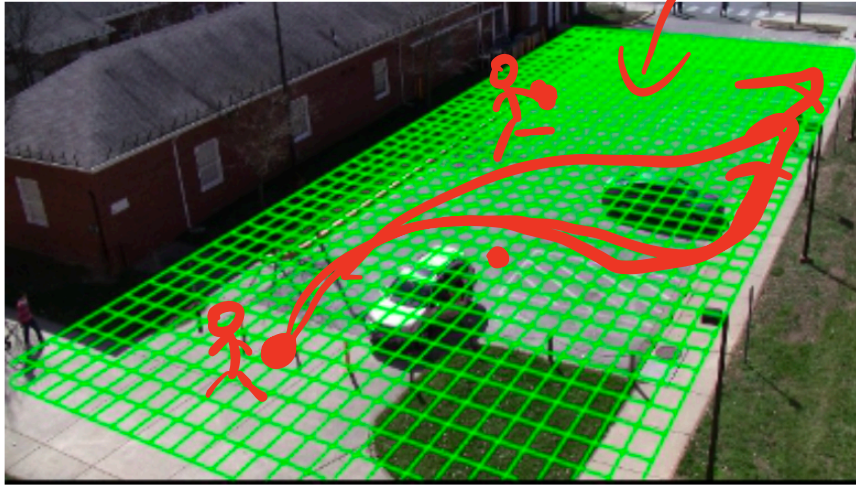


Fig. 4. Classifier feature response maps. Top left is the original image.

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