

Imitation Learning

ICML 2018 Tutorial
(Slides Available Online)

Yisong Yue



Hoang M. Le



yyue@caltech.edu

hmle@caltech.edu



@YisongYue

@HoangMinhLe



visongyue.com

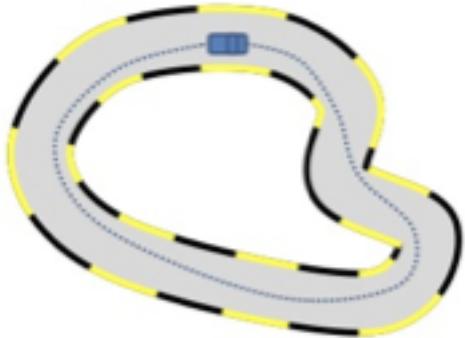
hoangle.info

Imitation Learning in a Nutshell

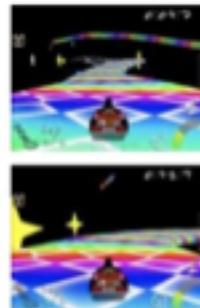
Given: demonstrations or demonstrator

Goal: train a policy to mimic demonstrations

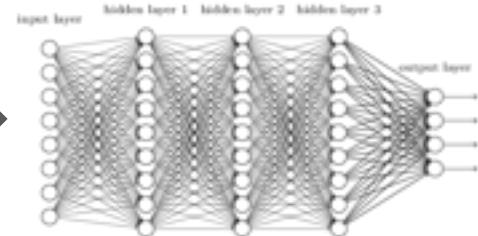
Expert Demonstrations



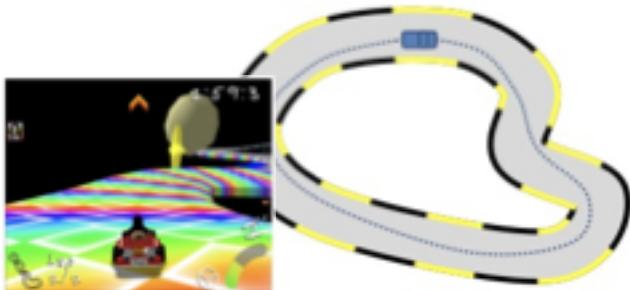
State/Action Pairs



Learning



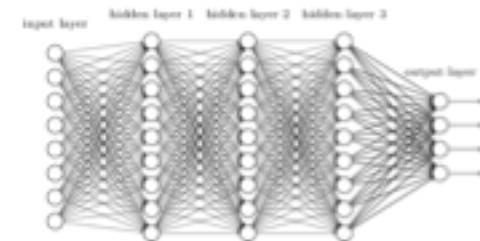
Ingredients of Imitation Learning



Demonstrations or Demonstrator



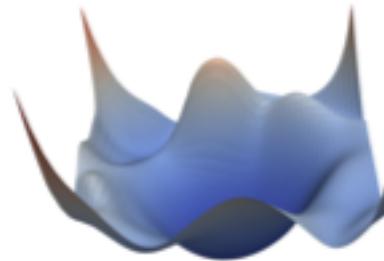
Environment / Simulator



Policy Class



Loss Function



Learning Algorithm

Tutorial Overview

Part 1: Introduction and Core Algorithms

- Teaser results
- Overview of ML landscape
- Types of Imitation learning
- Core algorithms of passive, interactive learning and cost learning

Part 2: Extensions and Applications

- Speech Animation
- Structured prediction and search
- Improving over expert
- Filtering and sequence modeling
- Multi-objective imitation learning
- Visual / few-shot imitation learning
- Domain Adaptation imitation learning
- Multi-agent imitation learning
- Hierarchical imitation learning
- Multi-modal imitation learning
- Weaker Feedback

ALVINN



Dean Pomerleau et al., 1989-1999

<https://www.youtube.com/watch?v=iIP4aPDTBPE>

Helicopter Acrobatics



Learning for Control from Multiple Demonstrations

Adam Coates, Pieter Abbeel, Andrew Ng, ICML 2008

An Application of Reinforcement Learning to Aerobatic Helicopter Flight

Pieter Abbeel, Adam Coates, Morgan Quigley, Andrew Y. Ng, NIPS 2006

<https://www.youtube.com/watch?v=0JL04Jjiocc>

Inferring Human Intent



Planning-based Prediction for Pedestrians

Brian Ziebart et al., IROS 2009

<https://www.youtube.com/watch?v=hiOteEd7gwE>

A Deep Learning Approach for Generalized Speech Animation

Sarah Taylor, Taehwan Kim, Yisong Yue et al., SIGGRAPH 2017

<https://www.youtube.com/watch?v=9zL7aeiW9fE>



Ghosting

(Sports Analytics)



Data Driven Ghosting using Deep Imitation Learning

Hoang M. Le et al., SSAC 2017

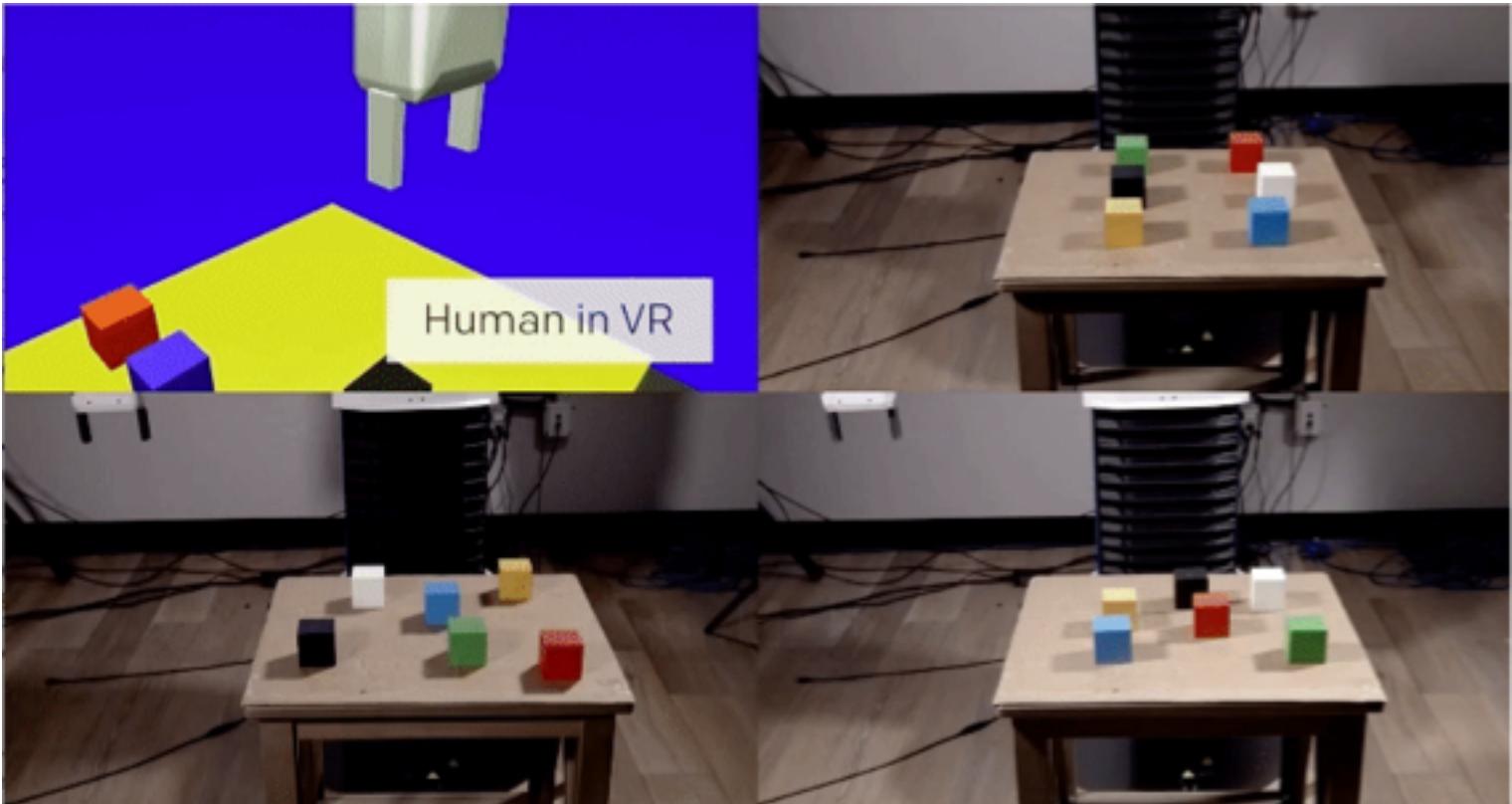
**English Premier League
2012-2013**

<https://www.youtube.com/watch?v=Wl-WL2ci0CA>

One Shot Imitation Learning

<https://www.youtube.com/watch?v=oM ZwklizzCM>

Duan et al., NIPS '17



Notation & Setup

State: s (sometimes x) (**state may only be partially observed)

Action: a (sometimes y)

Policy: π_θ (sometimes h)

- Policy maps states to actions: $\pi_\theta(s) \rightarrow a$
- ...or distributions over actions: $\pi_\theta(s) \rightarrow P(a)$

State Dynamics: $P(s'|s,a)$

- Typically not known to policy.
- Essentially the simulator/environment

Notation & Setup Continued

Rollout: sequentially execute $\pi(s_0)$ on an initial state

- Produce trajectory $\tau = (s_0, a_0, s_1, a_1, \dots)$

$P(\tau|\pi)$: distribution of trajectories induced by a policy

1. Sample s_0 from P_0 (distribution over initial states), initialize $t = 1$.
2. Sample action a_t from $\pi(s_{t-1})$
3. Sample next state s_t from applying a_t to s_{t-1} (requires access to environment)
4. Repeat from Step 2 with $t=t+1$

$P(s|\pi)$: distribution of states induced by a policy

- Let $P_t(s|\pi)$ denote distribution over t -th state
- $P(s|\pi) = (1/T)\sum_t P_t(s|\pi)$

Example #1: Racing Game

(Super Tux Kart)

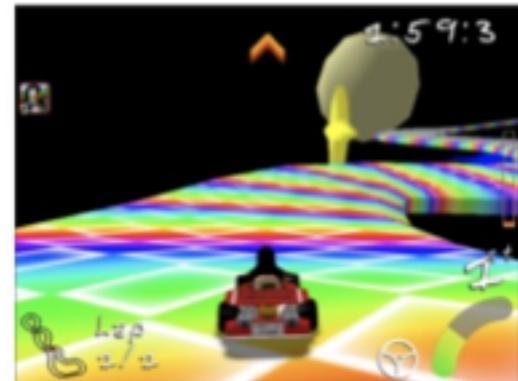
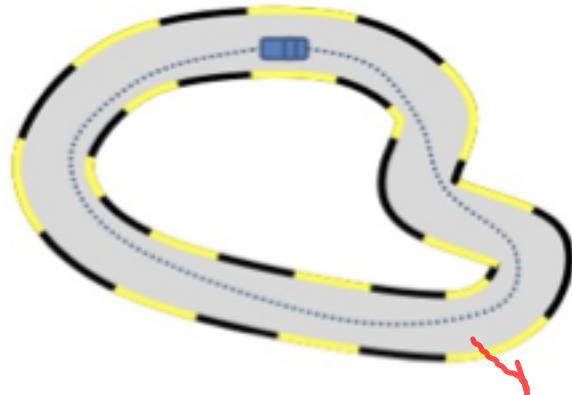
s = game screen

a = turning angle

Training set: $D=\{\tau:=(s,a)\}$ from π^*

- s = sequence of s
- a = sequence of a

Goal: learn $\pi_\theta(s) \rightarrow a$



Images from Stephane Ross

Example #2: Basketball Trajectories

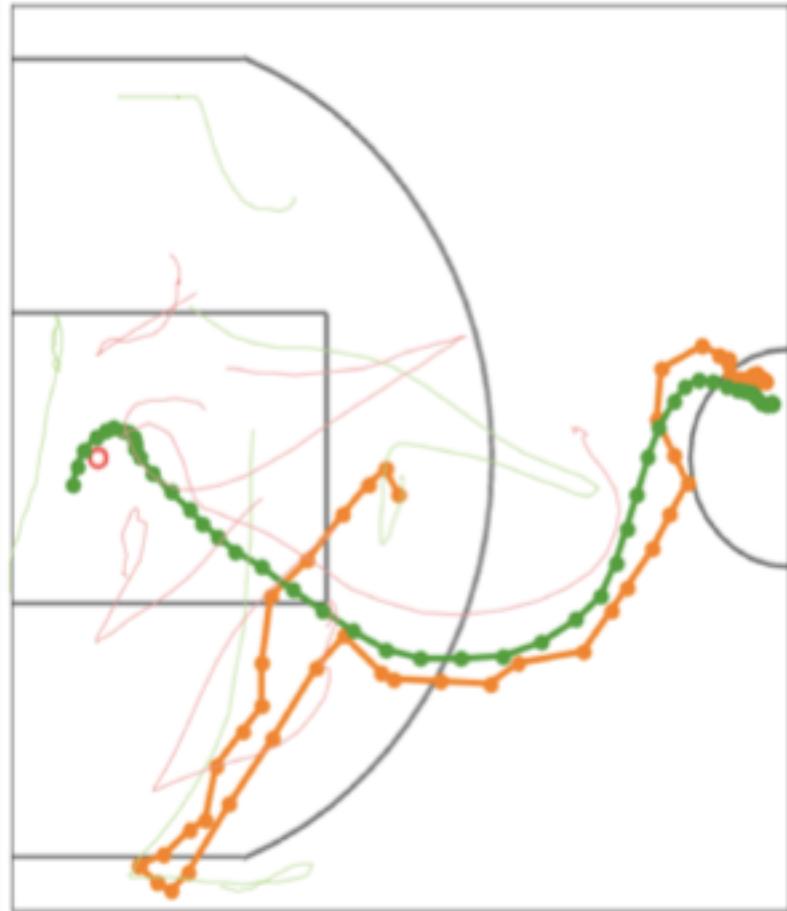
s = location of players & ball

a = next location of player

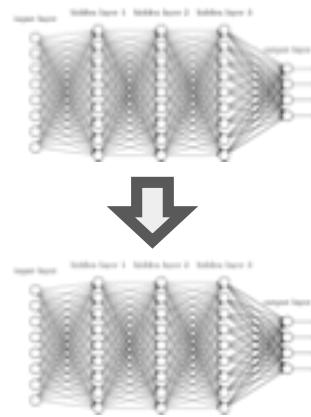
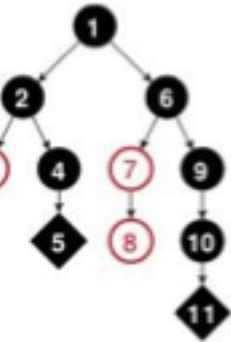
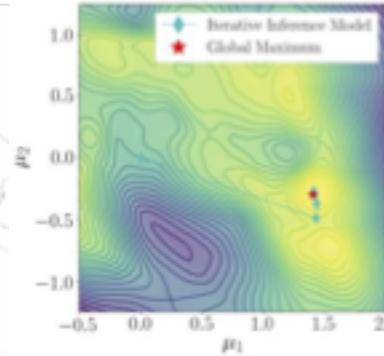
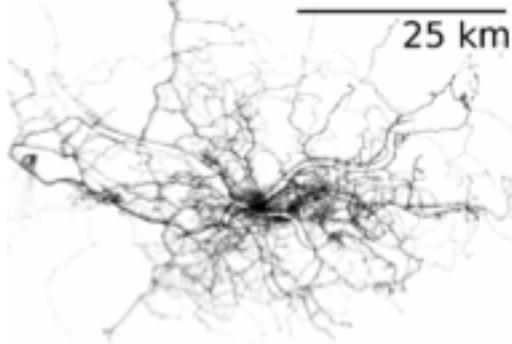
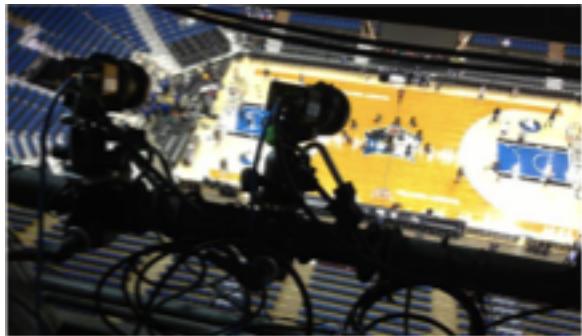
Training set: $D=\{\tau:=(s,a)\}$ from π^*

- s = sequence of s
- a = sequence of a

Goal: learn $\pi_\theta(s) \rightarrow a$



Data Sources



Outline of 1st Half

Behavioral Cloning (simplest Imitation Learning setting)

Compare with Supervised Learning

Landscape of Imitation Learning settings

Behavioral Cloning = Reduction to Supervised Learning

(Ignoring regularization for brevity.)

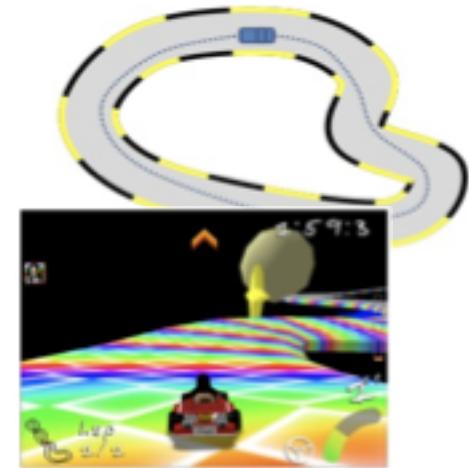
Define $P^* = P(s|\pi^*)$ (distribution of states visited by expert)

(recall $P(s|\pi^*) = (1/T)\sum_t P_t(s|\pi^*)$)

(sometimes abuse notation: $P^* = P(s, a^* = \pi^*(s)|\pi^*)$)

Learning objective:

$$\operatorname{argmin}_{\theta} E_{(s, a^*) \sim P^*} L(a^*, \pi_{\theta}(s))$$



Interpretations:

1. Assuming perfect imitation so far, learn to continue imitating perfectly
2. Minimize 1-step deviation error along the expert trajectories

(General) Imitation Learning vs Behavioral Cloning

(Ignoring regularization for brevity.)

Behavioral Cloning (Supervised Learning):

$$\operatorname{argmin}_{\theta} E_{(s,a^*) \sim P^*} L(a^*, \pi_{\theta}(s))$$

Distribution provided exogenously

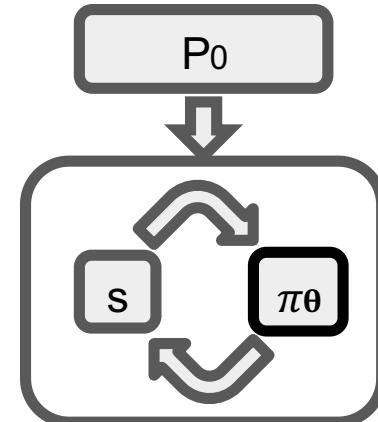
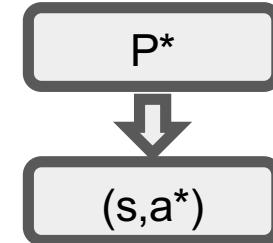
Training Loss

(General) Imitation Learning:

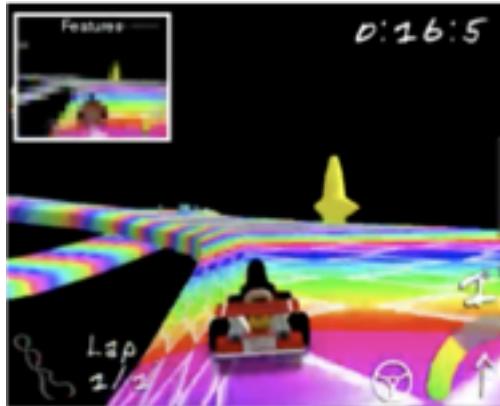
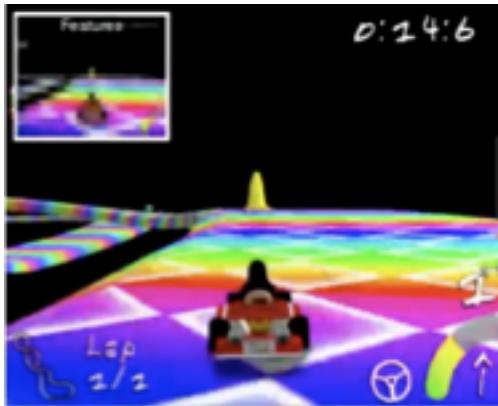
$$\operatorname{argmin}_{\theta} E_{s \sim P(s|\theta)} L(\pi^*(s), \pi_{\theta}(s))$$

Distribution depends on rollout.

