

# INVERSE REINFORCEMENT

LEARNING:

FROM

MAXIMUM MARGIN

TO

MAXIMUM ENTROPY



- SANJIBAN CHOUDHURY

# FORMALIZING INVERSE REINFORCEMENT LEARNING (IRL)

# OFF-TERRAIN NAVIGATION

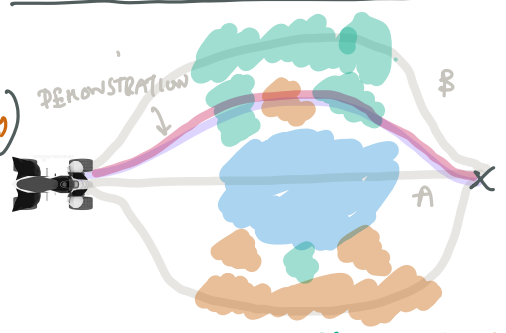
MDP:  $\langle S, A, T, C \rangle$   
 known known known UNKNOWN!

$$C_{\theta}(s,a) := w^1 \mathbb{1}(se\ WATER) + w^2 \mathbb{1}(se\ GRASS) + w^3 \mathbb{1}(se\ ROCKS)$$

$$:= w^1 f^1(s,a) + w^2 f^2(s,a) + w^3 f^3(s,a)$$

(feature map)

$$:= W^T f(s,a)$$



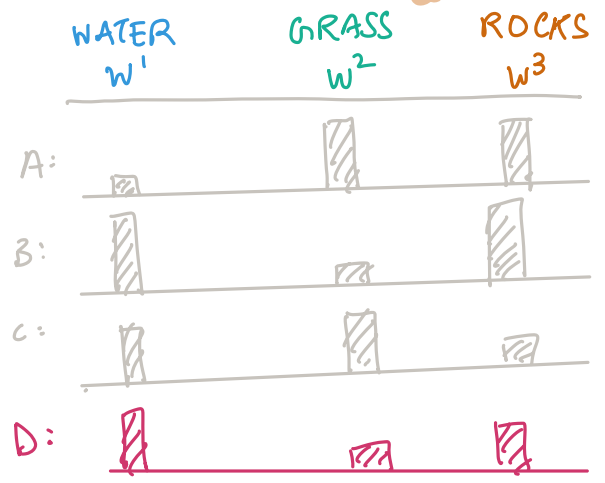
$C_{\theta}(s,a) \mapsto$  
 PLANNER  
 $\underset{a_1, \dots, a_T}{\operatorname{argmin}} \sum_{t=1}^T C_{\theta}(s_t, a_t)$ 
  $\rightarrow$  TRAJECTORY  $\Sigma = s_1, a_1, \dots, s_T, a_T$

$$C_{\theta}(\Sigma) = W^T f(\Sigma)$$

$$= \sum_{(s_t, a_t) \in \Sigma} w^T f(s_t, a_t)$$

## OBJECTIVE

Given (optimal) demonstration  $\Sigma^* = (s_1, a_1, \dots, s_T, a_T)$   
 find  $C_{\theta}(s,a)$  that generates it (find  $w$ )  
 Known as Inverse Optimal Control (IOC), IRL, ...





# JUST A (VERY VERY LARGE) MULTI-CLASS CLASSIFICATION PROBLEM?

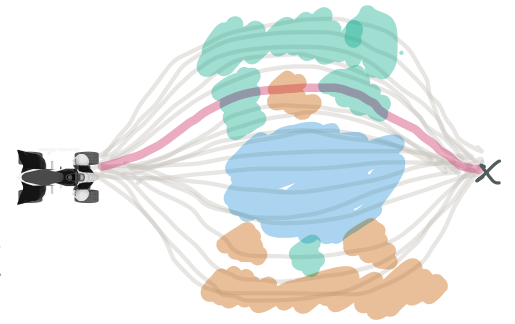
Given  $N$  datapoints  $\{ (\underbrace{\Phi_i}_{\text{MDP}}, \underbrace{\Sigma_i^h}_{\text{HUMAN TRAJ}}) \}_{i=1}^N$

Find  $\theta$

s.t.  $C_\theta(\Sigma_i^h, \Phi_i) \leq C_\theta(\Sigma, \Phi_i) \forall \Sigma, \forall i$

COST OF HUMAN TRAJ  
IN MDP

COST OF ALL TRajs  
IN MDP



( $i^{\text{th}}$  datapoint)

