# Lecture 24: Inverse RL

## 1) Inverse RL Motivation & Setting

Like imitation learning, which we studied last week, the setting of inverse RL seeks to learn from expert demonstrations. However, rather than attempt to directly learn the expert's policy, IRL tries to learn the <u>reward</u> <u>function</u>. There may be several motivations for doing so:

1) Scientific Inquiry (modelling human or animal behavior)

our focus → 2) Imitation learning via the reward function (reward function may be the most succinct & transferable information about a task, rather than a policy)

3) modelling other agents in a multiagent setting (adversarial or cooperative)

<u>Setting</u>: Finite Horizon MDP      (P known)

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

Reward function $r(s,a)$ unknown /signal $r_t$ unobserved But observe traces of expert policy which is optimal $\pi^*$.

<u>Basic Idea</u>: Find reward functions which are consistent with expert being optimal:

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

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

estimate from expert trajectories

Two Problems with this formulation:

1) As written, need to consider all possible policies ($A^S$)

2) Ambiguity: more than one reward function may satisfy this description (e.g. $r(s,a) = 0$)

Since our goal is to ultimately use the learned reward function for policy design, we can reframe the consistency as a property of policies:

"find a policy that is at least as good as the expert."

<u>Key assumption</u>: $r(s,a) = \Theta_*^T \phi(s,a)$

The reward function is linear with respect to a known feature mapping

ex $\phi(s,a) = \begin{bmatrix} \mathbb{P}(\text{building}) \\ \mathbb{P}(\text{sidewalk}) \\ \mathbb{P}(\text{road}) \\ \mathbb{P}(\text{person}) \end{bmatrix}$   and $\Theta_*$ encodes that staying on the road is good (positive weight), sidewalk is bad (negative), colliding with person very bad (large negative)

Then we can write the policy consistency problem as:

find $\pi$ s.t. $\underset{s,a \sim d^{\pi^*}_\mu}{\mathbb{E}}\left[\phi(s,a)\right] = \underset{s,a \sim d^{\pi}_\mu}{\mathbb{E}}\left[\phi(s,a)\right]$

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

$\underbrace{s,a \sim d^{\pi^*}_\mu}$ estimate from data with $\frac{1}{N}\sum_{i=1}^{N}\phi(s_i,a_i)$

By the linear cost assumption, this constraint implies that policies acheive the same reward.
    (Exercise: write a short proof of this fact)

However, this does not solve the <u>ambiguity</u> problem: many such policies may satisfy this constraint.

Idea: among all policies satisfying the consistency constraint, choose the one that is the most "uncertain"

We will study the Maximum Entropy IRL method:

$$\underset{\pi}{\max}\ \text{Entropy of } \pi$$

$$\text{s.t.} \quad \underset{s,a \sim d^{\pi^*}_\mu}{\mathbb{E}}\left[\phi(s,a)\right] = \underset{s,a \sim d^{\pi}_\mu}{\mathbb{E}}\left[\phi(s,a)\right]$$

# 2) Maximum Entropy Principle

our choice to choose the consistent policy with the maximum entropy follows from the "maximum entropy principle"

<u>Definition</u> (Entropy):

The entropy of a distribution $P \in \Delta(x)$ is defined as

$$Ent(P) = \mathbb{E}_{x \sim P}\left[-\log(P(x))\right] = -\sum_{x \in \mathcal{X}} P(x) \log(P(x))$$
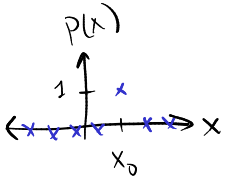
Because $P(x) \in [0,1]$, $Ent(P) \geq 0$.

Lower Entropy means lower uncertainty.

Ex: Deterministic distribution where $x = x_0$ w.p. 1.
$$P_{x_0}(x) = \mathbb{1}\{x = x_0\}.$$



$$Ent(P_{x_0}) = -\sum_{x \in \mathcal{X}} \mathbb{1}\{x=x_0\} \log\left(\mathbb{1}\{x=x_0\}\right)$$
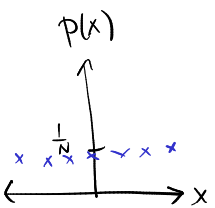
$$= -1 \cdot \log(1) + \sum_{x \neq x_0} 0 \cdot \log(0) = 0$$

Ex: Uniform distribution over $|\mathcal{X}| = N$ elements.
$$U(x) = \frac{1}{N}$$



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

$$= N \cdot \frac{1}{N} \cdot \log(N) = \log(N)$$

Exercise: argue that the uniform distribution has the highest possible entropy.