$P(s|\theta)$ = state distribution of π_{θ}



Limitations of Behavioral Cloning

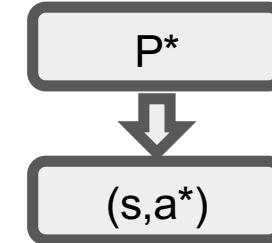


π_θ makes a mistake

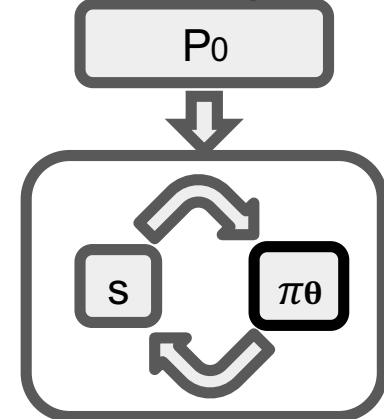
New state sampled not from P^* !

Worst case is catastrophic!

IID Assumption
(Supervised Learning)



Reality



Limitations of Behavioral Cloning



Expert Trajectories
(Training Distribution)



Behavioral Cloning
Makes mistakes, enters new states
Cannot recover from new states

When to use Behavioral Cloning?

Advantages

- Simple
- Simple
- Efficient

Disadvantages

- Distribution mismatch between training and testing
- No long term planning

Use When:

- 1-step deviations not too bad
- Learning reactive behaviors
- Expert trajectories “cover” state space

Don't Use When:

- 1-step deviations can lead to catastrophic error
- Optimizing long-term objective (at least not without a stronger model)

Outline of 1st Half

Behavioral Cloning (simplest Imitation Learning setting)

Compare with Supervised Learning

Landscape of Imitation Learning settings

Types of Imitation Learning

Behavioral Cloning

$$\operatorname{argmin}_{\theta} E_{(s,a^*) \sim P^*} L(a^*, \pi_{\theta}(s))$$

Works well when P^* close to P_{θ}

Inverse RL

Learn r such that:

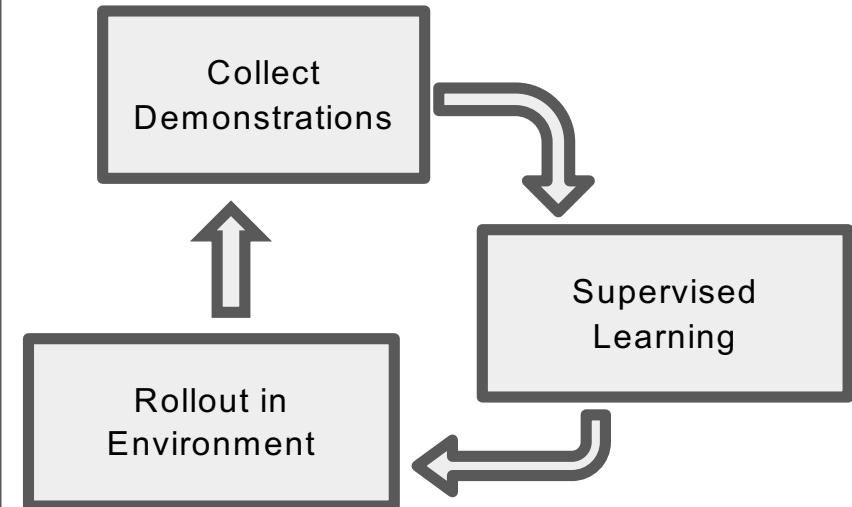
$$\pi^* = \operatorname{argmax}_{\theta} E_{s \sim P(s|\theta)} r(s, \pi_{\theta}(s))$$

RL problem

Assumes learning r is statistically easier than directly learning π^*

Direct Policy Learning

via Interactive Demonstrator



Requires Interactive Demonstrator
(BC is 1-step special case)

Types of Imitation Learning

	Direct Policy Learning	Reward Learning	Access to Environment	Interactive Demonstrator	Pre-collected Demonstrations
Behavioral Cloning	Yes	No	No	No	Yes
Direct Policy Learning (Interactive IL)	Yes	No	Yes	Yes	Optional
Inverse Reinforcement Learning	No	Yes	Yes	No	Yes

Other Considerations

other

Choice of imitation loss (e.g., generative adversarial learning)

Suboptimal demonstrations (e.g., outperform teacher)

Partial demonstrations (e.g., weak feedback)

Domain transfer (e.g., few-shot learning)

Structured domains (e.g., multi-agent systems, structured prediction)

Interactive Direct Policy Learning

Behavioral Cloning is simplest example

Beyond BC: using interactive demonstrator

Often analyzed via **learning reductions**

- Reduce “harder” learning problem to “easier” one
- Imitation Learning → Supervised Learning
- General Overview: <http://hunch.net/~il/projects/reductions/reductions.html>

Learning Reductions

Behavioral Cloning:

$$\mathbb{E}_{s \sim P(s|\theta)} L(a^*(s), \pi_\theta(s)) \rightarrow \mathbb{E}_{(s,a^*) \sim P^*} L(a^*, \pi_\theta(s))$$

A

B

What does learning well on B imply about A?

- E.g., can one lift PAC learning results from B to A?

Basic Results for Behavioral Cloning

Assume $L(a^*, a^*) = 0$

Suppose: $\varepsilon = E_{(s,a^*) \sim P^*} L(a^*, \pi_\theta(s))$

Then: $E_{s \sim P(s|\theta)} L(a^*(s), \pi_\theta(s)) = O(T\varepsilon)$

- Error on T -step trajectory is $O(T^2\varepsilon)$

*Paraphrased from

Theorem 2.1 in [Efficient Reductions for Imitation Learning - Ross & Bagnell, AISTATS 2010]

Lemma 3 in [A reduction from apprenticeship learning to classification - Syed & Schapire, NIPS 2010]

Direct Policy Learning vs Interactive Expert

Sequential Learning Reductions

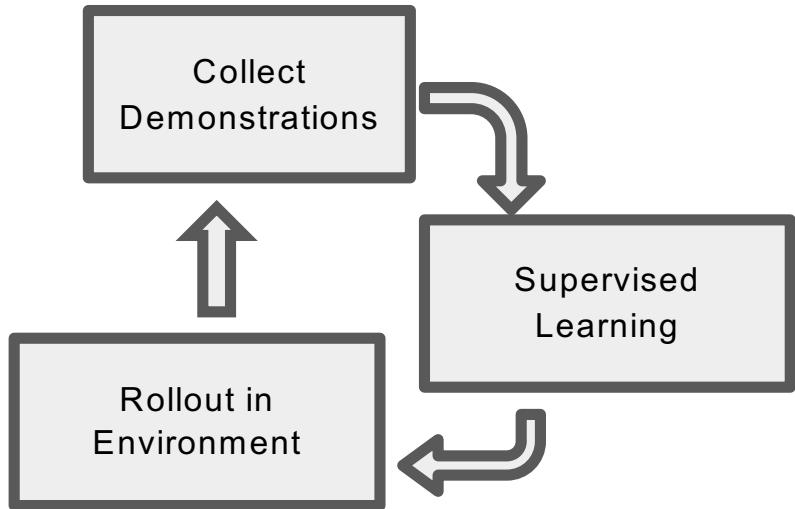
Sequence of distributions:

- $E_{s \sim P(m)} L(\pi^*(s), \pi(s))$
- Ideally converges to π^{OPT}

Best in policy class

Usually starting from:

- $E_{(s) \sim P^*} L(\pi^*(s), \pi(s))$



Requires Interactive Demonstrator
(BC is 1-step special case)

Interactive Expert

Can query expert at any state

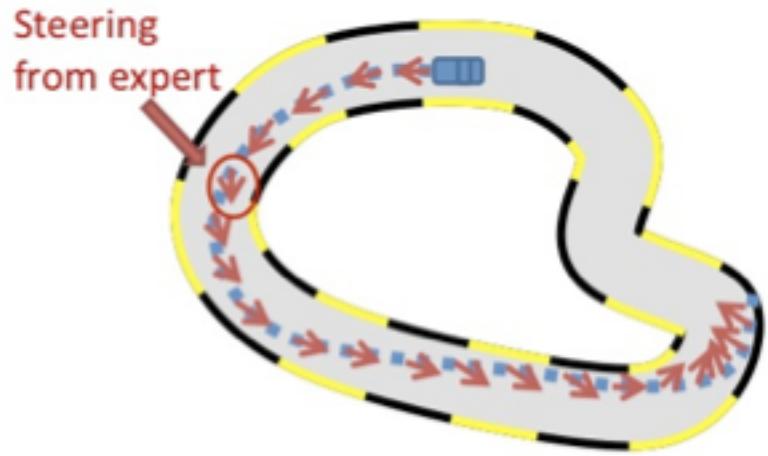
Construct loss function

- $L(\pi^*(s), \pi(s))$

Typically applied to rollout trajectories

- $s \sim P(s|\pi)$

Driving example: $L(\pi^*(s), \pi(s)) = (\pi^*(s) - \pi(s))^2$

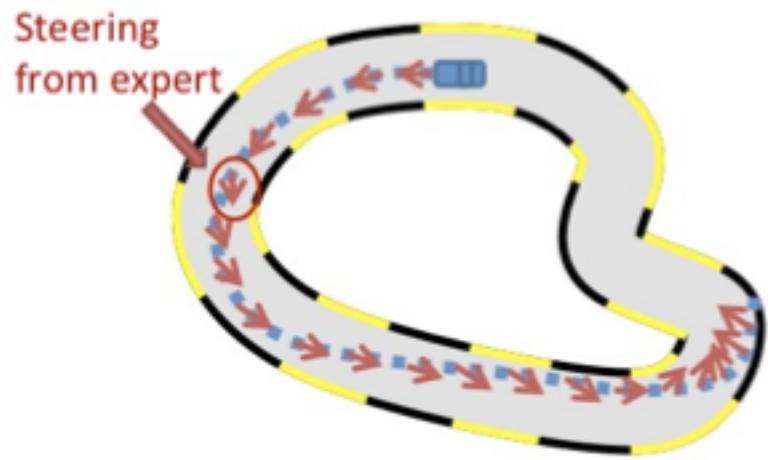


Example from Super Tux Kart
(Image courtesy of Stephane Ross)

Expert provides feedback on
state visited by policy

Alternating Optimization (Naïve Attempt)

1. Fix P , estimate π
 - Solve $\text{argmin}_{\theta} E_{s \sim PL}(\pi(s), \pi_\theta(s))$
2. Fix π , estimate P
 - Empirically estimate via rolling out π
3. Repeat



Not guaranteed to converge!

Sequential Learning Reductions

Initial predictor: π_0 ← **Initial expert demonstrations**

For $m=1$

- Collect trajectories τ via rolling out π_{m-1} ← **Typically roll out multiple times**
- Estimate state distribution P_m using $s \in \tau$
- Collect interactive feedback $\{\pi^*(s) | s \in \tau\}$
- **Data Aggregation** (e.g., DAgger)
 - Train π_m on $P_1 \cup \dots \cup P_m$
- **Policy Aggregation** (e.g., SEARN & SMILe)
 - Train π'_m on P_m
 - $\pi_m = \beta \pi'_m + (1-\beta) \pi_{m-1}$

← **Requires interactive expert**

Data Aggregation (DAgger)

Sequence of convex losses: $L_m(\pi) = E_{s \sim P(m)} L(\pi^*(s), \pi(s))$

Online Learning: find sequence π_m competitive with π^{OPT} :

$$R_M = \left(\frac{1}{M}\right) \sum_m L_m(\pi_m) - \left(\frac{1}{M}\right) \sum_m L_m(\pi^{OPT})$$

“Online Regret”



Follow-the-Leader: $\pi_m = \min_\theta \sum_{m'=1}^m L_{m'}(\pi_\theta)$

$$R_M = O(1/\sqrt{M}) \quad (M > T)$$

$$\exists \pi_m: \left(\frac{1}{M}\right) \sum_{m'} L_{m'}(\pi^{OPT}) \leq \left(\frac{1}{M}\right) \sum_{m'} L_{m'}(\pi_m) - O\left(\frac{1}{\sqrt{M}}\right)$$

Typically π_M

A reduction of imitation learning and structured prediction to no-regret online learning

Stephane Ross, Geoff Gordon, Drew Bagnell, AISTATS 2011

Policy Aggregation (SEARN & SMILE)

Train π'^m on $P_m \rightarrow \pi_m = \beta\pi'^m + (1-\beta)\pi_{m-1}$

$$\pi_m = (1-\beta)^M \pi_0 + \beta \sum_{m'} (1 - \beta)^{M-m'} \pi'^{m'}$$

π_0 is expert
(not available at test time)

π_m not much worse than π_{m-1}

- At most $\beta TL_m(\pi'^m) + \beta^2 T^2 / 2$ for T-step rollout

π_M not much worse than π_0

π_0 negligible in π_M

- At most $2T \log(T) (\frac{1}{M}) \sum_m L_m(\pi'^m) + O(1/T)$ for T-step rollout

Direct Policy Learning via Interactive Expert

Reduction to sequence of supervised learning problems

- Constructed from roll-outs of previous policies
- Requires interactive expert feedback

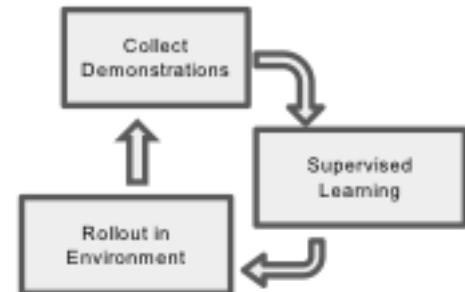
Two approaches: Data Aggregation & Policy Aggregation

- Ensures convergence
- Motivated by different theory

Not covered:

Depends on application

- What is expert feedback & loss function?



Types of Imitation Learning

Behavioral Cloning

$$\operatorname{argmin}_{\theta} E_{(s,a^*) \sim P^*} L(a^*, \pi_{\theta}(s))$$

Works well when P^* close to P_{θ}

Inverse RL

Learn r such that:

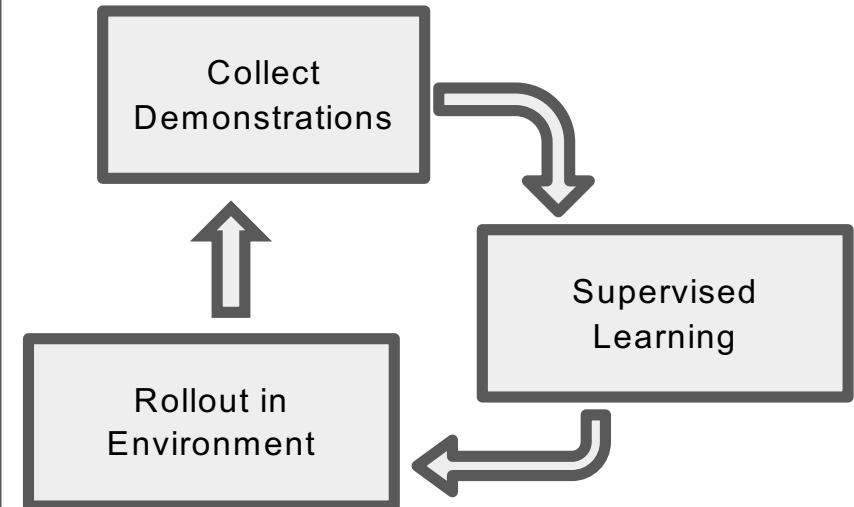
$$\pi^* = \operatorname{argmax}_{\theta} E_{s \sim P(s|\theta)} r(s, \pi_{\theta}(s))$$

RL problem

Assumes learning r is statistically easier than directly learning π^*

Direct Policy Learning

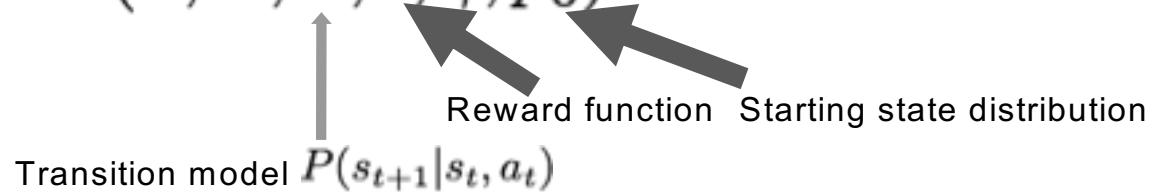
via Interactive Demonstrator



Requires Interactive Demonstrator
(BC is 1-step special case)

Learn Policy w/o Expert: Reinforcement Learning

MDP Formulation: $(S, A, P, r, \gamma, p_0)$



Goal: maximize cumulative rewards

$$\max_{\pi \in \Pi} V(\pi) \triangleq \max_{\pi \in \Pi} \mathbb{E}_{\pi} [r(s, a)] = \max_{\pi \in \Pi} \mathbb{E}_{s_0 \sim p_0} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| \pi \right]$$

$$Q^{\pi}(s, a) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| s_0 = s, a_0 = a \right]$$

Fully specified MDP: value & policy iteration

Learn Policy w/o Expert: Reinforcement Learning

MDP Formulation: $(S, A, \cancel{P}, r, \gamma, p_0)$

Optimize Value Function wrt Parameterized Policy π_θ

$$V(\pi_\theta) = \mathbb{E}_{s_0 \sim p_0} \left[\text{cumulative rewards} \mid \pi_\theta \right]$$

$$Q^{\pi_\theta}(s, a) = \mathbb{E} [\text{cumulative rewards} \mid \pi_\theta, s_0 = s, a_0 = a]$$

Policy Gradient Theorem (Sutton et al., ICML 1999)

$$\nabla_\theta V(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q^{\pi_\theta}(s, a)]$$

Learn Policy w/o Expert: Reinforcement Learning

MDP Formulation: $(S, A, \cancel{R}, r, \gamma, p_0)$

REINFORCE algorithm:

for each trajectory $\tau = \{s_0, a_0, s_1, a_1, \dots, s_T, a_T\} \sim \pi_\theta$:
for $t = 0, \dots, T - 1$:

$$\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(s_t, a_t) \hat{Q}^{\pi_\theta}(s_t, a_t)$$

Also check out...