EVERY TRAJ IS A CLASS  $\rightarrow$  A VERY LARGE CLASSIFICATION PROB

LINEAR SETTING: EVERY MDP  $\Phi_i$  INDUCES A FEATURE MAP  $f_i$

$$C_\theta(\Sigma, \Phi_i) = w^T f_i(\Sigma) = \sum_{(s,a) \in \Sigma} (w^1 f_i^1(s,a) + w^2 f_i^2(s,a) + \dots)$$

Find  $w$

s.t.  $w^T f_i(\Sigma_i^h) \leq w^T f_i(\Sigma) \forall \Sigma, \forall i$

$$w^T f_i(z_i^h) \leq w^T f_i(z) \quad \forall z$$

**FAIL**

$$\rightarrow \begin{cases} 0 & \text{if } z = z^h \\ 1 & \text{else} \end{cases}$$

$$w^T f_i(z_i^h) \leq w^T f_i(z) - \gamma_i(z) \quad \forall z$$

MARGIN

$$\leq \min_z [w^T f_i(z) - \gamma_i(z)]$$

**FAIL**

$$\min_w \|w\|^2 + \frac{1}{N} \sum_{i=1}^N \eta_i$$

$$\text{s.t. } w^T f_i(z_i^h) \leq \min_z [w^T f_i(z) - \gamma_i(z)] + \eta_i$$

(SLACK)

?  $w=0$  satisfies this opt  
TRIVIAL SOLUTIONS!

⚠ Exponential # of constraints ...

?? let  $w'$  be a weight s.t

$$w' f(z^h) \leq w' f(z) - \epsilon$$

(TRIVIAL)

Setting  $w = w' \times 10^{1000} \dots$

solves the problem!

(MULTIPLE SOLUTIONS!)

### Maximum Margin Planning

Nathan D. Ratliff  
 J. Andrew Bagnell  
 Robotics Institute, Carnegie Mellon University, Pittsburgh, PA. 15213 USA  
 Martin A. Zinkevich  
 Department of Computing Science, University of Alberta, Edmonton, AB T6G 2E

# EXTEND TO NON-LINEAR COST FUNCTIONS [LEARCH]

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \left( C_{\theta}(\xi_i^h, \phi_i) - \min_{\xi} [C_{\theta}(\xi, \phi_i) - r_i(\xi)] \right) + R(\theta)$$

(HUMAN COST)                      (PLANNER COST)                      (REG)

**Learning to Search:  
Functional Gradient Techniques  
for Imitation Learning**

Nathan D. Ratliff  
Robotics Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213  
ndr@ri.cmu.edu

David Silver  
Robotics Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213  
dsilver@ri.cmu.edu

ACRIB: 3 ways to do this (1) Boosting (2) Kernel gradient descent (3) Parametric functional gradient descent (e.g. deep learning)

Alg

for  $i=1 \dots N$  # loop over datapoints ( $\phi_i, \xi_i^h$ )

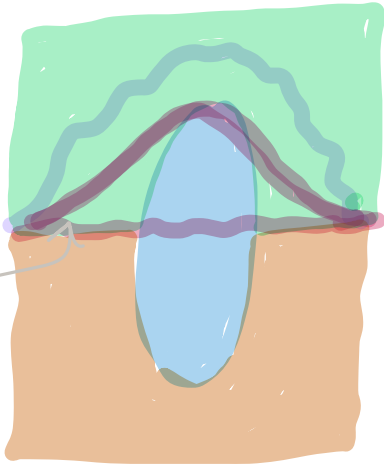
$\xi_i^* = \min_{\xi} [C_{\theta}(\xi, \phi_i) - r_i(\xi)]$  # Call planner to find optimal path

$\theta^+ = \theta - \alpha [\nabla_{\theta} C_{\theta}(\xi_i^h, \phi_i) - \nabla_{\theta} C_{\theta}(\xi_i^*, \theta)] + \nabla_{\theta} R(\theta)$





# PROBLEM: SUBOPTIMAL DEMONSTRATIONS



UNREALIZABLE EXPERT:  
 ANY  $C_0$  that results  
 in expert traj being optimal  
 $\rightarrow$  GRADIENT keep running  
 away!  
 $\rightarrow$  NO convergence!

INSTEAD OF:

FIND COST s.t  
 EXPERT COST  $\leq$  OPTIMAL TRAJ COST

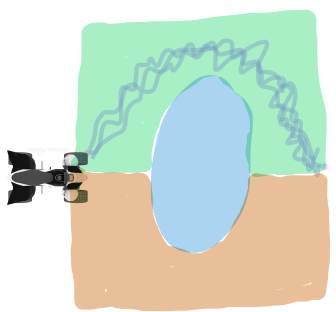
## PROBABILISTIC

FIND DISTRIBUTION over traj  $P(\xi)$  s.t  
 EXPECTED EXPERT COST = EXPECTED TRAJ COST | FOR ALL COST FUNCTIONS!

