# 1) Inverse RL

like imititation learning, we learn from expert demonstrations.

rather thatn learning the expert's policy, IRL tries to learn the _reward_ _function_

1) Imitation via reward fn
   - reward fn. is more succint/transferrable

2) Scientific Inquiry
   - modelling human/animal behavior

3) Multiagent seting - model other agents

## Setting:

$$M = \{ S, A, P, r, H, \gamma \} \qquad (P \text{ known})$$

- Reward function $r(s,a)$ is unknown and signal $r_t$ is unobserved

- observe trajectories from expert w/ optimal policy $\pi^*$

## Basic Idea: Find a reward function which is consistent with the optimality of the expert policy

$$\text{find } r \quad s.t. \quad \underset{s,a \sim d_\gamma^{\pi^*}}{\mathbb{E}[r(s,a)]} \geq \underset{s,a \sim d_M^{\pi}}{\mathbb{E}[r(s,a)]} \quad \forall \; \pi$$

$r: S \times A \to [0,1]$

$$\underset{s \sim M}{\mathbb{E}[V^{\pi^*}(s)]}$$

estimate from expert trajectories

## Problems w/ Formulation:

1) Need to consider all $A^{SH}$ policies

2) Ambiguity: more than one reward function may satisfy $(r=0)$

Reframe: Find a policy that is as good as the expert $\leftarrow$ unknown

Key assumption: $\mathbb{E}[r(s,a)] = \Theta_*^T \phi(s,a)$ $\leftarrow$ known
linear reward wrt features

-ex- $\phi(s,a) = \begin{bmatrix} \mathbb{P}(\text{building}) \\ \mathbb{P}(\text{sidewalk}) \\ \mathbb{P}(\text{road}) \\ \vdots \end{bmatrix}$

$\Theta_*$ weighs negatives (driving on sidewalk, hitting building) & positives (road)

We can write the policy consistency problem:

find $\pi$ s.t. $\underset{s,a \sim d_\gamma^{\pi^*}}{\mathbb{E}}[\phi(s,a)] = \underset{s,a \sim d_\gamma^\pi}{\mathbb{E}}[\phi(s,a)]$

$\pi: S \times \mathcal{R} \to \Delta(S)$

estimate from $\frac{1}{N}\sum_{i=1}^{N} \phi(s_i,a_i)$

To solve the ambiguity problem, we will use the "maximum entropy principle"

The Max Entropy IRL method:

$\underset{\pi}{\max}$ Entropy $\pi$ $\leftarrow$
s.t. $\underset{d_\gamma^{\pi^*}}{\mathbb{E}}[\phi(s,a)] = \underset{d_\gamma^\pi}{\mathbb{E}}[\phi(s,a)]$

## 2) Maximum Entropy Principle

Def (Entropy) Distribution $P(x) \in \Delta(X)$
$\text{Ent}(P) = \underset{x \sim P}{\mathbb{E}}[-\log(P(x))] = -\sum_{x \in X} P(x)\log(P(x))$

$P(x) \in [0,1]$, Entropy is positive

ex — deterministic distribution $X = x_0$ w.p. 1

$$P_{x_0}(x) = \mathbb{1}\{x = x_0\}$$

$$\text{Ent}(P_{x_0}) = -\sum_{x \neq x_0} 0 \cdot \overrightarrow{\log(0)}^{0} + -1 \cdot \overrightarrow{\log(1)}^{0}$$

$$= 0$$

ex — uniform distribution over $|X| = N$ elements

$$\text{Ent}(u) = -\sum_{x \in X} \frac{1}{N} \log(1/N) = -\log(1/N)$$

$$= \log N$$

## Max Ent Principle:

"Among consistent distributions,
choose the one w/ the
most uncertainty, ie.
the highest entropy.

w/ cmtraints arising
from observation,
mean, variance

The max-ent IRL approach:

$$\max_{\substack{\pi \\ s \sim d_M^\pi}} \mathbb{E}\left[\text{Ent}(\pi(\cdot | s))\right] = \max_{\pi} - \mathbb{E}_{s \sim d_M^\pi}\left[\mathbb{E}_{a \sim \pi}\left[\log(\pi(a | s))\right]\right]$$

s.t. contraints

$$= \min_{\substack{\pi \\ s, a \sim d_M^\pi}} \mathbb{E}\left[\log(\pi(a | s))\right]$$

$$\min_{\substack{\pi \\ s, a \sim d_M^\pi}} \mathbb{E}\left[\log \pi(a | s)\right]$$

$$\text{s.t. } \mathbb{E}_{s, a \sim d_M^{\pi^*}} \phi(s, a) = \mathbb{E}_{s, a \sim d_M^\pi} \phi(s, a)$$

$$\frac{1}{N}\sum_{i=1}^{N} \phi(s_i^*, a_i^*)$$

# 3) Constrained optimization

Consider the constrained optimization problem:

$$x^* = \arg\min_x \left[ f(x) \text{ s.t. } g(x) = 0 \right]$$



## Lagrange Formulation:

$$\min_{x \in \mathbb{R}^d} \left[ \max_{w \in \mathbb{R}} f(x) + w \cdot g(x) \right]$$
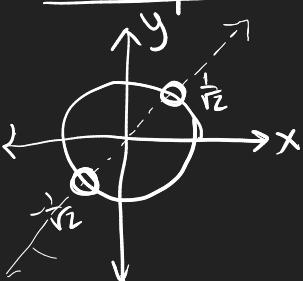
if $g(x) \neq 0$, $w \to \pm\infty$  so inner max is $\infty$

if $g(x) = 0$, inner maximization: $f(x)$

$$\max_w f(x) + w \cdot g(x) = \begin{cases} \infty & g(x) \neq 0 \\ f(x) & g(x) = 0 \end{cases}$$

$$x^* = \arg\min_x \max_w f(x) + w \cdot g(x)$$

## example:



$$\min \quad x+y \qquad \text{s.t.} \quad x^2 + y^2 = 1$$

$$\Downarrow$$

$$\min_{x,y} \max_w \underbrace{x + y + w(x^2 + y^2 - 1)}_{\mathcal{L}(x,y,w)}$$

$$\nabla_x \mathcal{L} = 1 + 2xw \overset{0}{\Longrightarrow} x^* = \frac{-1}{2w_*} -$$

$$\nabla_y \mathcal{L} = 1 + 2yw \overset{0}{\Longrightarrow} y^* = \frac{-1}{2w_*} -$$

$$\rightarrow \nabla_w \mathcal{L} = x^2 + y^2 - 1 \overset{0}{\Longrightarrow} (x^*)^2 + (y^*)^2 = 1 -$$

$$\left(\frac{-1}{2w_*}\right)^2 + \left(\frac{-1}{2w_*}\right)^2 = 1 \Longrightarrow w_* = \pm\sqrt{1/2} = \frac{\sqrt{2}}{2}$$

critical points: $\left(-\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}\right)$ and $\left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}\right)$

# Iterative Procedure

initialize $w_0$

For $t = 0, \dots, T-1$

$$\begin{cases} x_t = \arg\min_x f(x) + w_t \, g(x) & \text{(Best response)} \\ w_{t+1} = w_t + \eta \, g(x_t) & \text{(incremental update)} \end{cases}$$

Return $\bar{x} = \frac{1}{T} \sum_{t=0}^{T-1} x_t$, $\quad \bar{x} \to x^*$ as $T \to \infty$

if $f, g$ are convex