- **Introduction to Reinforcement Learning with Function Approximation** – Richard Sutton – NIPS 2015 Tutorial
- **Deep Reinforcement Learning** – David Silver – ICML 2016 Tutorial
- **Deep Reinforcement Learning, Decision Making, and Control** – Sergey Levine & Chelsea Finn – ICML 2017 Tutorial
- **Deep Reinforcement Learning through Policy Optimization** – Pieter Abbeel & John Schulman – NIPS 2016 Tutorial

Challenges with Reward Engineering

MDP Formulation: $(S, A, P, r, \gamma, p_0)$  assumed given



Inverse Reinforcement Learning

MDP Formulation: (S, A, P, γ, p_0)

Given: $\mathcal{D} = \{\tau_1, \dots, \tau_m\} = \{(s_0^i, a_0^i, s_1^i, a_1^i \dots)\} \sim \pi^*$

Goal: Learn a reward function r^* so that

$$\pi^* = \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_\pi [r^*(s, a)]$$

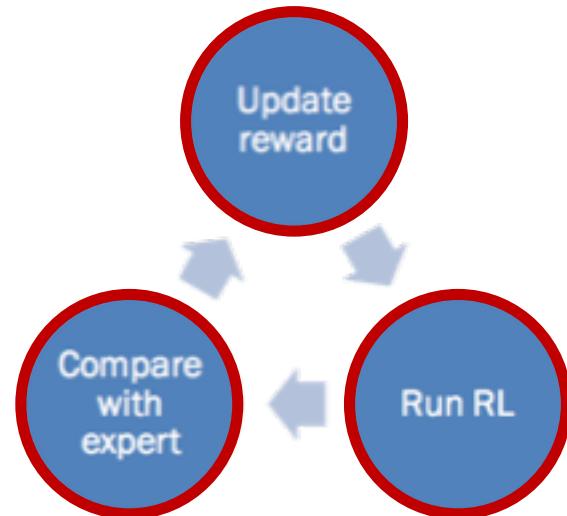


can also be $\mathbb{E}_\pi [r^*(s)]$

Inverse Reinforcement Learning

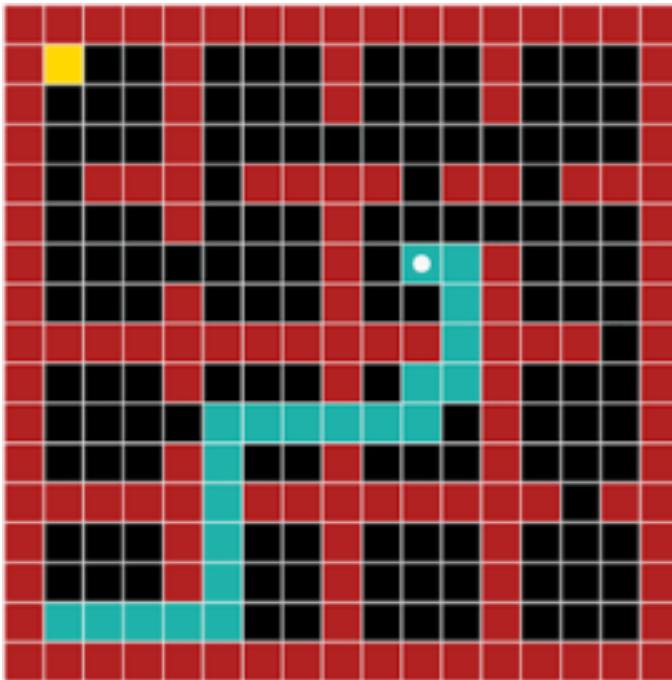
■ Inverse RL high-level recipe:

- Expert demonstrations: $\mathcal{D} = \{\tau_1, \dots, \tau_m\}$
- Learn reward function: $r_\theta(s, a) = \theta^\top \phi(s, a)$
- Learn policy given reward function (RL)
- Compare learned policy with expert



Reward Learning is Ambiguous

1. Many reward functions correspond to the same policy



Imitation Learning via Inverse RL (model-given)

Abbeel & Ng, ICML '04

Goal: find reward function r

$$\max_{\pi \in \Pi} \mathbb{E}_{\pi} [r(s, a)] \geq \mathbb{E}_{\pi^*} [r^*(s, a)] - \epsilon$$



Different from “idealized” IRL

Game-Theoretic Inverse RL (model-given)

Syed & Schapire, NIPS '07

Goal: find π performing better than π^* over a class of rewards

$$\max_{\pi \in \Pi} \min_{r \in \mathcal{R}} \mathbb{E}_{\pi} [r(s, a)] - \mathbb{E}_{\pi^*} [r(s, a)]$$



Different from “idealized” IRL

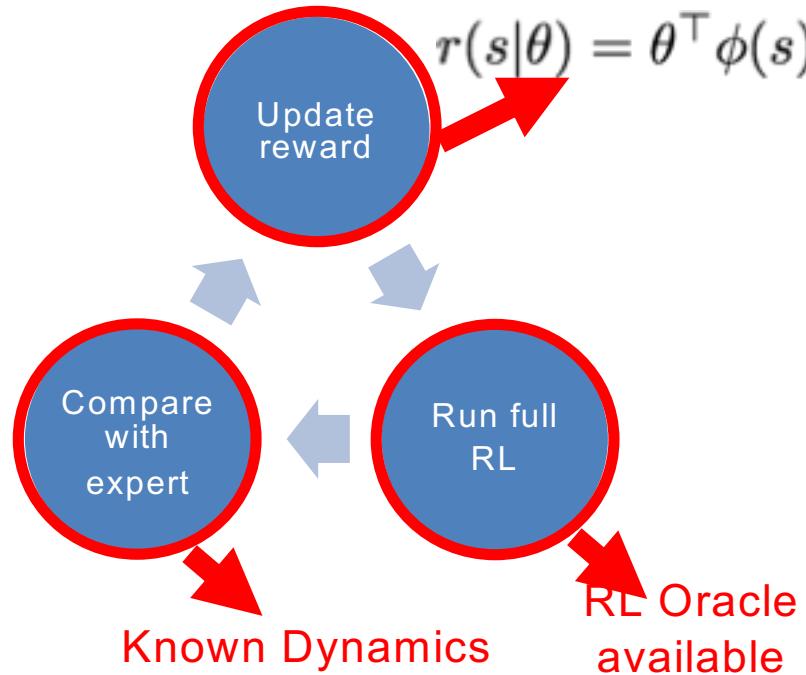
Imitation Learning via Inverse RL (model-given)

Abbeel & Ng, ICML '04

Syed & Schapire, NIPS '07

Assumptions:

$$(S, A, P, \cancel{r}, \gamma, p_0)$$



Syed & Schapire, NIPS '07

Abbeel & Ng, ICML '04

$$\|\theta\|_1 = 1, \theta \succeq 0$$

$$\|\theta\|_2 \leq 1$$

Linear Reward \rightarrow Feature Expectations

Abbeel & Ng, ICML '04
Syed & Schapire, NIPS '07

Assume: $r(s) = \theta \cdot \phi(s)$

Value of a policy expressed in terms of feature expectation

$$\begin{aligned} V(\pi|s_0) &= \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \theta \cdot \phi(s_t) | \pi \right] \\ &= \theta \cdot \underbrace{\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t) \right]}_{\mu(\pi)} \end{aligned}$$

feature expectations 

Feature Matching → Optimality

Abbeel & Ng, ICML '04
Syed & Schapire, NIPS '07

Feature Expectation:

$$\mu(\pi) = \mathbb{E} [\text{visited state features} \mid \pi]$$

If reward r is linear:

$$\mu(\pi) = \mu(\pi^*) \Rightarrow V(\pi) = V(\pi^*)$$



Don't know exactly

Feature Matching in Inverse RL (model-given)

Abbeel & Ng, ICML '04

Find π s.t $\|\mu(\pi) - \mu(\pi^*)\|_2 \leq \epsilon$

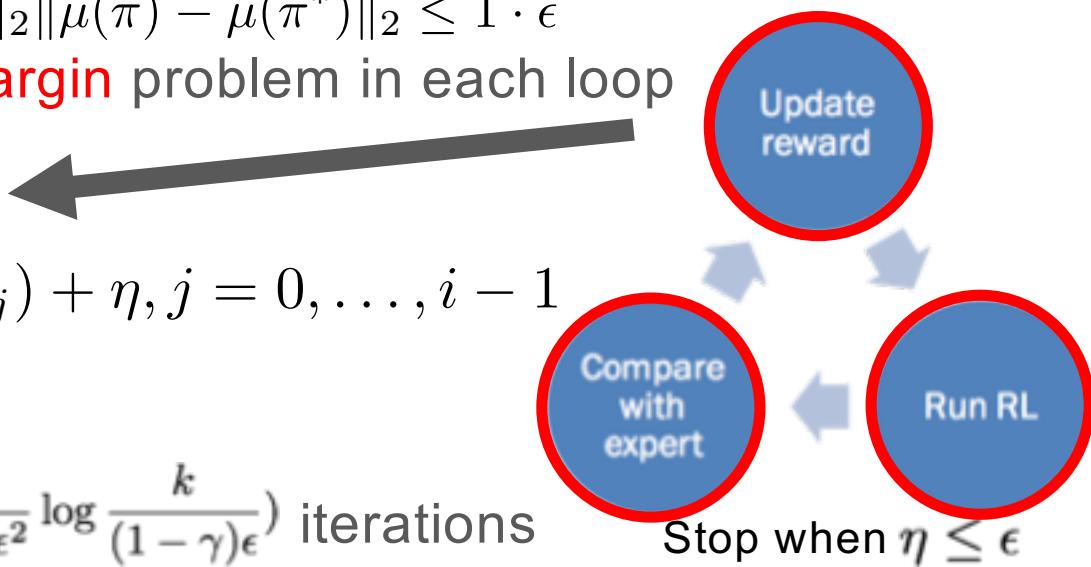
Why? $|V(\pi) - V(\pi^*)| = |\theta^\top \mu(\pi) - \theta^\top \mu(\pi^*)|$
 $\leq \|\theta\|_2 \|\mu(\pi) - \mu(\pi^*)\|_2 \leq 1 \cdot \epsilon$

Algorithm: solve **max-margin** problem in each loop

$$\max_{\eta, \theta} \quad \eta$$

s.t. $\theta^\top \mu(\pi^*) \geq \theta^\top \mu(\pi_j) + \eta, j = 0, \dots, i - 1$
 $\|\theta\|_2 \leq 1$

Theory: at most $O\left(\frac{k}{(1-\gamma)^2 \epsilon^2} \log \frac{k}{(1-\gamma)\epsilon}\right)$ iterations

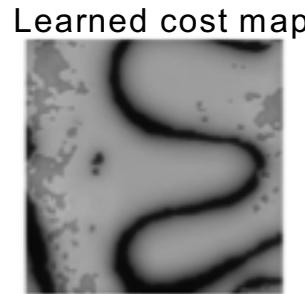


See also...

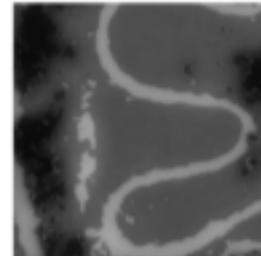
Maximum Margin Planning - Ratliff et al., ICML 2006

- Generalizes to multiple MDPs
- Agnostic about a real underlying MDP or reward function

Mode 1: follow the road



Mode 2: "hide" in the trees



Experiment: learning to plan based on satellite color imagery

See also...

A Game-Theoretic Approach to Apprenticeship Learning – Syed & Schapire, NIPS 2007

- Restricted to $\|\theta\|_1 = 1, \theta \succeq 0$

- Game-theoretic formulation

$$\max_{\pi \in \Pi} \min_{\theta} [\theta^\top \mu(\pi) - \theta^\top \mu(\pi^*)]$$

- Iteratively learn θ by Multiplicative Weight Update

Feature Expectation Matching

Want to find π s.t discounted state visitation features match the expert $\mu(\pi) \approx \mu(\pi^*)$

Ill-posed problem → regularization

Another regularization method: maximum entropy

Policy Learning is Still Ambiguous

1. Many reward functions correspond to the same policy
2. Many stochastic mixtures of policies correspond to the same feature expectation

$$\left. \begin{array}{l} \mu(\pi_1) \approx \mu(\pi^*) \\ \mu(\pi_2) \approx \mu(\pi^*) \end{array} \right\} \mu(\alpha\pi_1 + (1 - \alpha)\mu_2) \approx \mu(\pi^*)$$

Maximum Entropy Principle

Policy π induces distribution over trajectories $P(\tau)$

Feature matching:

$$\sum_{\tau} P(\tau) \mu(\tau) = \mu(\pi^*)$$

$$\sum_{\tau} P(\tau) = 1$$

Maximum entropy principle: The probability distribution which best represents the current state of knowledge is the one with largest entropy (E.T. Jaynes 1957)

Maximum Entropy Principle

Choose trajectory distribution that satisfies feature matching constraints without over-committing

$$\begin{aligned} \max_P \quad & - \sum_{\tau} P(\tau) \log P(\tau) && \xrightarrow{\text{As uniform as possible}} H(\pi) \\ \text{s.t.} \quad & \sum_{\tau} P(\tau) \mu(\tau) = \hat{\mu}^* && \xrightarrow{\frac{1}{m} \sum_{\tau_i \in \mathcal{D}} \mu(\tau_i)} \\ & \sum_{\tau} P(\tau) = 1 \end{aligned}$$

Maximum Entropy Principle

Constrained Optimization

Lagrangian

Maximize
Unconstrained
Objective

$$r_\theta(s) = \theta^\top \phi(s)$$

$$\Rightarrow P(\tau|\theta) \propto e^{\theta^\top \phi(\tau)}$$

The distribution that maximizes entropy given linear constraints is in the exponential family

MaxEnt Inverse RL (model-given)

Ziebart et al., AAAI '08

MaxEnt formulation: $P(\tau|\theta) = \frac{1}{Z(\theta)} e^{\sum_{s_t \in \tau} \theta^\top \phi(s_t)}$

$$Z(\theta) = \int e^{r(\tau|\theta)} d\tau$$

reward inference: max log likelihood

$$\theta^* = \operatorname{argmax}_{\theta} L(\theta) = \operatorname{argmax}_{\theta} \sum_{\tau_i \in \mathcal{D}} \log P(\tau_i|\theta)$$

MaxEnt Inverse RL (model-given)

Ziebart et al., AAAI '08

Gradient descent over log-likelihood

$$\begin{aligned}\nabla_{\theta} L(\theta) &= \frac{1}{m} \sum_{\tau_i \in \mathcal{D}} \underbrace{\mu(\tau_i)}_{\text{expert state feature}} - \sum_s \underbrace{d_s^{\theta}}_{\text{state occupancy measure}} \phi(s) \\ &= \frac{1}{T} \sum_{s' \in \tau_i} \phi(s')\end{aligned}$$

Dynamic programming: state occupancy measure (visitation freq)

$$d_{t+1,s'} = \sum_a \sum_s d_{t,s} \pi_{\theta}(a|s) P(s'|s, a)$$

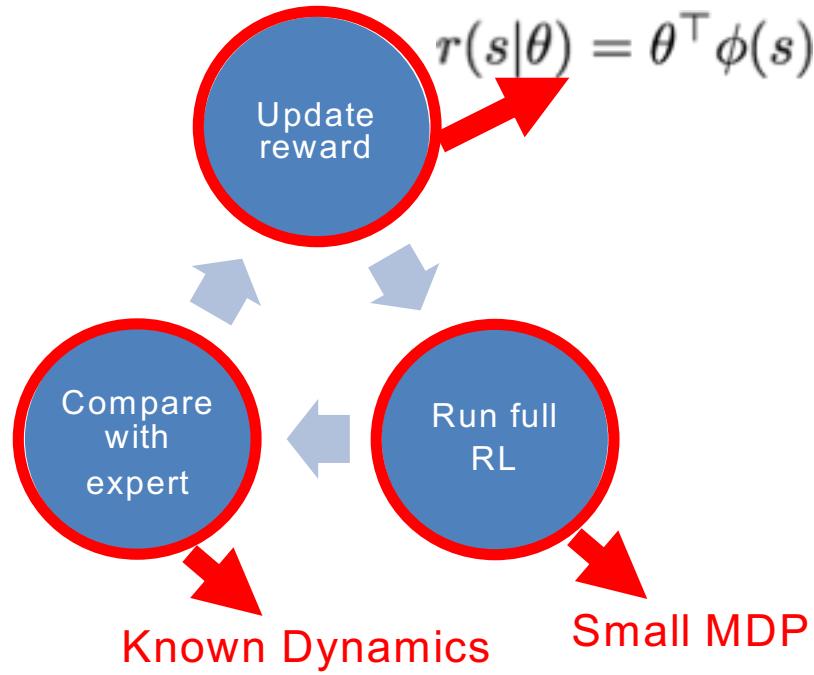
Solve forward RL w.r.t. θ

Given from the model

Inverse RL So Far...

Model-given:

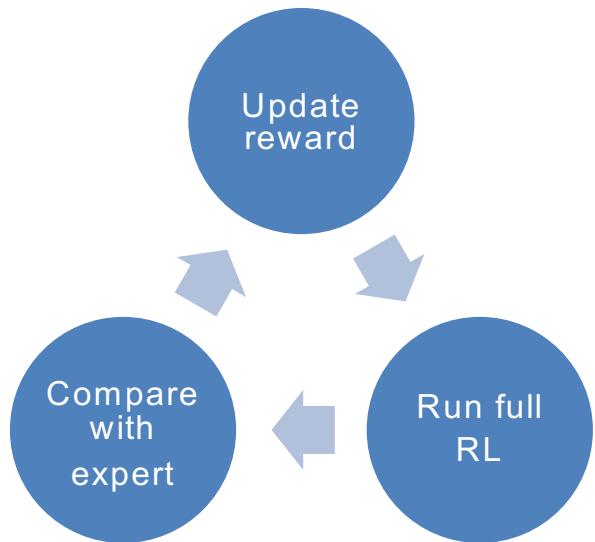
$$(S, A, P, \cancel{X}, \gamma, p_0)$$



Next ...

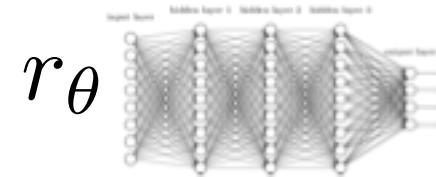
Model-given:

$(S, A, P, \cancel{R}, \gamma, p_0)$



Model-free:

$(S, A, \cancel{P}, \gamma, p_0)$



Interact with
environment / simulator
(Model-free)

Large / continuous
space

MaxEnt Deep IRL + Model-Free Optimization

Guided Cost Learning: Deep Inverse Optimal Control via Policy Optimization - Finn et al., ICML '16

reward $r_\theta(s_t, a_t)$ parameterized by neural net

Recall MaxEnt formulation $P(\tau|\theta) = \frac{1}{Z(\theta)} e^{r_\theta(\tau)}$

$$Z(\theta) = \int e^{r(\tau|\theta)} d\tau$$

Max likelihood objective

$$\begin{aligned}\theta^* &= \operatorname{argmax}_\theta L(\theta) = \operatorname{argmax}_\theta \sum_{\tau_i \in \mathcal{D}} \log P(\tau_i|\theta) \\ &= \operatorname{argmax}_\theta \frac{1}{|\mathcal{D}|} \sum_{\tau_i \in \mathcal{D}} r_\theta(\tau_i) - \log Z(\theta)\end{aligned}$$



Need to sample to evaluate gradients

MaxEnt Deep IRL + Model-Free Optimization

Finn et al., ICML '16

Approximate $Z(\theta) = \int e^{r(\tau|\theta)} d\tau$

Use “proposal” distribution $q(\tau)$
to sample trajectories $\mathcal{D}_{\text{samp}}$

$$Z(\theta) \approx \text{average}_{\tau \in \mathcal{D}_{\text{samp}}} \left(\frac{e^{r_\theta(\tau)}}{q(\tau)} \right)$$



Reward Optimization

Optimal dist $q(\tau) \propto e^{r_\theta(\tau)}$,
which depends on optimal
policy w.r.t. unknown r_θ



Policy Optimization

MaxEnt Deep IRL + Model-Free Optimization

Finn et al., ICML '16

Estimated log likelihood:

$$L(\theta) \approx \frac{1}{m} \sum_{\tau_i \in \mathcal{D}_{\text{demo}}} r_\theta(\tau_i) + \log \frac{1}{n} \sum_{\tau_j \in \mathcal{D}_{\text{samp}}} \frac{e^{r_\theta(\tau_j)}}{q(\tau_j)}$$

with gradient:

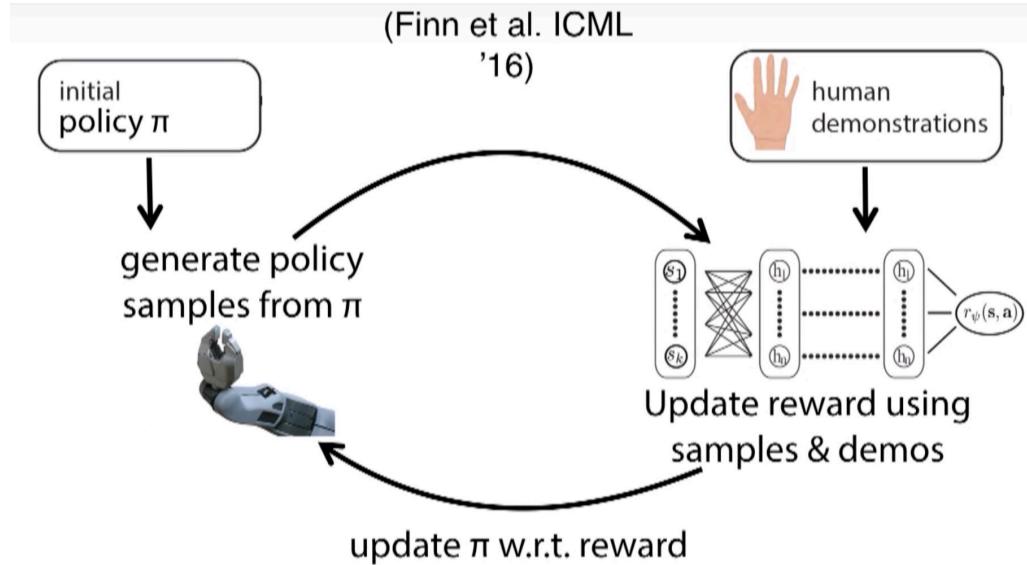
 sampling distribution
(also depends on θ) $q(\tau) \propto e^{r_\theta(\tau)}$

$$\nabla_\theta L(\theta) \approx \frac{1}{m} \sum_{\tau_i \in \mathcal{D}_{\text{demo}}} \nabla_\theta r_\theta(\tau_i) + \frac{1}{Z} \sum_{\tau_j \in \mathcal{D}_{\text{samp}}} w_j \nabla_\theta r_\theta(\tau_j)$$

$$\text{where } w_j = \frac{e^{r_\theta(\tau_j)}}{q(\tau_j)} \text{ and } Z = \sum_j w_j$$