PARADIGM 2: PROBABILISTIC MODELING VIA MOMENT MATCHING

Let's return to linear case, i.e.  $C(\xi, \theta) = \mathbf{w}^T f(\xi, \theta) = [\omega^1 \omega^2 \dots] \begin{bmatrix} f^1(\xi, \theta) \\ \vdots \\ f^k(\xi, \theta) \end{bmatrix}$

Assume: Human traj distribute  $P(\xi^h)$



Find  $P(\xi)$

$$\max - \int P(\xi) \log P(\xi) d\xi \quad (\text{ENTROPY})$$

s.t.  $E_{\xi \sim P(\xi^h)} f^k(\xi^h) = E_{\xi \sim P(\xi)} f^k(\xi) \quad (\text{MOMENT MATCHING})$   
 $\forall k$

... PROBLEM ...



Ex. Avg time expert spends on grass  $f^1$  = Avg time learner spends on grass } ANY  
 " on water = " on water } LINEAR  
 " on rocks = " on rocks } COMBINATION!  
 (If all moments matched, EXPERT = LEARNER)  
 COST COST

# PRINCIPLE OF MAXIMUM ENTROPY

## Information Theory and Statistical Mechanics

E. T. JAYNES

Department of Physics, Stanford University, Stanford, California

(Received September 4, 1956; revised manuscript received March 4, 1957)

(Pick a distribution that is least committal)

$$\max_{\mathcal{Z}} - \sum P(\mathcal{Z}) \log P(\mathcal{Z})$$

$$\text{s.t.} \quad \sum_{\mathcal{Z}} P(\mathcal{Z}) f^k(\mathcal{Z}) = F^k$$

(EXPECTED FEATURE COUNT)

(EXPERT FEATURE COUNT)

Eg Avg limu on water  $F^1 = 0.0$ , Avg limu on gas  $F^2 = 0.8$ , Avg limu on rocks  $F^3 = 0.2$

Sketch of solution: First write out LAGRANGIAN

$$\mathcal{L}(\dots) = - \sum P(\mathcal{Z}) \log P(\mathcal{Z}) - \sum_k \lambda^k \left( \sum_{\mathcal{Z}} P(\mathcal{Z}) f^k(\mathcal{Z}) - F^k \right)$$

Take gradient & set to 0 and solve! Solution:

$$P_{\lambda}(\mathcal{Z}) = \frac{1}{\mathcal{Z}} \exp \left( - \sum \lambda^k f^k(\mathcal{Z}) \right)$$

THINK OF LAGRANGE MULTIPLIER  $\lambda^k$  as being weights  $w^k$ !  
 $P(\mathcal{Z}) \propto \exp(-\text{COST}) \Rightarrow$  LOW COST near HIGH PROB!

FINAL step: Plug in solution to optimization to solve  $\lambda^k$

$$\max_{\lambda^k} \frac{1}{N} \sum_{i=1}^N \log P_{\lambda}(\mathcal{Z}_i^k)$$

↳ MAXIMIZE LOG LIK OF EXPERT TRAJ.  $\rightarrow$  GENERATIVE MODELING!





# MAXIMUM ENTROPY INVERSE OPTIMAL CONTROL

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N -\log P_{\theta}(\Sigma_i^h | \phi) \text{ where}$$

$$P_{\theta}(\Sigma | \phi) = \frac{1}{Z(\theta, \phi)} \exp(-C_{\theta}(\Sigma, \phi))$$

(HOPEL)

## Maximum Entropy Inverse Reinforcement Learning

Brian D. Ziebart, Andrew Maas, J. Andrew Bagnell, and Anind K. Dey  
School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213  
bziebart@cs.cmu.edu, amaas@andrew.cmu.edu, dbagnell@ri.cmu.edu, anind@cs.cmu.edu

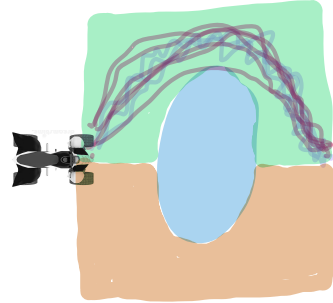
Plug in MODEL in objective to get.

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N C_{\theta}(\Sigma_i^h, \phi) + \log Z(\theta, \phi)$$

$$\int \frac{\exp(-C_{\theta}(\Sigma, \phi))}{Z} d\Sigma$$

EXPERT COST

SOFT MIN OF LEARNER COST



CONVERGES!!!