# Maximum entropy Principle:

Among distributions consistent with constraints (ie. observed data, mean, variance) choose the one with the most uncertainty, ie. highest entropy.

This can be seen as an application of Occam's razor since we are in some sense making the fewest assumptions about the distribution.

ex - Gaussian distribution.

$$\max \ \text{Ent}(P) \quad \text{s.t.} \quad \mathbb{E}_{x \sim P}[x] = \mu \quad \text{and} \quad \mathbb{E}[xx^T - \mu\mu^T] = \Sigma$$

The solution is a Gaussian Distribution
(proof / Derivation out of scope)

The max-ent IRL approach solves:

$$\max_{\pi} \ \mathbb{E}_{s \sim d_\mu^\pi}\left[\text{Ent}(\pi(\cdot|s))\right] \longrightarrow \min_{\pi} \mathbb{E}\left[\log(\pi(a|s))\right]$$

$$\text{s.t.} \quad \mathbb{E}_{s,a \sim d_\mu^{\pi^*}}[\phi(s,a)] = \mathbb{E}_{s,a \sim d_\mu^\pi}[\phi(s,a)]$$

we can simplify

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

# 3) Constrained Optimization

We've considered many optimization algorithms throughout the semester but not many with constraints (exception: trust regions)

Consider the constrained optimization problem:

$$x^* = \arg \min_x f(x) \qquad s.t. \quad g(x) = 0 \qquad \text{(primal)}$$

To solve this problem we consider the Lagrange formulation, which converts it into an unconstrained:

$$\min_x \left[ \max_w f(x) + w g(x) \right] \qquad \text{(Lagrange)}$$

Now if $g(x) \neq 0$ (i.e. $x$ is infeasible for primal), then
$$\max_w f(x) + w \cdot g(x) = \infty \quad (w^* = \infty).$$

And if $g(x) = 0$ (i.e. $x$ is feasible) then
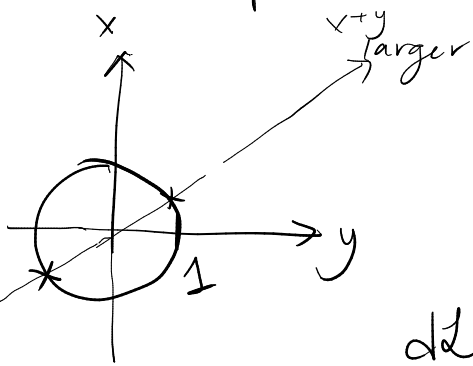$$\max_w f(x) + w g(x) = f(x).$$

Then the Lagrange formulation encodes

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

since $\infty$ is undesirable with respect to the outermost minimization, solving the Lagrange formulation is equivalent!

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

# Example:



$$\min \ x+y \quad s.t. \quad x^2+y^2=1$$

$$\left( \begin{array}{l} \mathcal{L}(x,y,w) = x+y+ w(x^2+y^2-1) \\ \min\limits_{x,y} \ \max\limits_{w} \ \mathcal{L}(x,y,w) \end{array} \right.$$

$$\frac{d\mathcal{L}}{dx} = 1+2wx \overset{=0}{\Longrightarrow} x^* = \frac{-1}{2w}$$

$$\frac{d\mathcal{L}}{dy} = 1+2wy \overset{=0}{\Longrightarrow} y^* = \frac{-1}{2w}$$

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

solving system of three equations:

$$\frac{1}{4w^2} + \frac{1}{4w^2} = 1 \iff w = \pm\sqrt{1/2}$$

The critical points are

$$(-1/\sqrt{2}, -1/\sqrt{2}) \qquad (1/\sqrt{2}, 1/\sqrt{2})$$

$\swarrow$ minimum $\qquad\qquad \swarrow$ maximum

# Iterative Procedure:

For $t=0, \cdots, T-1$:

$$x^t = \arg\min_x \ f(x) + w^t g(x) \qquad \text{(Best response)}$$

$$w^{t+t} = w^t + \eta g(x) \qquad \text{(Incremental update)}$$

Return $\bar{x} = \frac{1}{T}\sum\limits_{t=0}^{T-1} x_t$

Informal Theorem: $\bar{x} \to x^*$ as $T \to \infty$ if $f, g$ convex