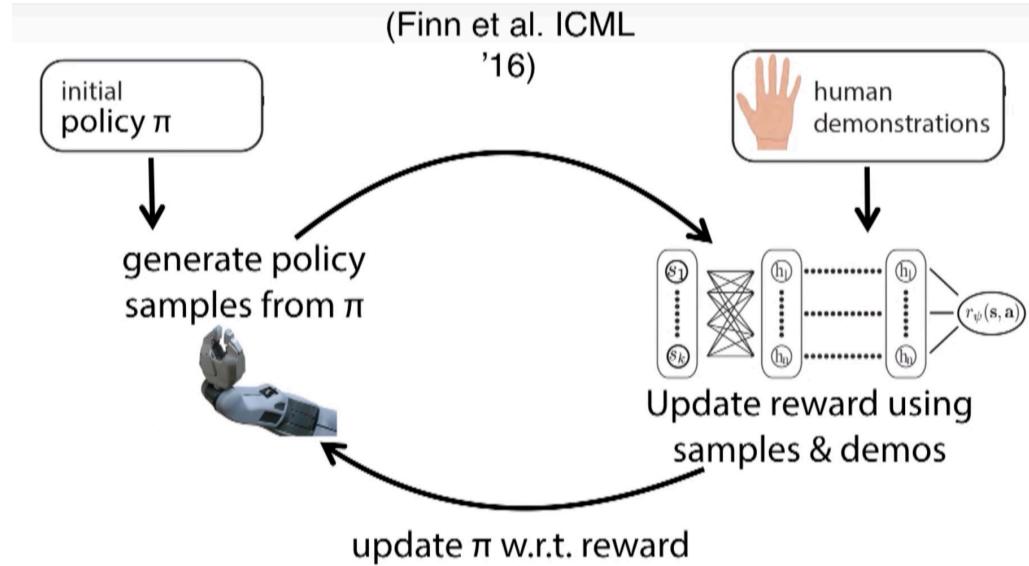
MaxEnt Deep IRL + Model-Free Optimization

- Generate samples by preventing the policy from changing too rapidly
- Update policy: take only 1 policy optimization step



Connection to Generative Adversarial Learning*

- Policy: Generator
- Reward Function: Discriminator



*Generative Adversarial Imitation Learning – Ho & Ermon – NIPS ‘16

IL via Occupancy Measure Matching

Apprenticeship Learning Using Linear Programming - Syed, Bowling & Schapire, ICML '08

Occupancy measure

$$d_{sa}^\pi = \text{visitation frequency of } (s, a) \Big| \text{following } \pi$$

$$d_{sa}^\pi = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{1}_{(s_t = s, a_t = a)} \Big| \pi \right]$$

Fact:

$$V(\pi) = \mathbb{E}_\pi [r(s, a)] = \sum_{s,a} r(s, a) d_{sa}^\pi$$

IL Using Linear Programming (model-given)

Syed, Bowling & Schapire, ICML '08

Imitation learning \Rightarrow occupancy measure matching

$$d_{sa}^\pi = d_{sa}^* \Rightarrow \mathbb{E}_\pi [r(s, a)] = \mathbb{E}_{\pi^*} [r(s, a)]$$

Direct imitation learning by solving dual LP of original MDP



Occupancy Measure Matching + MaxEnt

$$d_{sa}^{\pi} = d_{sa}^{*} \Rightarrow \mathbb{E}_{\pi} [r(s, a)] = \mathbb{E}_{\pi^*} [r(s, a)]$$

for any reward function $r(s, a)$

Ill-posed occupancy matching \rightarrow MaxEnt principle again:

$$\max_{\pi} H(\pi)$$

$$\text{s.t. } \text{distance}(d^{\pi}, d^*) \leq \epsilon$$

Again, turning into unconstrained optimization:

$$\min_{\pi} \text{distance}(d^{\pi}, d^*) - \lambda H(\pi)$$

Occupancy Matching + Model-free Optimization

$$\min_{\pi} \text{distance}(d^\pi, d^*) - \lambda H(\pi)$$

Generative Adversarial Imitation Learning - Ho & Ermon. NIPS '16:

$$\begin{aligned}\text{distance}(d^\pi, d^*) &= \max_{D \in (0,1)^{S \times A}} \mathbb{E}_\pi [\log(D(s, a))] + \mathbb{E}_{\pi^*} [\log(1 - D(s, a))] \\ &\approx D_{KL}(d^\pi \|(d^\pi + d^*)/2) + D_{KL}(d^* \|(d^\pi + d^*)/2)\end{aligned}$$



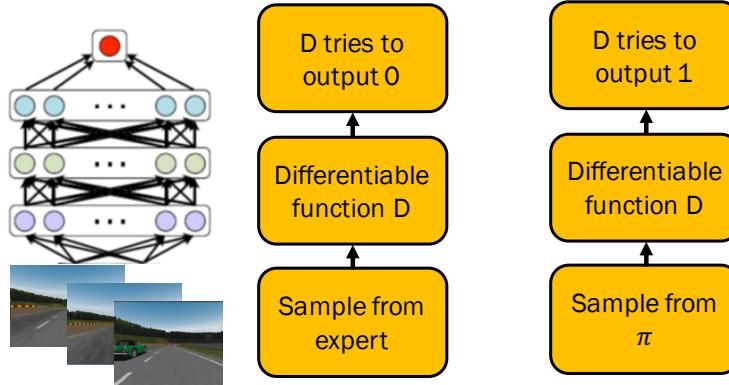
Jensen-Shannon divergence

Generative Adversarial Imitation Learning

Find saddle point (π, D)

Ho & Ermon, NIPS '16

$$\min_{\pi} \max_{D \in (0,1)^{S \times A}} \mathbb{E}_{\pi} [\log(D(s, a))] + \mathbb{E}_{\pi^*} [\log(1 - D(s, a))] - \lambda H(\pi)$$



Discriminator Update

$$\hat{\mathbb{E}}_{\tau_w} [\nabla_w \log(D_w(s, a))] + \hat{\mathbb{E}}_{\tau^*} [\log(1 - D_w(s, a))]$$

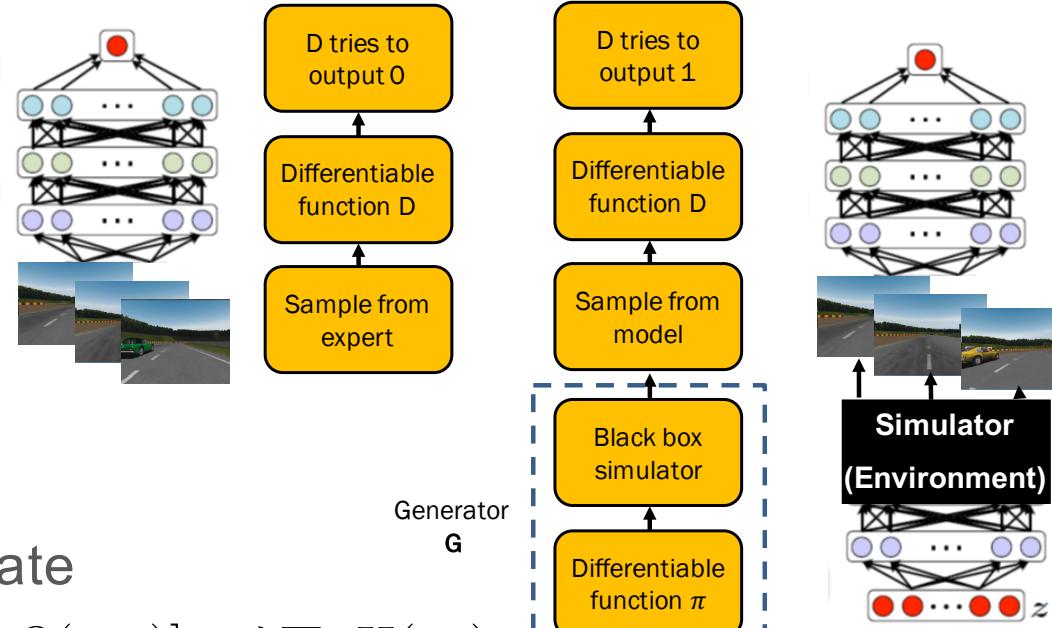
Thanks to materials from Stefano Ermon

Generative Adversarial Imitation Learning

Find saddle point (π, D)

Ho & Ermon, NIPS '16

$$\min_{\pi} \max_{D \in (0,1)^{S \times A}} \mathbb{E}_{\pi} [\log(D(s, a))] + \mathbb{E}_{\pi^*} [\log(1 - D(s, a))] - \lambda H(\pi)$$



Generator Update

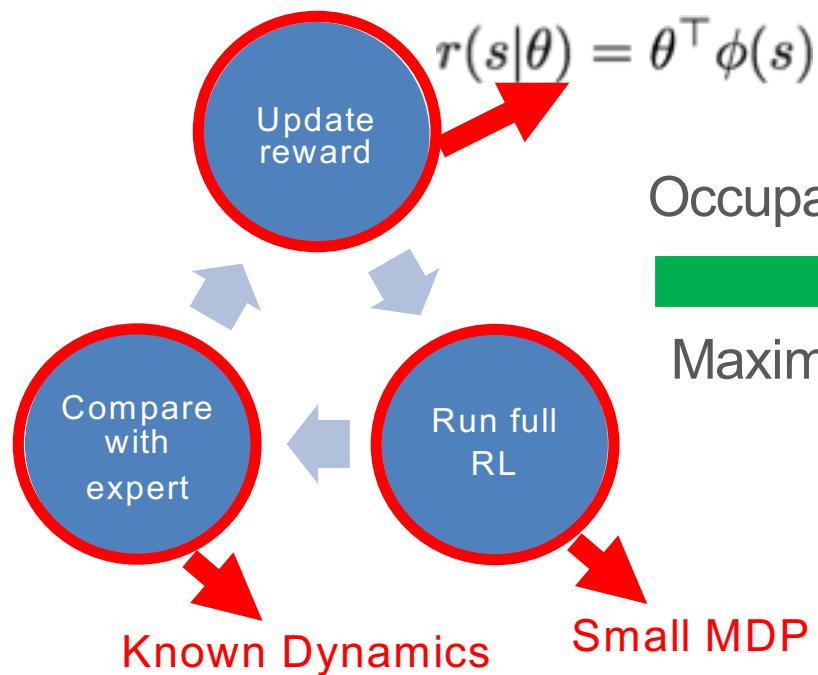
$$\hat{\mathbb{E}}_{\tau_\pi} [\nabla_\theta \log \pi_\theta(a|s) Q(s, a)] - \lambda \nabla_\theta H(\pi_\theta)$$

Thanks to materials from Stefano Ermon

Inverse Reinforcement Learning

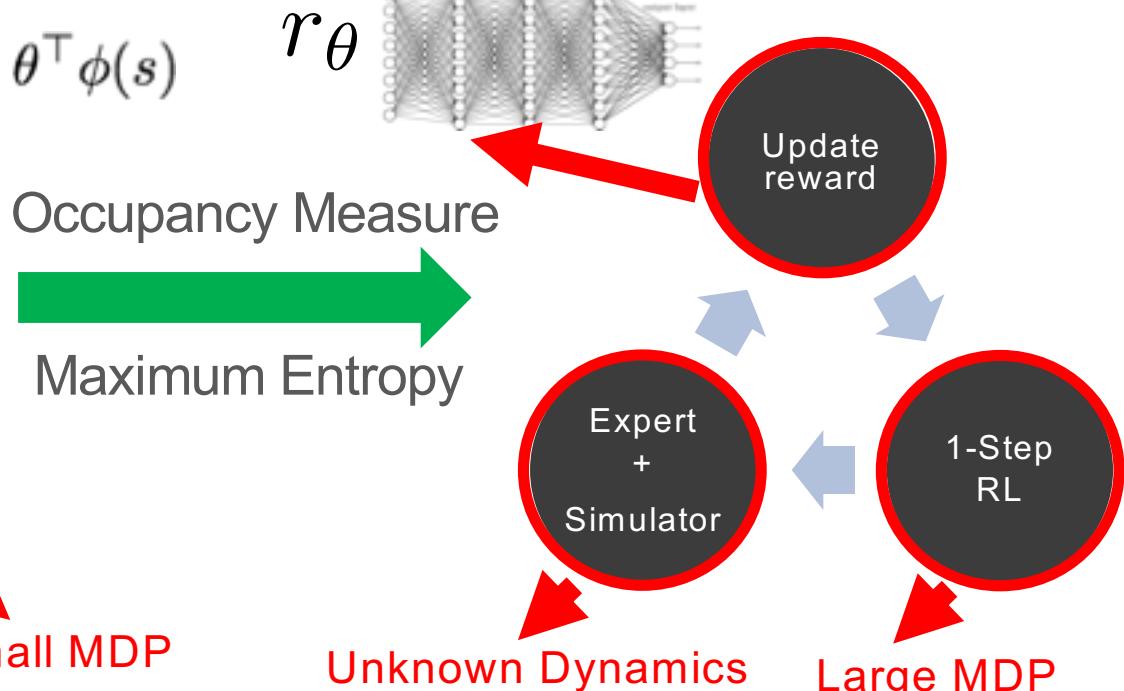
Model-given:

$$(S, A, P, \cancel{r}, \gamma, p_0)$$



Model-free:

$$(S, A, \cancel{P}, \gamma, p_0)$$



List of Core Papers Mentioned

Search-based structured prediction - Daume, Langford & Marcu, Machine Learning 2009

A reduction of imitation learning and structured prediction to no-regret online learning - Ross, Gordon & Bagnell, AISTATS 2011

Efficient Reductions for Imitation Learning - Ross & Bagnell, AISTATS 2010

A reduction from apprenticeship learning to classification - Syed & Schapire, NIPS 2010

Apprenticeship learning via inverse reinforcement learning - Abbeel & Ng, ICML 2004

A Game-Theoretic Approach to Apprenticeship Learning - Syed & Schapire, NIPS 2007

Apprenticeship learning using linear programming - Syed, Bowling & Schapire, ICML 2008

Maximum entropy inverse reinforcement learning - Ziebart, Mass, Bagnell & Dey, AAAI 2008

Generative adversarial imitation learning - Ho & Ermon, NIPS 2016

Guided cost learning: deep inverse optimal control via policy optimization - Finn, Levine & Abbeel, ICML 2016

Imitation Learning

ICML 2018 Tutorial

Break

Yisong Yue



Hoang M. Le



yyue@caltech.edu

hmle@caltech.edu



@YisongYue

@HoangMinhLe



visongyue.com

hoangle.info

Tutorial Overview

Part 1: Introduction and Core Algorithms

- Teaser results
- Overview of ML landscape
- Types of Imitation learning
- Core algorithms of passive, interactive learning and cost learning

Part 2: Extensions and Applications

- **Speech Animation**
- **Structured prediction and search**
- **Improving over expert**
- **Filtering and sequence modeling**
- **Multi-objective imitation learning**
- **Visual / few-shot imitation learning**
- **Domain Adaptation imitation learning**
- **Multi-agent imitation learning**
- **Hierarchical imitation learning**
- **Multi-modal imitation learning**
- **Weaker Feedback**

Speech Animation

A Deep Learning Approach for Generalized Speech Animation

[Taylor et al., SIGGRAPH 2017]

Input sequence $X = \langle x_1, \dots, x_T \rangle$

Output sequence $Y = \langle y_1, \dots, y_T \rangle$ $y \in R^D$

Goal: learn predictor $\pi : X \rightarrow Y$

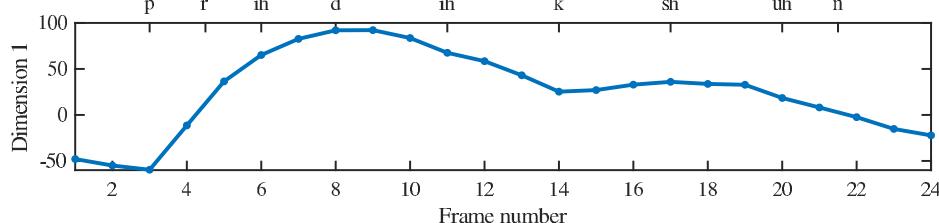
X

Frame	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
Token	-	p	p	r	ih	ih	d	d	ih	ih	ih	ih	ih	k	k	sh	sh	sh	uh	uh	n	-



Phoneme sequence

Y



Sequence of face configurations



Train using
Behavioral Cloning

A Deep Learning Approach for Generalized Speech Animation

Sarah Taylor, Taehwan Kim, Yisong Yue et al., SIGGRAPH 2017

<https://www.youtube.com/watch?v=9zL7aeiW9fE>



German

A Deep Learning Approach for Generalized Speech Animation

Sarah Taylor, Taehwan Kim, Yisong Yue et al., SIGGRAPH 2017

<https://www.youtube.com/watch?v=9zL7aeiW9fE>





Input Audio

s s s s s ih ih ih g g r r ae ae ae ae fff



Speech Recognition



Speech Animation

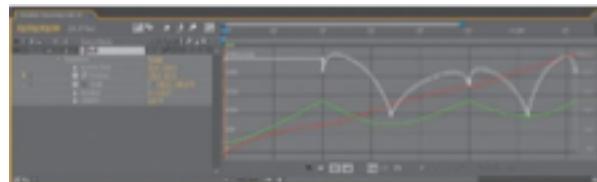


(chimp rig courtesy of Hao Li)



Retargeting

E.g., [Sumner & Popovic 2004]



Editing

Structured Prediction & Learning to Search

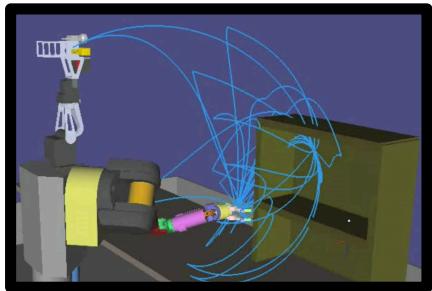
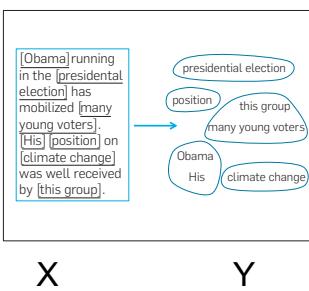
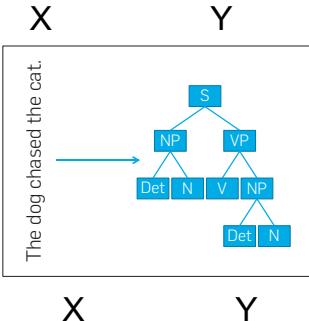
Structured Prediction

Learn mapping from structured space X to structured space Y

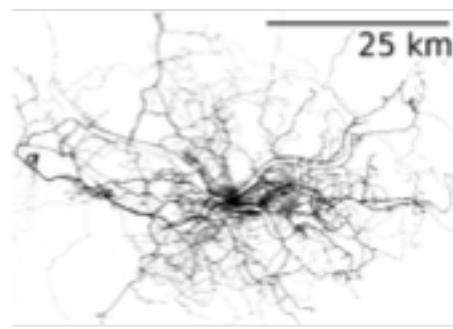
Typically requires solving optimization problem at prediction

- E.g., Viterbi algorithm for MAP inference in HMMs

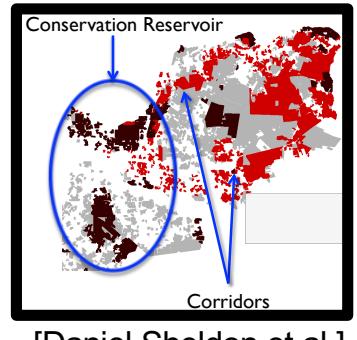
Structured Prediction



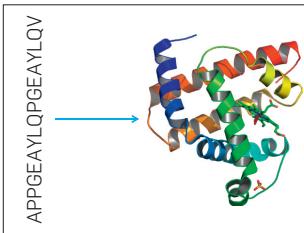
[Debadeepta Dey et al.]
[Nathan Ratliff et al.]



[Brian Ziebart et al., 2008]



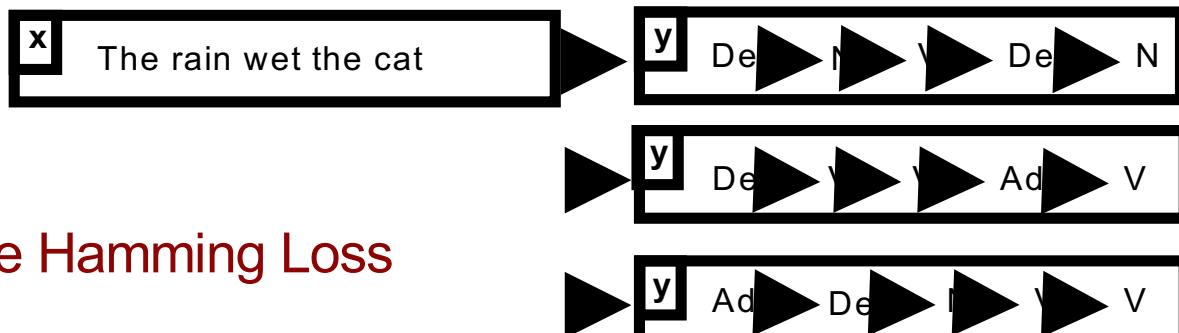
[Daniel Sheldon et al.]
[Daniel Golovin et al.]