### ALGORITHM

- for  $i = 1 \dots N$
- Sample  $\Sigma^i \sim \exp(-C_{\theta}(\Sigma, \phi))$  (# SOFT MIN (Sample low cost traj) (MCMC, Laplace approx))

$$\theta^+ \leftarrow \theta - \alpha \left[ \nabla_{\theta} C_{\theta}(\Sigma^i, \phi) - \mathbb{E}_{\Sigma \sim p} \nabla_{\theta} C_{\theta}(\Sigma, \phi) \right]$$

# MAX-MARGIN PLANNER

$$\theta^+ \leftarrow \theta - \alpha \left[ \nabla_{\theta} C_{\theta}(\Sigma^h; \Phi_i) - \nabla_{\theta} C_{\theta}(\Sigma^*, \Phi_i) \right]$$

HUMAN
OPTIMAL COST TRAJ

# MAXIMUM ENTROPY INVERSE OPT CONTROL

$$\theta^+ \leftarrow \theta - \alpha \left[ \nabla_{\theta} C_{\theta}(\Sigma^h; \Phi) - \mathbb{E}_{\Sigma \sim \text{sample}} \nabla_{\theta} C_{\theta}(\Sigma; \Phi_i) \right]$$

(HUMAN)
SAMPLE LOW COST TRAJ



TIE-BREAKING: MAX MARGIN

TIE-BREAKING: MAX ENTROPY.

REGULARIZATION ON COST (L2)

ENTROPIC REGULARIZATION ON PLANNER

OPTIMAL PLANNER (A<sub>eff</sub><sup>\*</sup>)

"SOFT" PLANNER (SAC e.g.)

BOTH SOLVING THE SAME GAME

NO REGRET  
+  
BEST RESPONSE

max	min	$C_{\theta}(\Sigma^h) - C_{\theta}(\Sigma)$
$C_{\theta}$	$\Sigma$	

NO REGRET  
+  
ENTROPY REG BEST RESPONSE

# KEY CHALLENGE: LEARN COST FUNCTION THAT EXPLAINS EXPERT DEMONSTRATIONS

\* METHOD 1: Recover costs that make expert optimal

↳ Maximum margin planning.

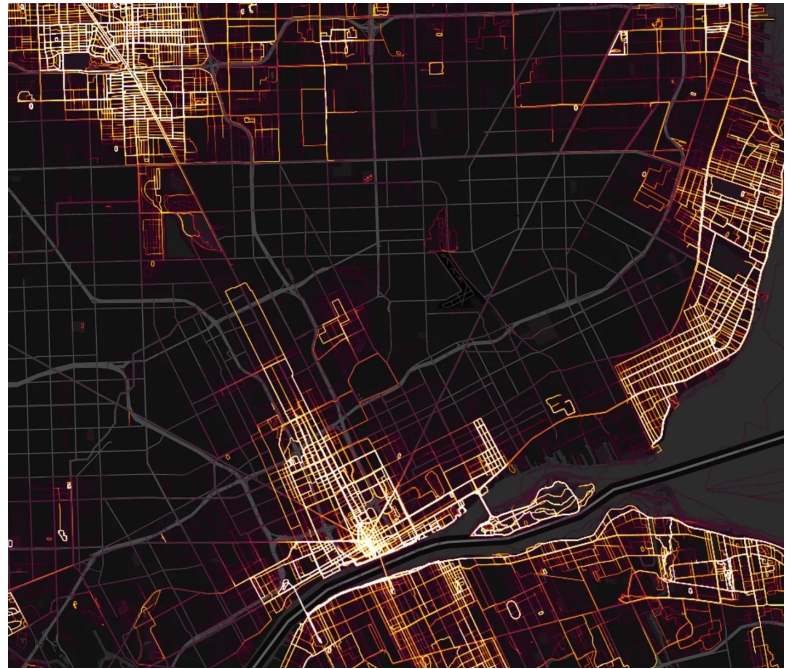
↳ PROBLEM: No suboptimal expert

\* METHOD 2: Learn a distribution over traj that matches all cost function moments

↳ Maximum entropy

↳ Can deal with any expert dist!

\* Both methods are solving the same game in different ways!



It's all a game between

DISCRIMINATOR  
(COST FUNCTION / ADVANTAGE  
function)

&

GENERATOR  
(PLANNER / POLICE)

Hey everyone,

Welcome to six<sup>th</sup> lecture in our series on imitation learning.

I am Sayan Chaudhry, a research scientist at Arora  
& San to be assistant professor at  
Cornell

Today we'll talk about a class of Imitation Learning methods that we want to try to recover the underlying reward or cost function. We will look at the problem from 2 different perspectives - #1 by treating this as a planning problem and #2 by treating this as a stochastic process. And in the end, we will see that they both converge at a unified game theoretic framework. Let's get started.