Example: Sequence Prediction

Part-of-Speech Tagging

- Given a sequence of words x
- Predict sequence of tags y (dynamic programming)



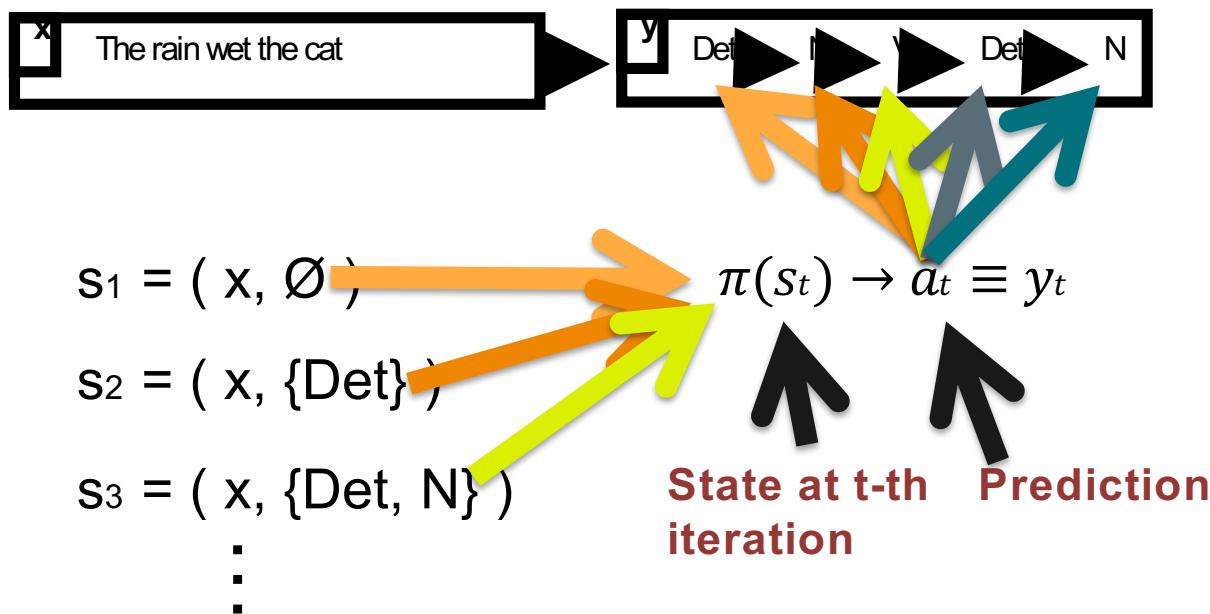
Goal: Minimize Hamming Loss

Goal: Learn Decision Function

(In Inference/Optimization Procedure)

- Sequentially predict outputs:

Interactive Feedback via Structured Labels
Supervised Learning on (s, a)



Data Aggregation (DAgger) for Sequence Labeling (Basic Version)



- Iteration 1: Rollout π_0 (memorize training set)

$$P_1 \left\{ \begin{array}{ll} s_1 = (x, \emptyset) & a_1 = \text{Det} \\ s_2 = (x, \{\text{Det}\}) & a_2 = \text{N} \\ s_3 = (x, \{\text{Det}, \text{N}\}) & a_3 = \text{V} \\ \vdots & \vdots \end{array} \right. \quad \text{Train } \pi_1 \text{ on } P_1$$

Data Aggregation (DAgger) for Sequence Labeling (Basic Version)



- Iteration 2: Rollout π_1

$$\mathcal{P}_2 \left\{ \begin{array}{ll} s_1 = (x, \emptyset) & a_1 = \text{Det} \\ s_2 = (x, \{\text{N}\}) & a_2 = \text{N} \\ s_3 = (x, \{\text{N}, \text{N}\}) & a_3 = \text{V} \\ \vdots & \vdots \end{array} \right. \quad \text{Train } \pi_2 \text{ on } \mathcal{P}_1 \text{ and } \mathcal{P}_2$$

Data Aggregation (DAgger) for Sequence Labeling (Basic Version)



- Iteration 3: Rollout π_2

$$\mathcal{P}_3 \left\{ \begin{array}{ll} s_1 = (x, \emptyset) & a_1 = \text{Det} \\ s_2 = (x, \{\text{Det}\}) & a_2 = \text{N} \\ s_3 = (x, \{\text{Det}, \text{V}\}) & a_3 = \text{V} \\ \vdots & \vdots \end{array} \right. \quad \text{Train } \pi_3 \text{ on } \mathcal{P}_1 \text{ and } \mathcal{P}_2 \text{ and } \mathcal{P}_3$$

(Basic) Ingredients for Structured Prediction

- Deal with Distributional Drift (otherwise just use behavioral cloning)
 - Related approaches include scheduled sampling:
Scheduled Sampling for Sequence Prediction with Recurrent Neural Networks
[Bengio et al., NIPS 2015]
- Imitation Loss compatible w/ Structured Prediction Loss
 - Hamming Loss \rightarrow 0/1 per-step loss (Hamming Loss is additive)
 - More generally: reduce to cost-sensitive classification loss
- Efficiency considerations
 - **Efficiently Programmable Learning to Search** [Daume et al., 2014]
 - **Deeply aggregated: Differentiable imitation learning for sequential prediction**
[Sun et al., 2017]

Learning Policies for Contextual Submodular Optimization

(Example of Cost-Sensitive Reduction)

Stephane Ross et al., ICML 2013

Training set: (x, F_x)

Context

Monotone Submodular Function

Goal: learn π that sequentially maps x to action set A to maximize F_x

Already selected actions

Learning Reduction: at state $s=(x,A)$ create cost for each action

- Define: $f_{\max} = \max_a F_x(A+a) - F_x(A)$
- $c_s(a) = f_{\max} - (F_x(A+a) - F_x(A))$
- Supervised training example: (s, c_s)

Greedy
Submodular
Regret

Can prove convergence to $(1-1/e)$ -optimal policy!

Structured Prediction \approx Learning to Search

Compile input description x into search problem

Train π to effectively navigate search space

High level supervised reward/cost

- E.g., routing cost, BLEU score, Hamming loss, submodular utility, Jaccard, integer program value, log likelihood of graphical model, etc.

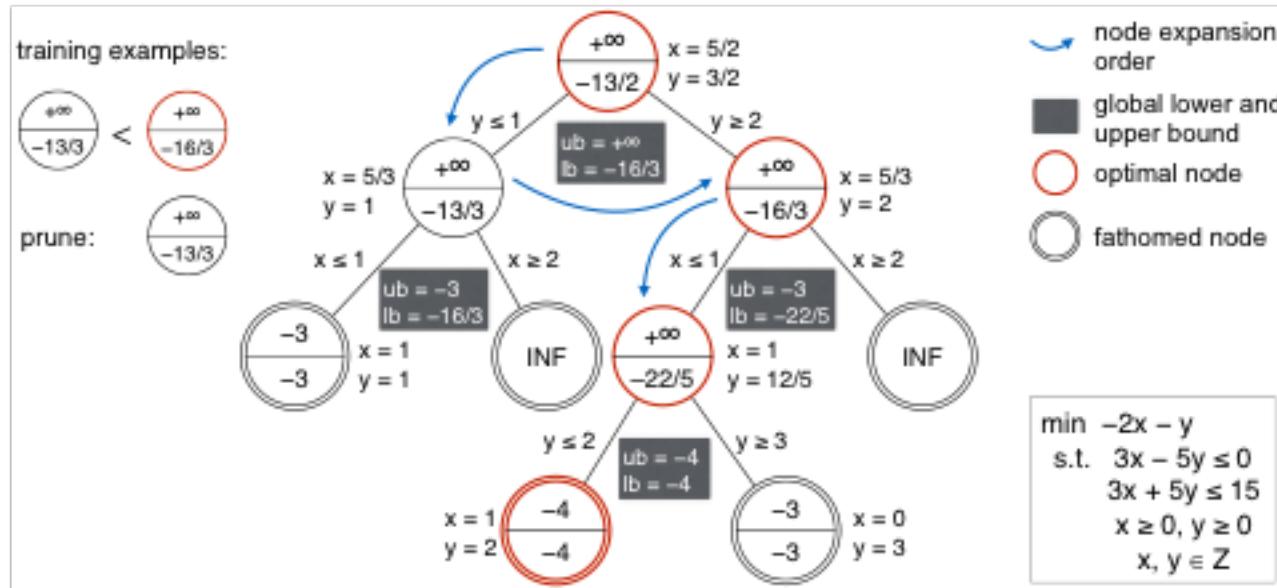
Expert might be weak (search space exponentially large)

Learning to Search in Branch and Bound Algorithms

[He et al., NIPS 2014]

Used to solve Integer Programs

- Two local decisions: ranking and pruning
- Train using solved instances (e.g., by Gurobi)



Sub-Optimal Expert Imitation Learning → Reinforcement Learning

- Query Access to Environment (simple exploration)
 - **Learning to Search via Retrospective Imitation** [Song et al., 2018]
 - **Learning to Search Better than Your Teacher** [Chang et al., ICML 2015]
 - **Reinforcement and imitation learning via interactive no-regret learning**
[Ross & Bagnell, 2014]
- More Sophisticated Exploration (also Query Access to Environment)
 - **Residual Loss Prediction: Reinforcement Learning With No Incremental Feedback**
[Daume et al., ICLR 2018]
- Actor-Critic Style Approaches (also Query Access to Environment)
 - **Truncated Horizon Policy Search: Combining Reinforcement Learning and Imitation Learning** [Sun et al., ICLR 2018]
 - **Sequence Level Training with Recurrent Neural Networks** [Ranzato et al., ICLR 2016]

Learning to Search Better than Your Teacher

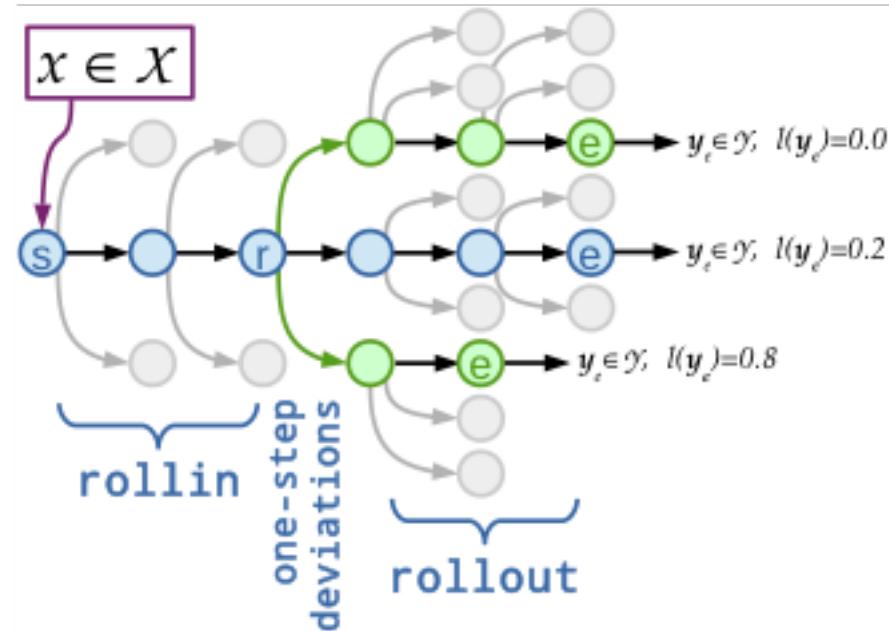
[Chang et al., ICML 2015]

Roll-in: execute policy w/o learning
Roll-out: execute policy for learning

roll-out →	Reference	Mixture	Learned
↓ roll-in	Inconsistent		
Reference			
Learned	Not locally opt.	Good	RL



SEARN [Daume et al., 2009]



Improving Multi-step Prediction of Learned Time Series Models

[Venkatraman et al., AAAI 2015]

$$x_{t+1} = E[f(x_t)]$$

dynamics
state

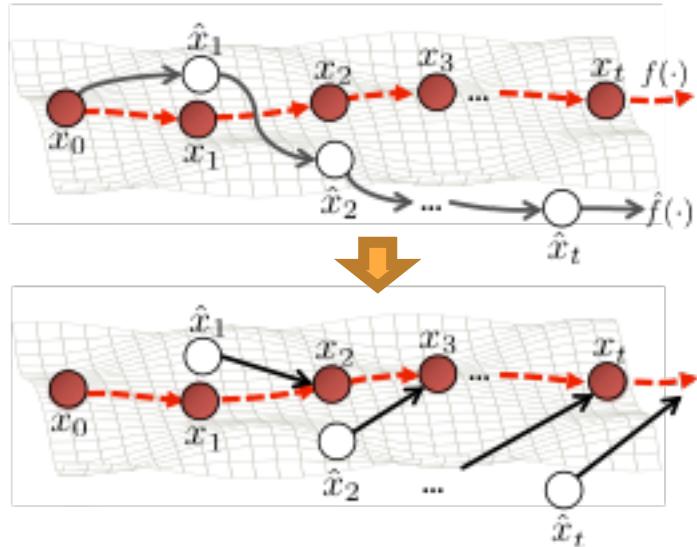
Learn π that maps x_t to x_{t+1}

- Input is previous prediction
- Sequence prediction special case

Can also extend to latent-variables:

- **Learning to Filter with Predictive State Inference Machines** [Sun et al., 2016]

Interactive Imitation Learning



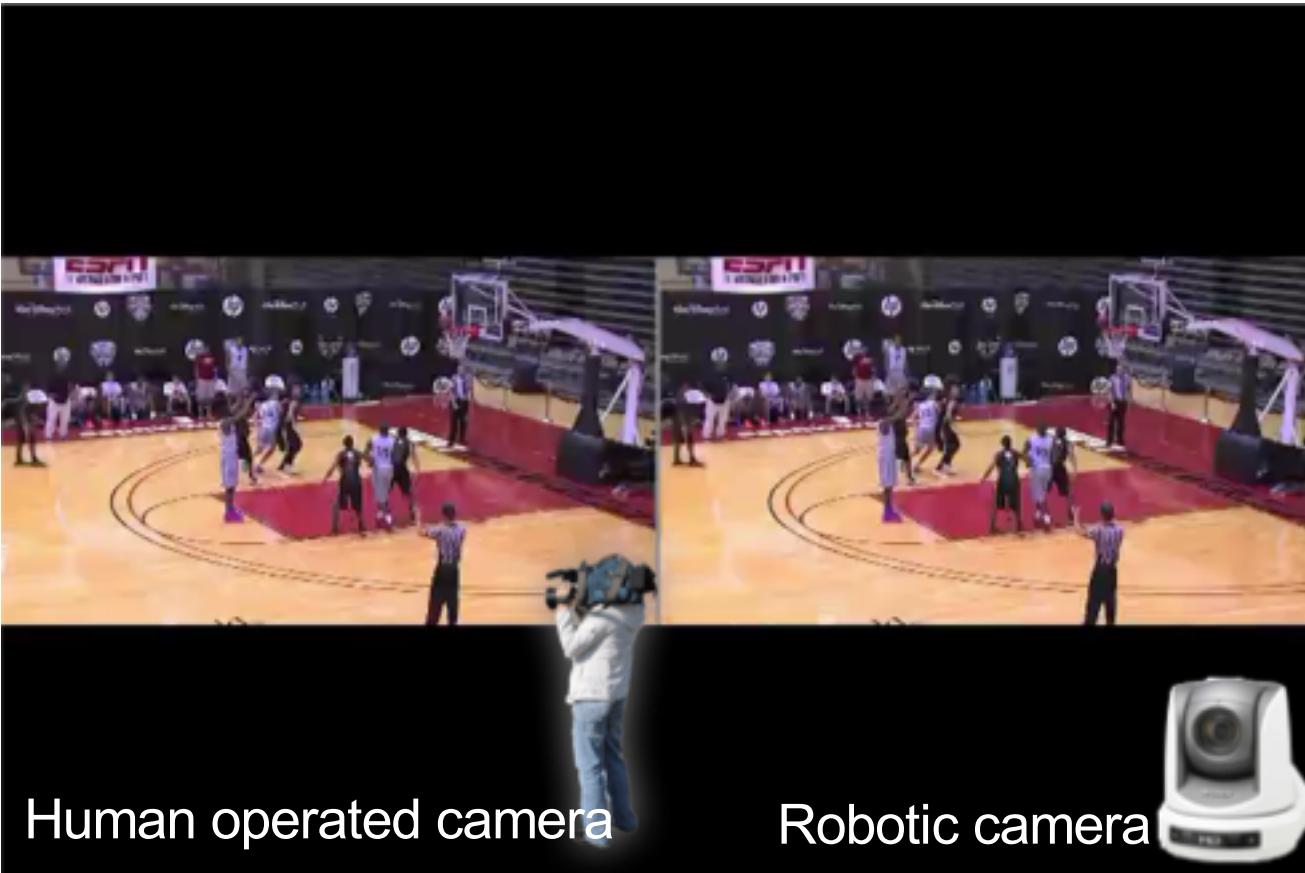
Imitation Learning for Structured Prediction

- Learn decision function to navigate search space
- Reduce to interactive imitation learning
 - Expert feedback from computational oracle
- Careful choice of loss function
- Suboptimal experts
 - Due to computational complexity of solving search space

Imitation Learning with Multiple Objectives

Smooth Imitation Learning

Le et al., ICML '16



Smooth Imitation Learning

Le et al., ICML '16

Input: sequence of context x_t

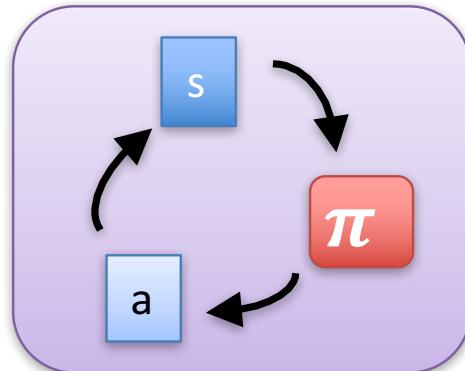
- e.g., noisy player detection

State: $s_t = (x_{t:t-K}, a_{t-1:t-K})$

- recent context and actions

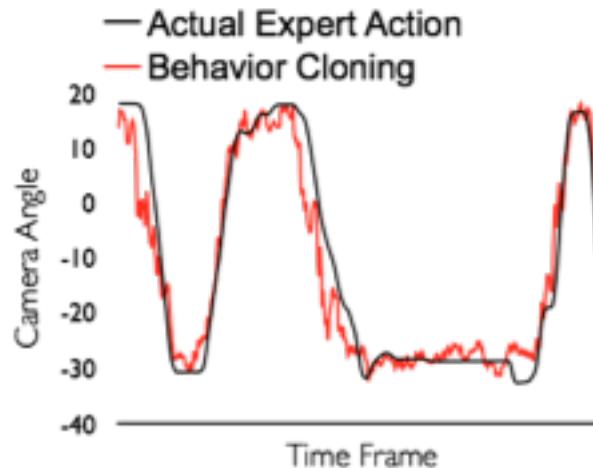
Goal: learn $\pi(s_t) \rightarrow a_t$

- Imitate expert



Smooth Imitation Learning

Le et al., ICML '16

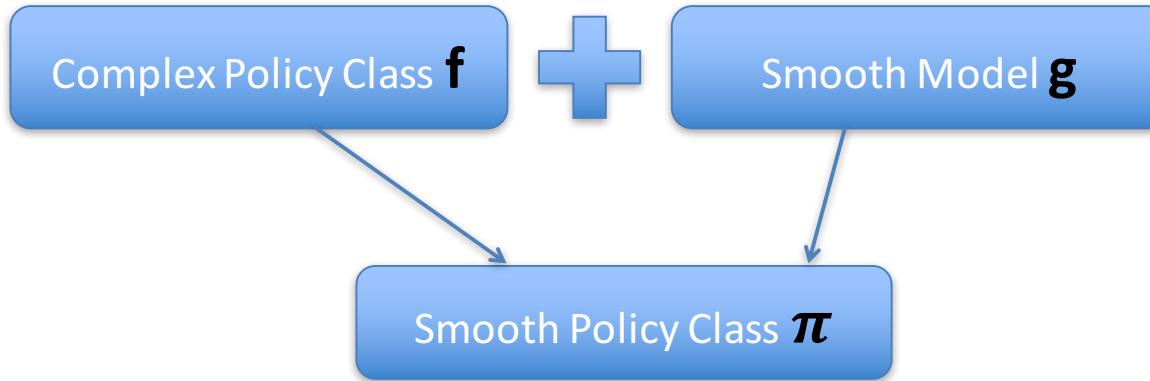


need: Imitation Accuracy + Smoothness

Smooth Imitation Learning

Le et al., ICML '16

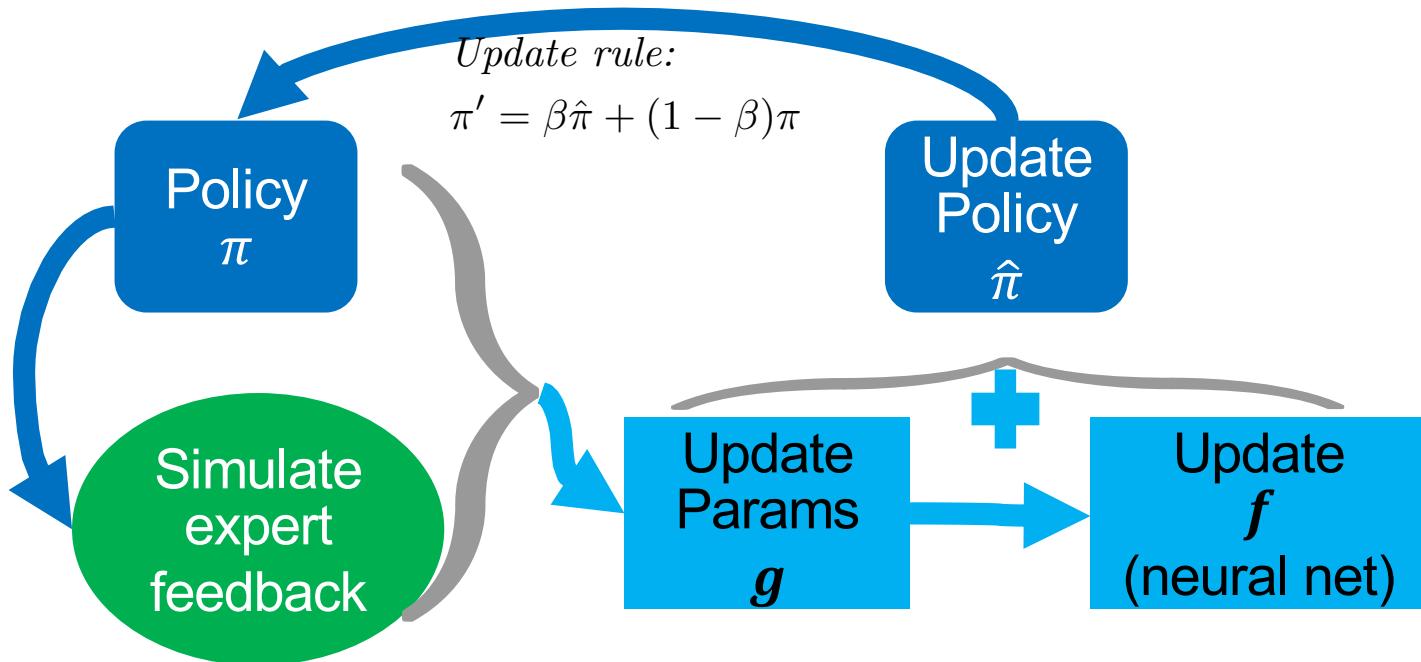
Smooth Imitation by Functional Regularization



$$\begin{aligned}\pi(s) &= \underset{a}{\operatorname{argmin}} (f(s) - a)^2 + \lambda(g(s) - a)^2 \\ &= \frac{f(s) + \lambda g(s)}{1 + \lambda}\end{aligned}$$

Smooth Imitation Learning

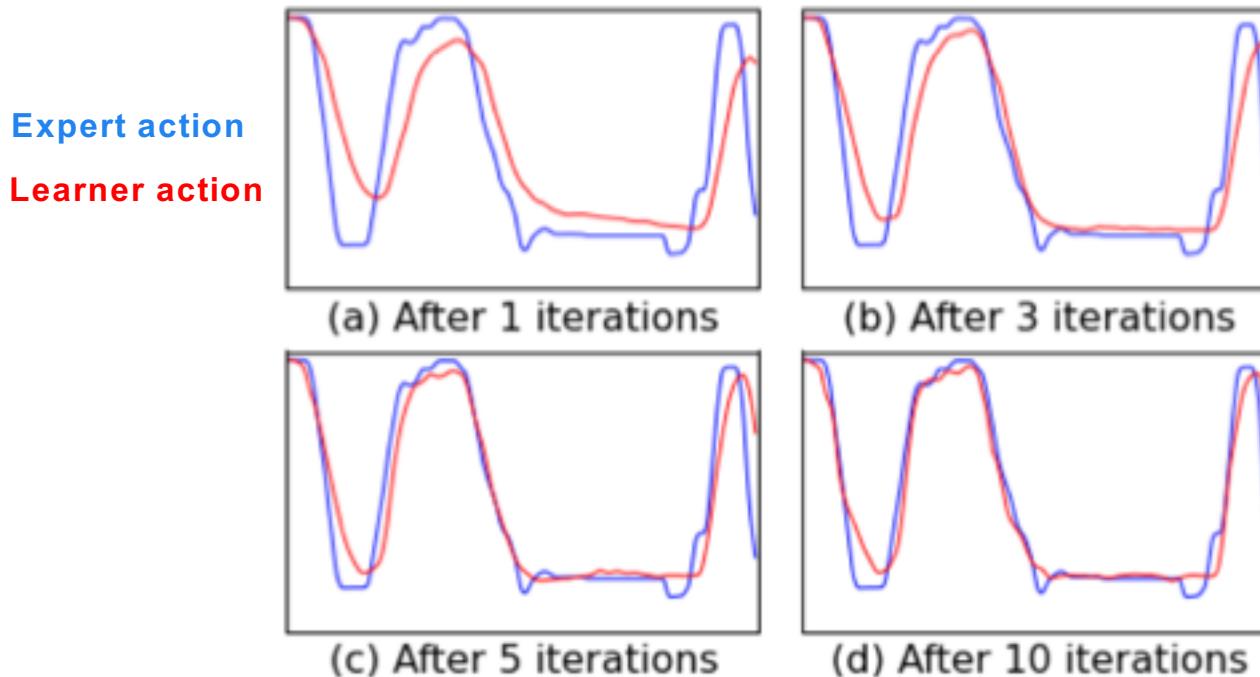
Algorithm: Iterative Imitation Learning for Continuous Spaces



Smooth Imitation Learning

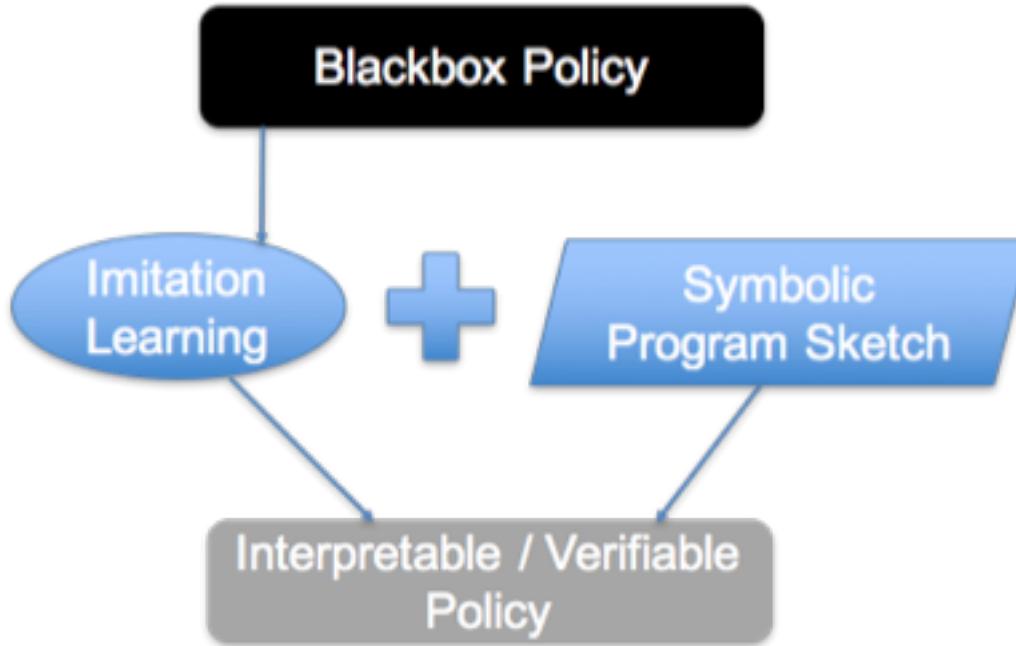
Le et al., ICML '16

Theory: Monotonic policy improvement, adaptive learning rate



Programmatically Interpretable Reinforcement Learning

Verma et al., ICML '18



```
if (0.001 - hd(hTrackPos) > 0) and (0.001 + hd(hTrackPos) > 0)
    then 1.96 + 4.92 * (0.44 - hd(hRPM)) + 0.89 * fold(+, hRPM) + 49.79 * (hd(tl(hRPM)) - hd(hRPM))
else 1.78 + 4.92 * (0.40 - hd(hRPM)) + 0.89 * fold(+, hRPM) + 49.79 * (hd(tl(hRPM)) - hd(hRPM))
```

Safe Imitation Learning for Autonomous Driving

Zhang & Cho, AAAI '17

π = steering angle policy

Safe Frames



Unsafe Frames



TORCS Simulation

Safe Imitation Learning for Autonomous Driving

Zhang & Cho, AAAI '17

c_{safe} = whether π likely to
deviates from π^*

$$c_{\text{safe}}^*(\pi, s) = \begin{cases} 0, & \text{if } \|\pi(s) - \pi^*(s)\| > \epsilon \\ 1, & \text{otherwise} \end{cases}$$

If $c_{\text{safe}}(\pi, s) = 1$, drive with π , otherwise fall back to π^*

Safe Imitation Learning for Autonomous Driving

Zhang & Cho, AAAI '17

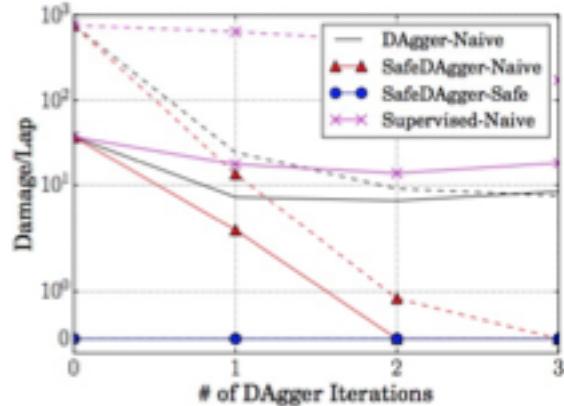
SafeDAgger:

Safety classifier in the loop

Only query “unsafe” states

Results:

- Reduced number of damages
- Faster learning - can be viewed as active DAgger

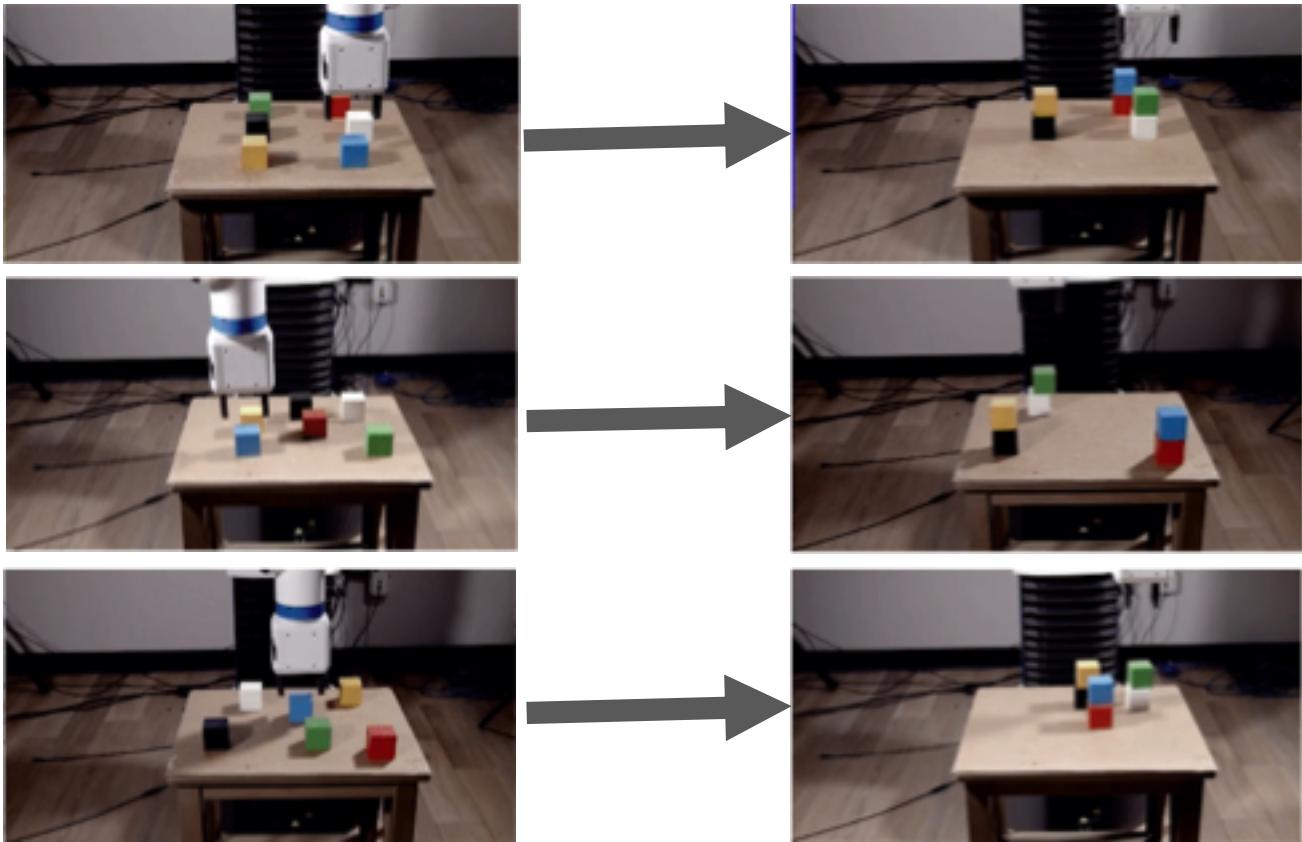


Multi-task, Meta and Transfer Imitation Learning

One Shot Imitation Learning



3 instances of an example task:



Duan et al., NIPS '17

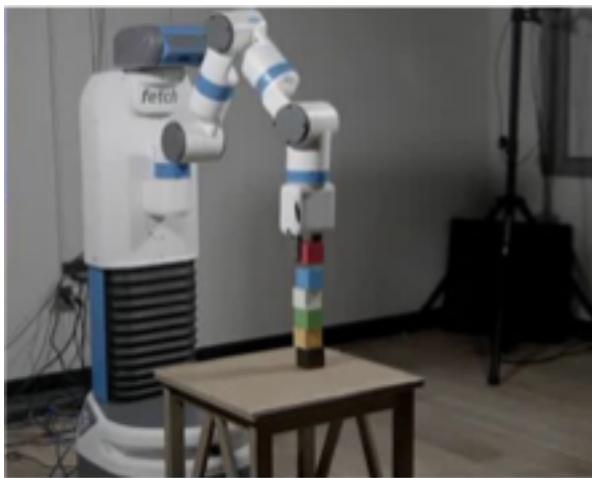
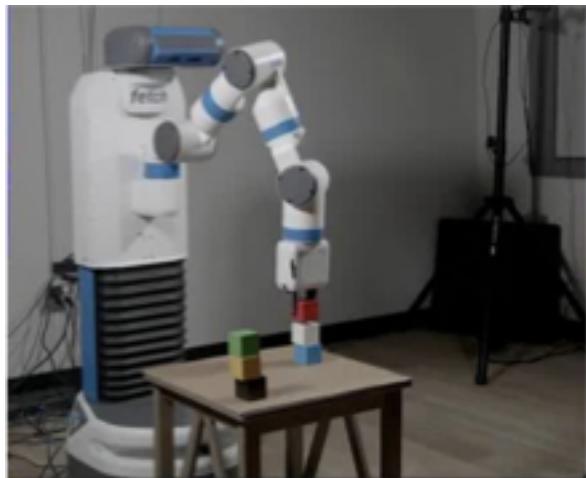
Tobin et al., IROS '17

One Shot Imitation Learning

How about other related tasks?

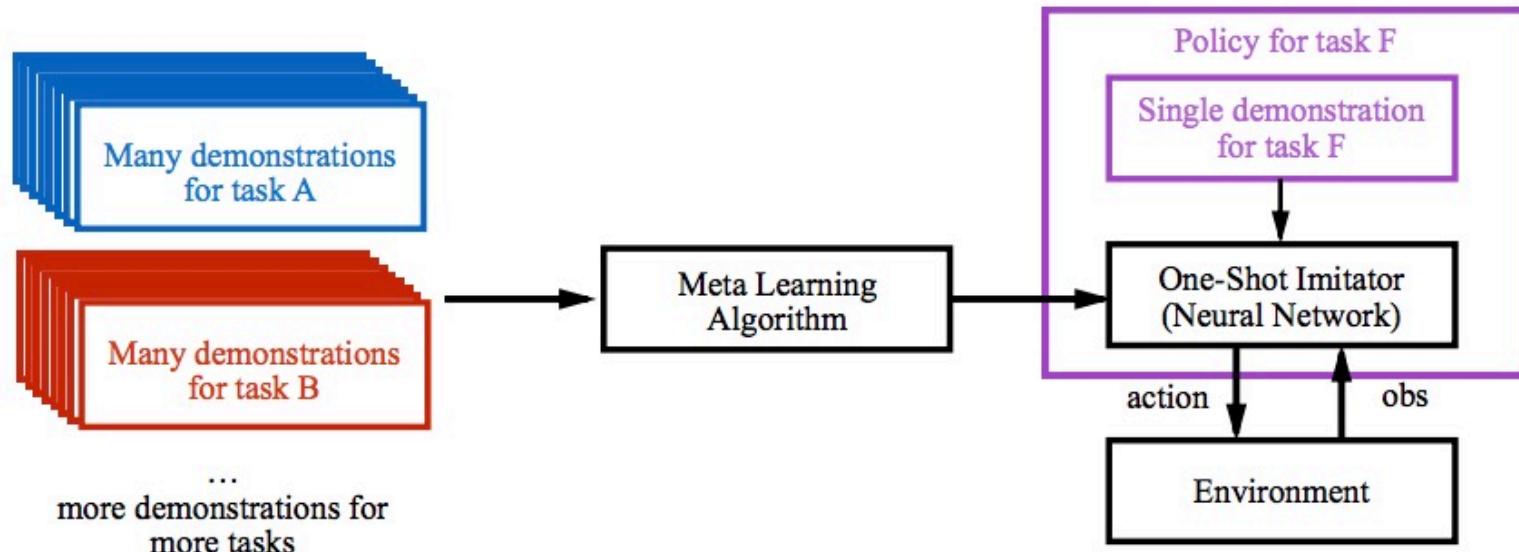
Duan et al., NIPS '17

Tobin et al., IROS '17



One Shot Imitation Learning

Duan et al., NIPS '17

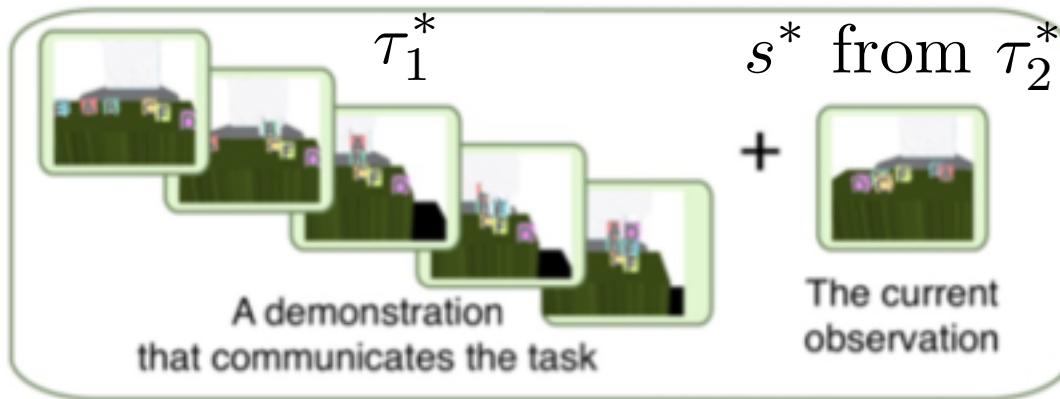
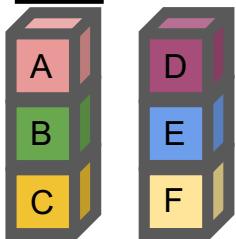


(b) One-Shot Imitation Learning

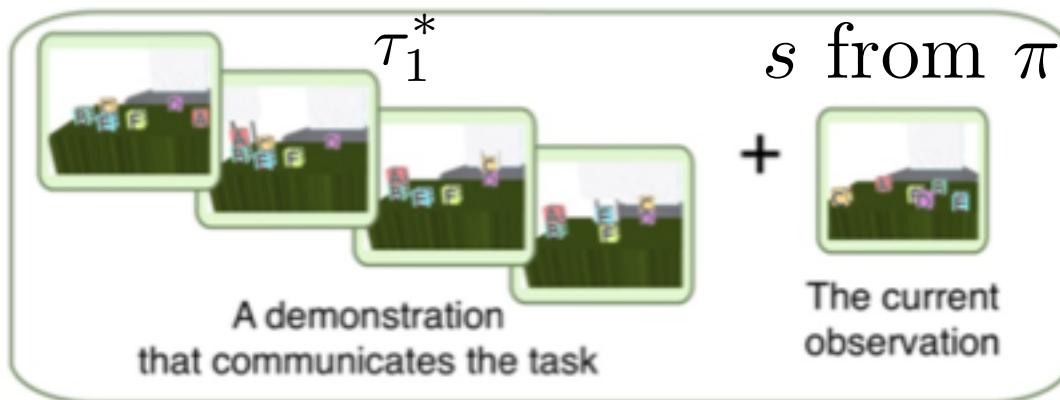
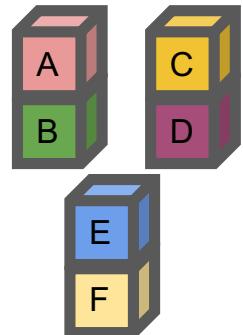
One Shot Imitation Learning

Duan et al., NIPS '17

Training task



Test task



One Shot Visual Imitation Learning - Transfer

Tobin et al., IROS '17

Trained in simulation only
using domain randomization

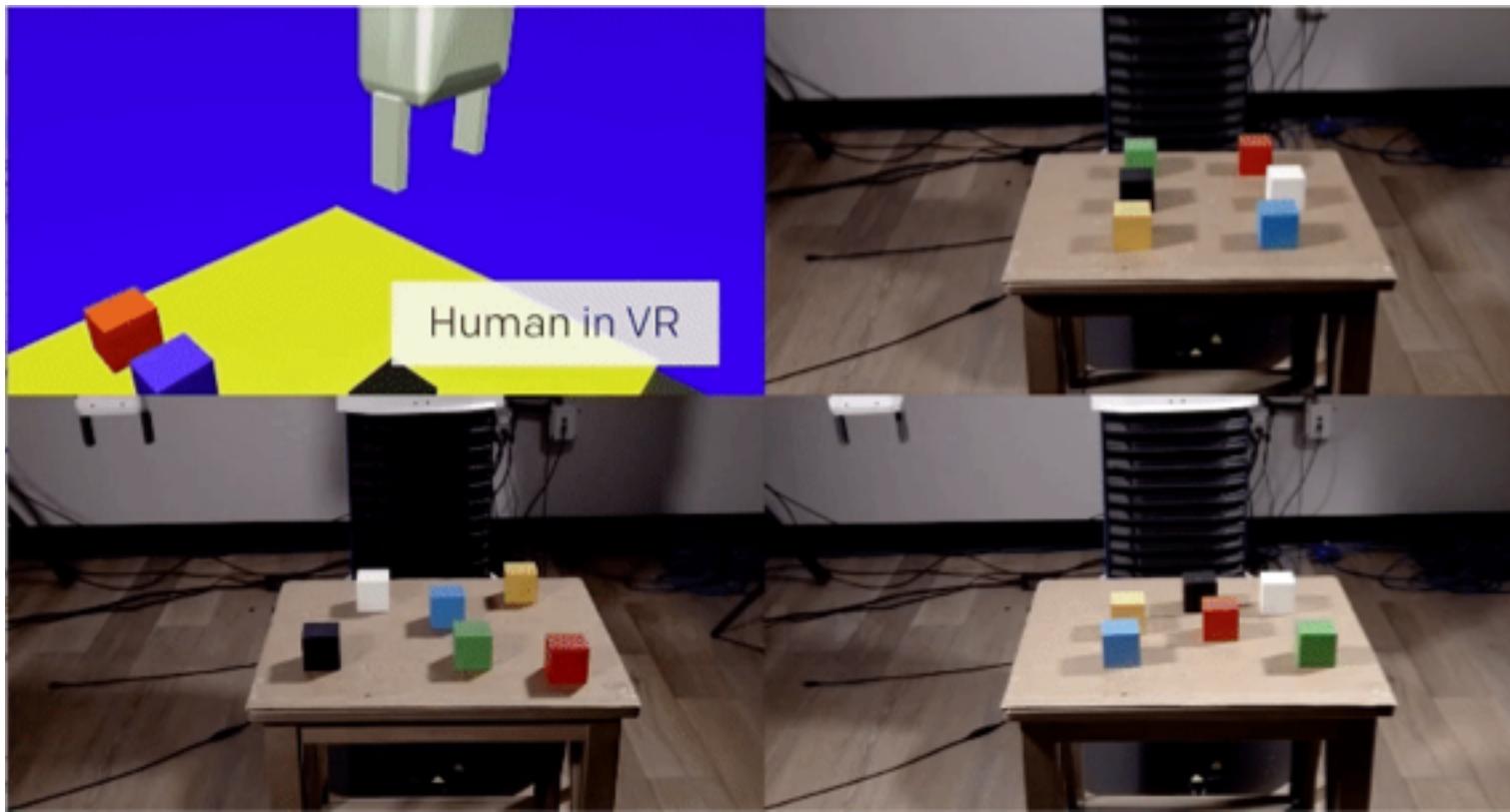
Test in the real world



One Shot Visual Imitation Learning - Transfer

Tobin et al., IROS '17

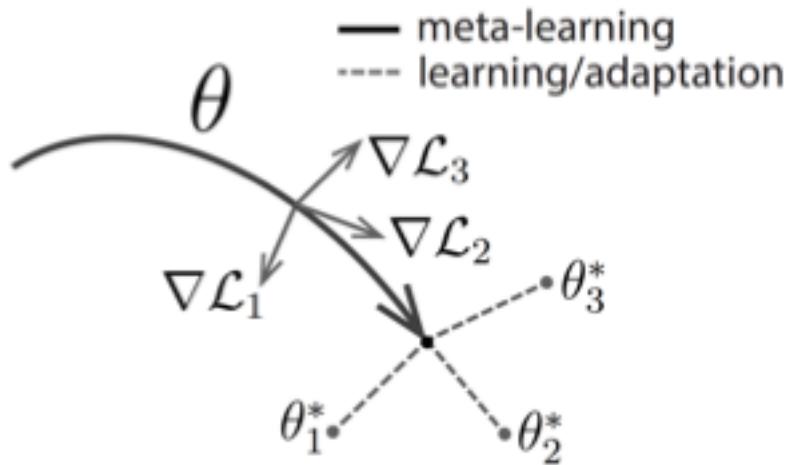
Training entirely in simulation, expert demonstrations through VR



Visual Meta-Imitation Learning

Meta Learning (MAML)

Finn et al., ICML '17



Meta-Imitation Training

Finn et al., CoRL '17

1. Sample task $\mathcal{T}_i \sim p(\mathcal{T})$
2. Using one demonstration for \mathcal{T}_i
$$\theta_i \leftarrow \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(\pi_{\theta})$$
3. Using another demonstration for \mathcal{T}_i
$$\theta \leftarrow \theta - \beta \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(\pi_{\theta_i})$$



Fast adaptation update

Meta-objective update

One Shot Imitation from Watching Videos

Yu et al., RSS '18



pick up a novel object and place it into a previously unseen bowl

Repeated Inverse Reinforcement Learning

Amin et al., NIPS '17

Multiple tasks from multiple environments $\{E_t, t = 1, 2, \dots\}$ sequentially

Unknown reward function θ_*  *latent human preferences*

Human behaves optimally wrt θ_* in each E_t  *External*

Goal: learn θ_* from demonstrations from as few tasks as possible and generalize to new ones, under two settings:

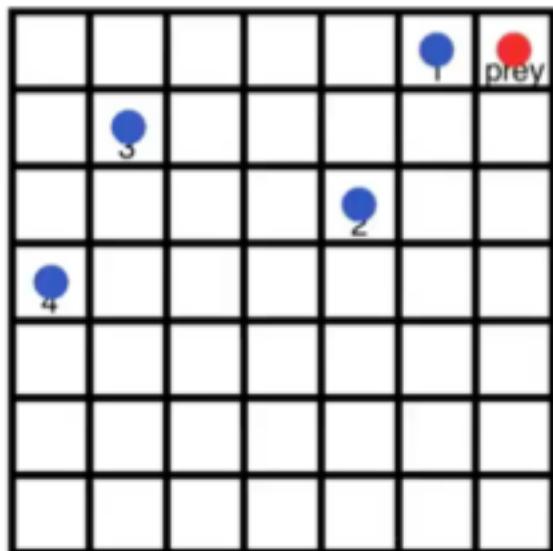
1. Agent chooses the task actively
2. Nature chooses the task adversarially

Multi-Agent / Multi-Modal Imitation Learning

Coordinated Multi-Agent Imitation Learning

Le et al., ICML '17

Coordination via implicit, latent roles information



English Premier League
2012-2013

Match date: 04/05/2013

Coordinated Multi-Agent Imitation Learning

Le et al., ICML '17

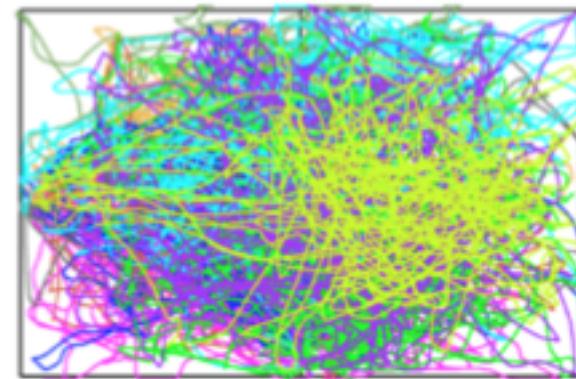
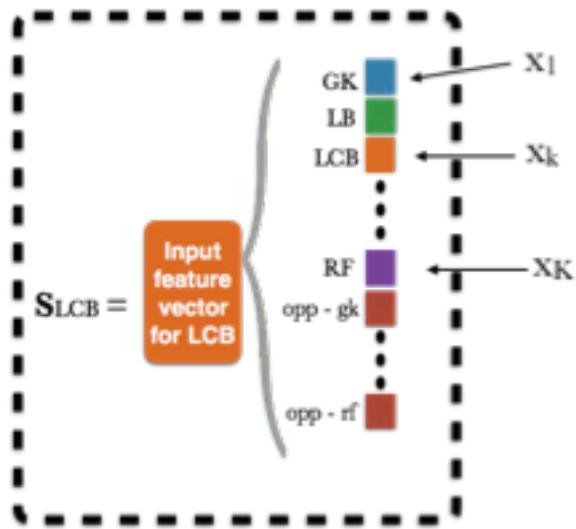
Learning needs semantically consistent input representation

Want:

$$s_k = \phi(x_1, \dots, x_k, \dots, x_K, a_1, \dots, a_K)$$



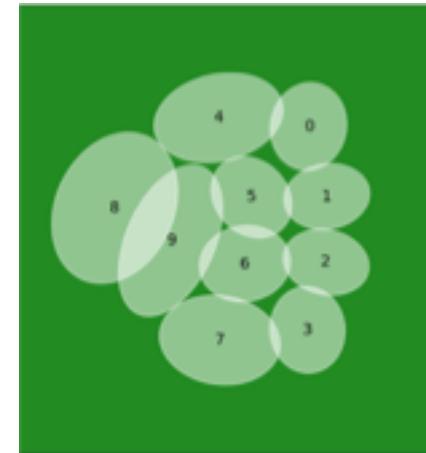
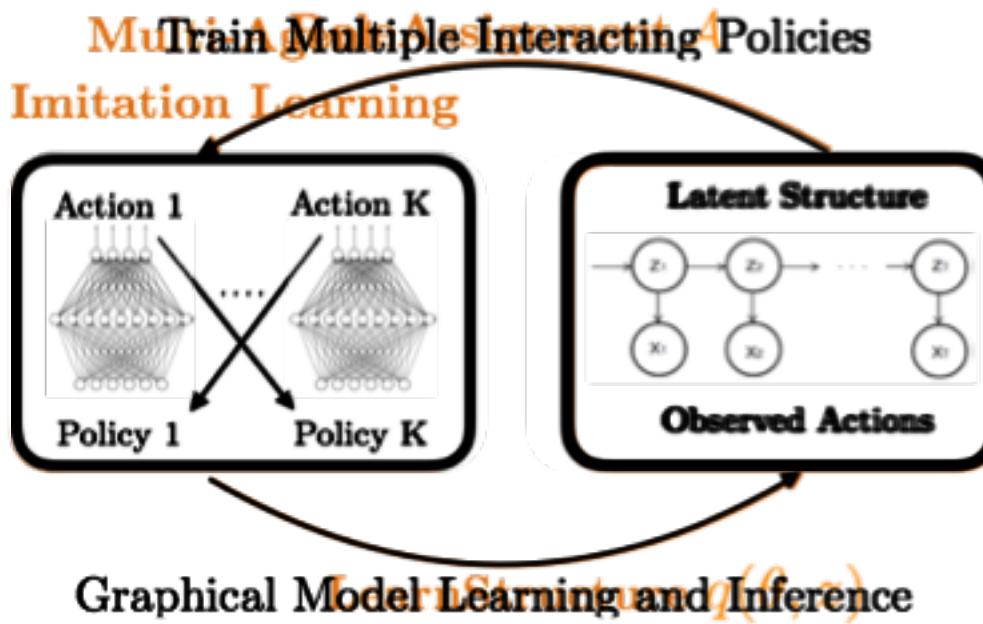
Reality:



Coordinated Multi-Agent Imitation Learning

Le et al., ICML '17

Imitation Learning + Graphical Model Learning and Inference



Coordinated Multi-Agent Imitation Learning

Le et al., ICML '17



English Premier League
2012-2013

Match date: 04/05/2013

Un-coordinated



English Premier League
2012-2013

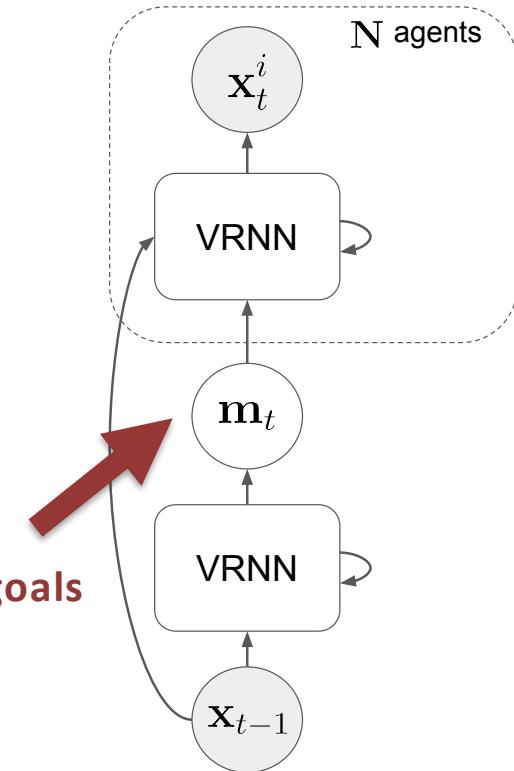
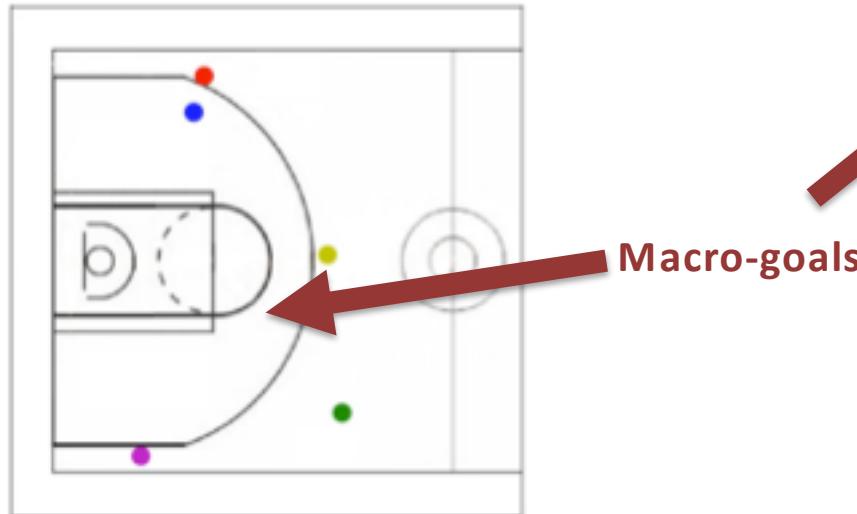
Match date: 04/05/2013

Coordinated

Generative Multi-agent Behavioral Cloning

Zhan et al., arXiv '18

- **Generative Multi-agent Imitation Learning**
 - No single “correct” action
- **Hierarchical**
 - Make predictions at multiple resolutions



Multi-modal Imitation Learning - InfoGAIL

infoGAIL: Interpretable Imitation Learning from Visual Demonstrations – Li et al., NIPS ‘17



Passing Left

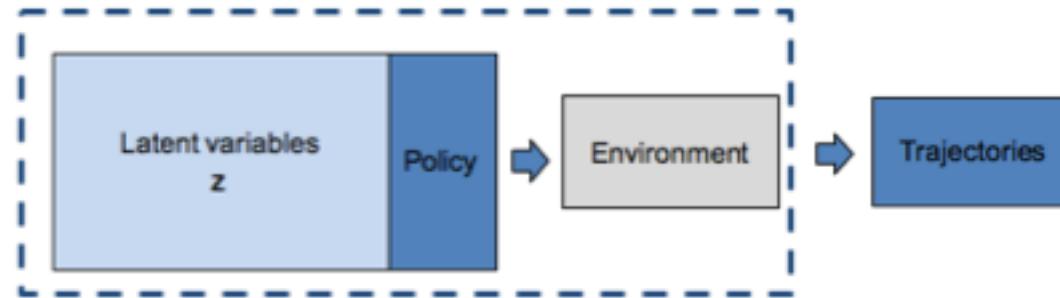


Passing Right

Thanks to materials from Stefano Ermon

Multi-modal Imitation Learning - InfoGAIL

Li et al., NIPS '17



Pass left ($z=0$)



Pass right ($z=1$)

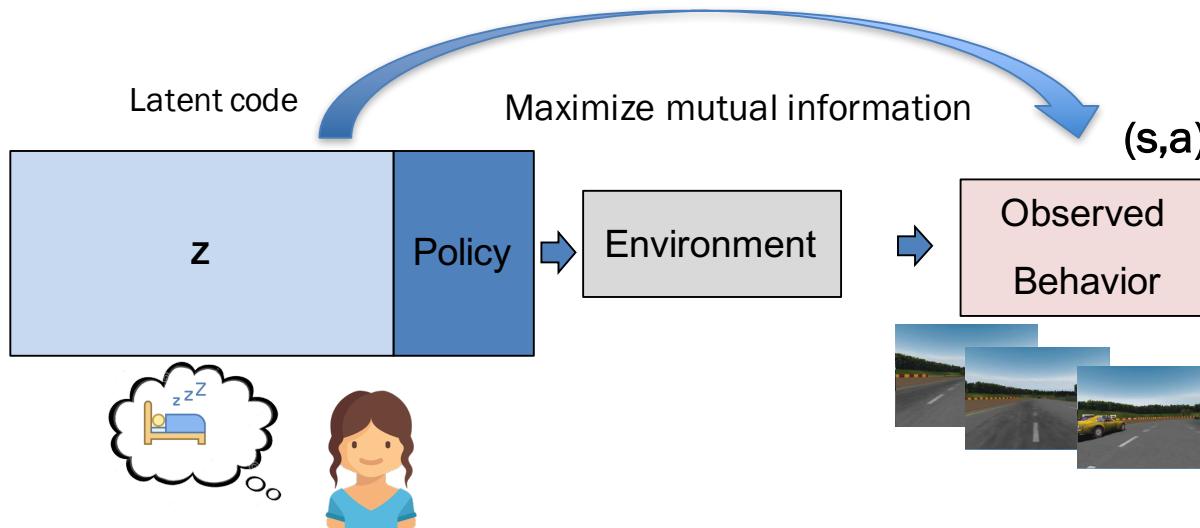


Multi-modal Imitation Learning - InfoGAIL

Li et al., NIPS '17

Further regularization: maximize the mutual information between z and (s, a)

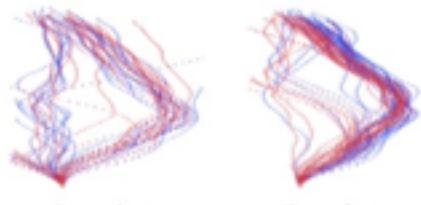
$$\begin{aligned} L_I(\pi, P) &= \mathbb{E}_{z \sim p(z), a \sim \pi(\cdot | s, z)} [\log P(z | s, a)] + H(z) \\ &\leq I(z; s, a) \end{aligned}$$



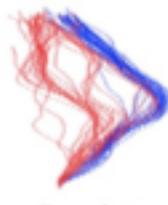
Multi-modal Imitation Learning - InfoGAIL

Li et al., NIPS '17

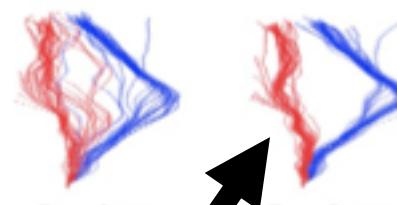
Training using Wasserstein GAN



Epoch 5



Epoch 9



Epoch 13

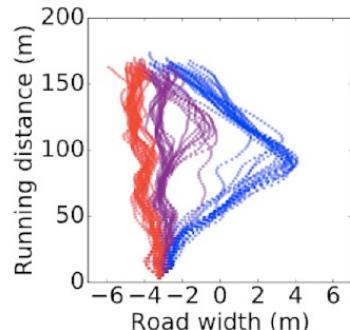
Passing Right

Passing Left

Passing Right

Mode Generation

Z=0
Z=0.5
Z=1



Thanks to materials from Stefano Ermon

See also...

- **Multi-Modal Imitation Learning from Unstructured Demonstrations using Generative Adversarial Nets -**
Hausman et al., NIPS '17
- **Multi-agent Generative Adversarial Imitation Learning**
- Song et al., arXiv '18

Inverse Reinforcement “Teaching”

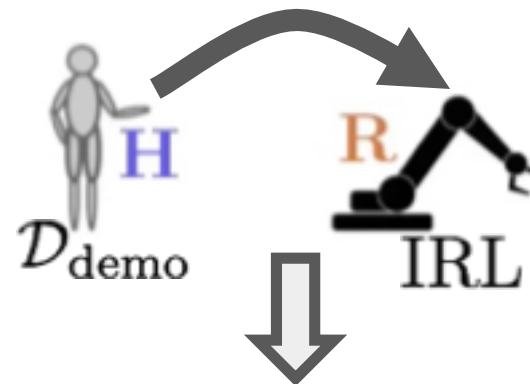
Cooperative Inverse Reinforcement Learning

Hadfield-Menell et al., NIPS '16

How should humans demonstrate the tasks?



Inverse RL: demonstrated behavior is (near) optimal as it accomplishes the task efficiently

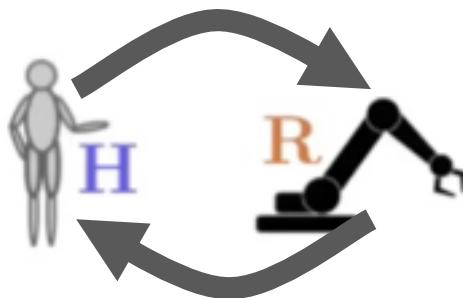


precludes useful teaching behaviors

Cooperative Inverse Reinforcement Learning

Hadfield-Menell et al., NIPS '16

- Human and robot play cooperative games



- Two players:

- **H** observes the actual reward signal; **R** only knows a prior distribution on reward functions

H may choose to accept less reward on a particular action in order to convey more information to **R**

See also...

- **An Efficient, Generalized Bellman Update for Cooperative Inverse Reinforcement Learning** - Malik et al., ICML 2018

- More efficient POMDP solver
- (Partially) relax assumption of human rationality

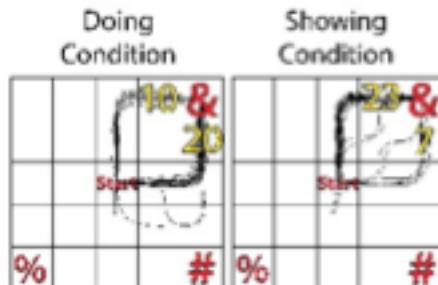
Showing vs Doing: Teaching by Demonstration

Ho et al., NIPS '16

Pedagogical Inverse Reinforcement Learning

*Teaching depends on expert's inferences about
a learner's inferences about Demonstration*

Human Experiment:



Hierarchical Learning / Other Types of Expert Feedback

Hierarchical Imitation and Reinforcement Learning

Le et al., ICML '18

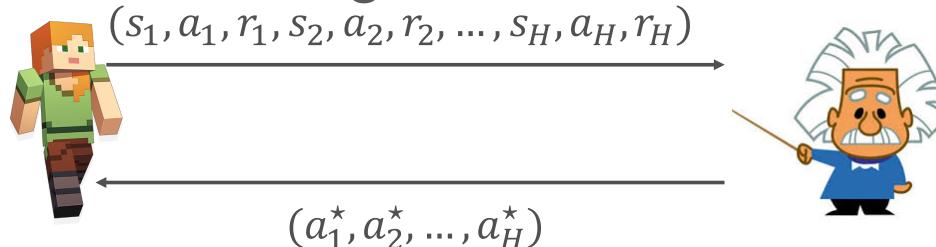
How can we most effectively leverage teacher's feedback?

Hierarchical feedback is more natural



Restroom: End of hall → Left →
Corridor → Left → Go to restroom

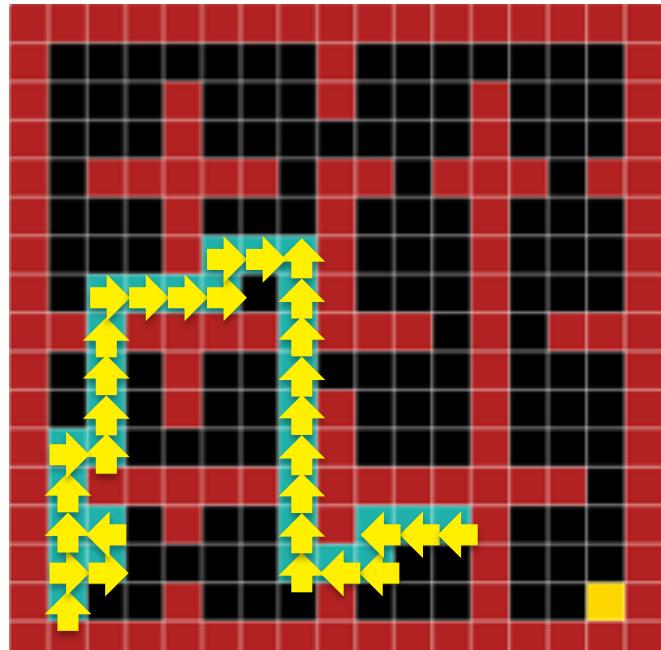
≠ Flat Imitation Learning / Reinforcement Learning



Hierarchical Imitation and Reinforcement Learning

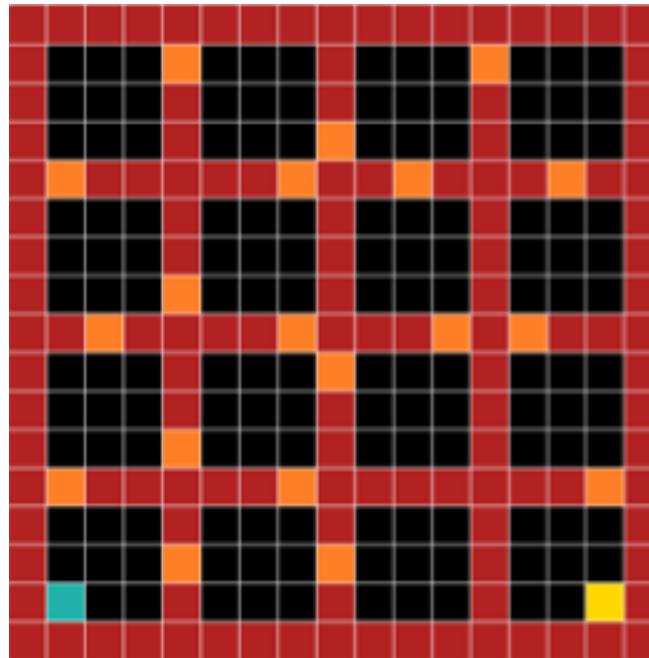
Le et al., ICML '18

(Flat) Imitation Learning



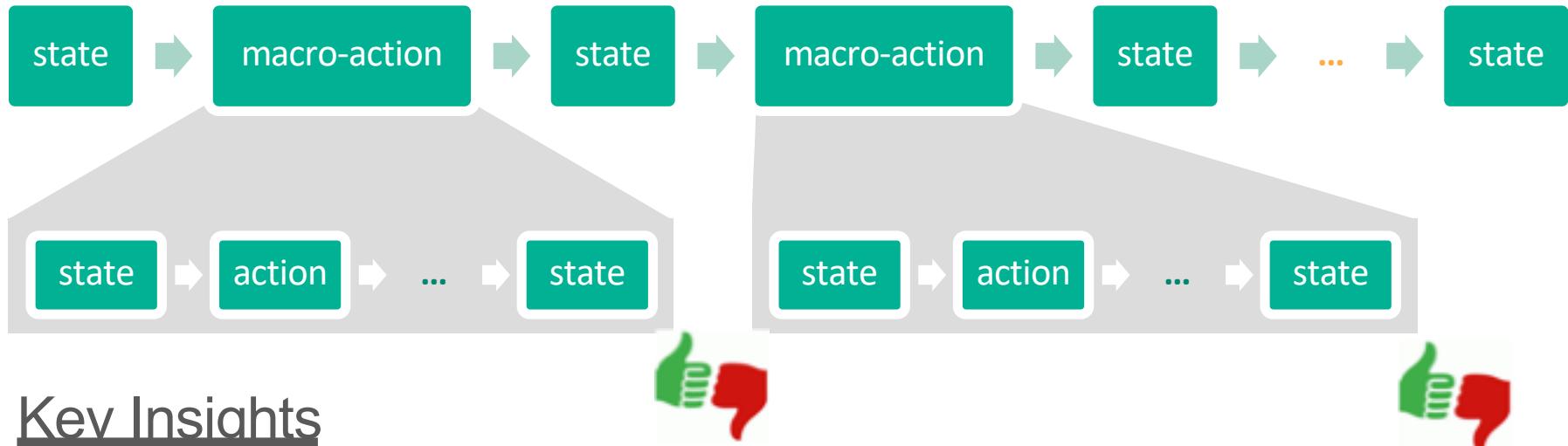
Feedback effort proportional to
task horizon

Hierarchical Imitation Learning



Hierarchical Imitation and Reinforcement Learning

Le et al., ICML '18



Key Insights

Verifying is cheaper than labeling low-level actions

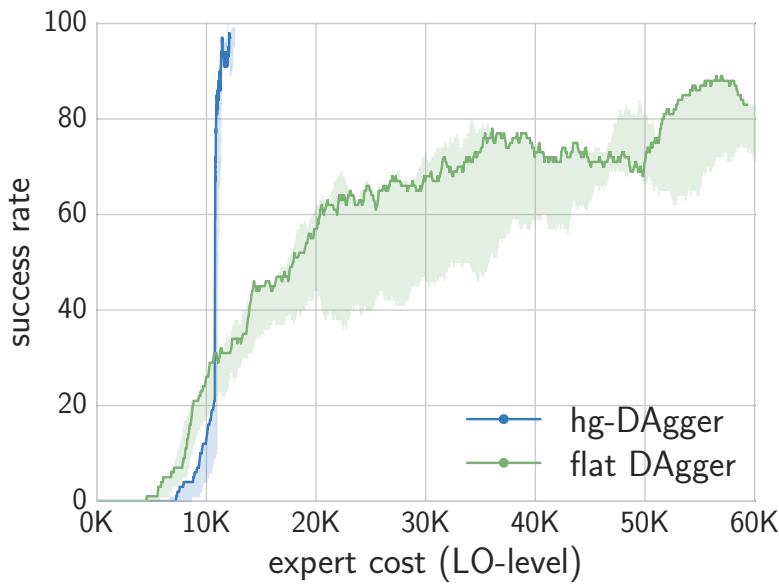
Expert provides **high-level feedback**
and zooms in to **low-level only when needed**

Weak feedback

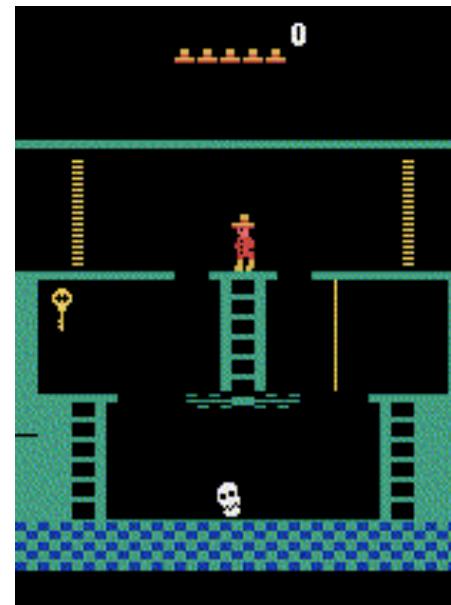
Hierarchical Imitation and Reinforcement Learning

Le et al., ICML '18

Hierarchical Imitation Learning



Hybrid Imitation-Reinforcement



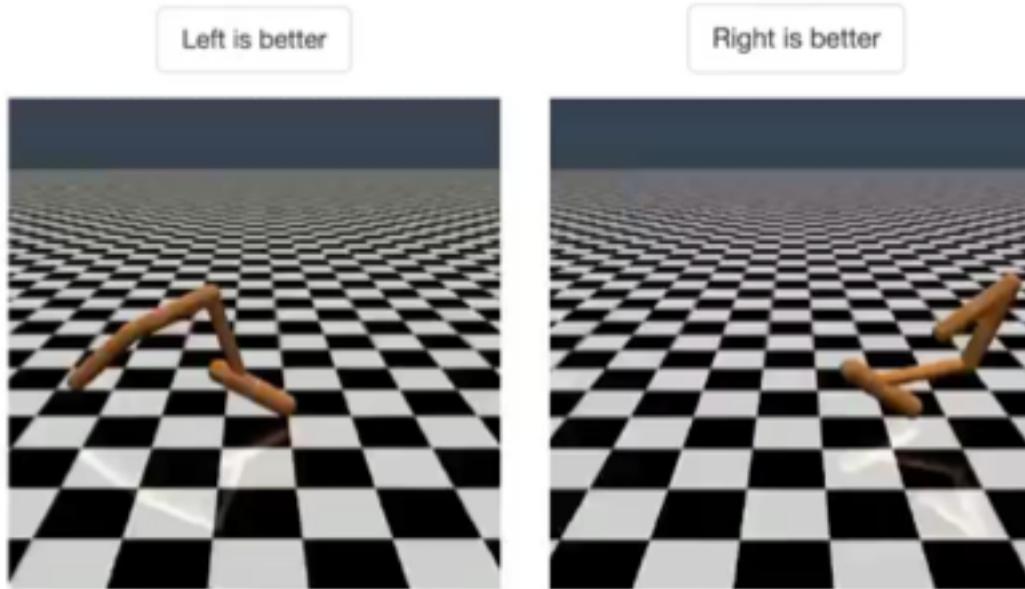
Significantly less expert-costly
than flat imitation learning

Significantly faster learning
than pure RL

Learning from Human Preferences

Christiano et al., NIPS '17

Preference feedback is more natural to specify than exact numeric rewards



Learning from Human Preferences

Christiano et al., NIPS '17



Assume Bradley-Terry preference model: $P [\tau_1 \succ \tau_2] = \frac{e^{r(\tau_1)}}{e^{r(\tau_1)} + e^{r(\tau_2)}}$

Minimize cross-entropy loss:

$$\text{loss}(r) = - \sum_{(\tau_1, \tau_2, \text{pref}) \in \mathcal{D}} \text{pref}(\tau_1) \log P[\tau_1 \succ \tau_2] + \text{pref}(\tau_2) \log P[\tau_2 \succ \tau_1]$$

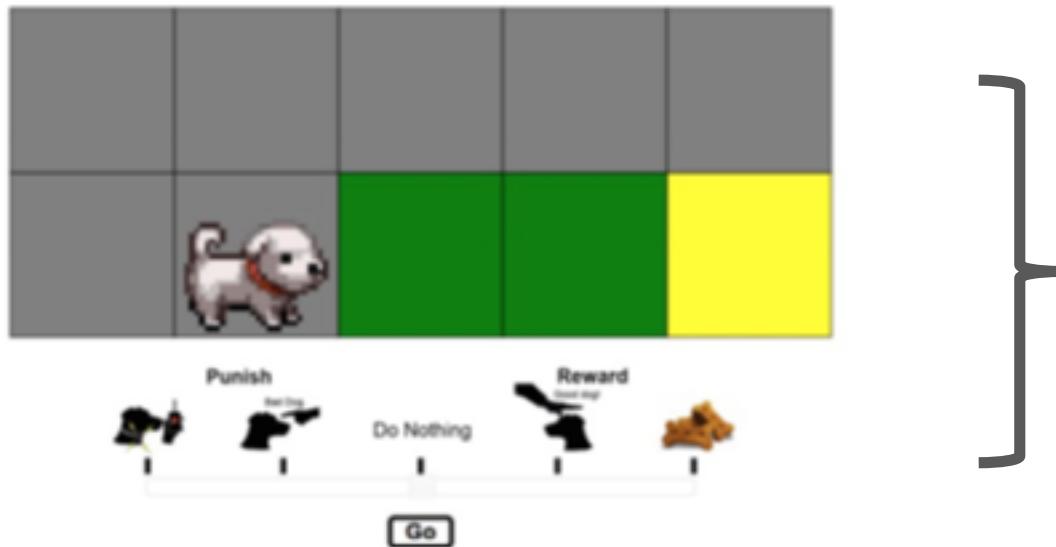
See also previous work...

- **Programming by Feedback** - Akrour et al., ICML 2014
- **A Bayesian Approach for Policy Learning from Trajectory Preference Queries** - Wilson et al., NIPS 2012
 - Employed similar feedback elicitation mechanism to Christiano et al.
- **Learning Trajectory Preferences for Manipulators via Iterative Improvement** - Jain et al., NIPS 2013
 - Co-active feedback (feedback is only an incremental improvement, not necessarily the optimal result)

Interactive Learning from Policy-Dependent Human Feedback (COACH)

MacGlashan et al., ICML '17

How Do Humans Train Animal?



Humans tend to give
policy-dependent
feedback

training interface shown to AMT users

Reward / punishment for "good" dog / "alright" dog / "bad" dog

Interactive Learning from Policy-Dependent Human Feedback (COACH)

MacGlashan et al., ICML '17

- *Policy-dependent feedback*
- *Differential feedback*
- *Diminishing return*



Advantage function is good feedback model

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$$



Treat human feedback as approximation to advantage function

Interactive Learning from Policy-Dependent Human Feedback (COACH)

MacGlashan et al., ICML '17

- Recall Policy Gradient Theorem

$$\nabla_{\theta} V(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) A^{\pi_{\theta}}(s, a)]$$

- Update Rule:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \pi_{\theta}(s, a) \frac{\text{human reward}}{\pi_{\theta}(s, a)}$$

discretized rewards

where $\mathbb{E}[\text{human reward}] = A^{\pi}(s, a)$

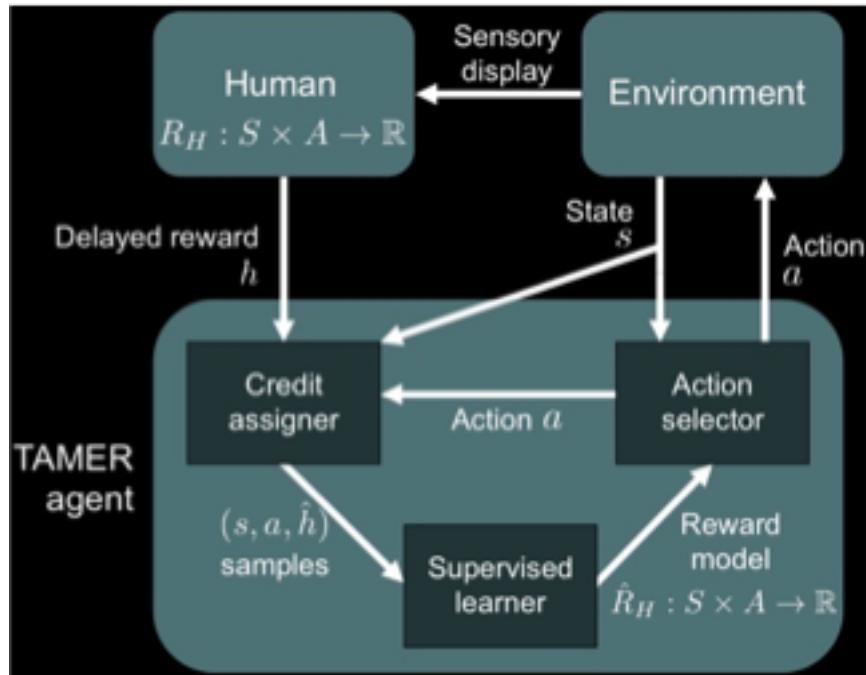
Interactive Learning from Policy-Dependent Human Feedback (COACH)

MacGlashan et al., ICML '17



See also...

- Interactively shaping agents via human reinforcement: the TAMER framework - Knox & Stone, K-CAP '09



<http://www.cs.utexas.edu/~bradknox/TAMER.htm>

Summary

1. Broad research area with many interesting questions
2. Practical approaches to sequential decision making
 - Especially vs RL: sample complexity, cost of experimentation, simulator availability, generalization, etc.
3. Diverse range of applications



Imitation Learning

ICML 2018 Tutorial

Q & A

Yisong Yue



Hoang M. Le



yyue@caltech.edu

hmle@caltech.edu



@YisongYue

@HoangMinhLe



visongyue.com

hoangle.info