# Inverse Reinforcement Learning
## (Inverse Optimal Control)

**TODO**
↳ put in MDP lingo as well

optimal control: given $U$, find $q^* = \text{argmin } U[q]$

$\underline{IOC}$: given $\xi_D$, find $U: \square_{LI} \to \mathbb{R}^+$ s.t. $U[\xi_D] \leq U[\xi], \forall \xi \in \square_{S-g}$

find a cost function(al) that <u>explains</u> the demonstration

why? e.g. $g \to \hat{g}$    $\hat{\xi} = \underset{\xi \in \square_{S-\hat{g}}}{\text{argmin }} U[\xi]$

→ so that we can generalize to new problems

example applications : - driving

- predictable motion

- anything where it's harder to write down tradeoffs than demonstrate behavior

MMP - Maximum Margin Planning

$U[\xi_D] \leq U[\xi] \; \forall \xi \in \square_{S-g}$

$U[\xi_D] \leq \underset{\xi}{\min} U[\xi] \quad (1)$

Problem: $U[\xi] = k, \forall \xi$ solves $(1)$

Solution: make $U[\xi_D]$ better by a margin

$U[\xi_D] \leq \underset{\xi}{\min} [ U[\xi] - \ell[\xi, \xi_D] ]$

with $\ell$ <u>smaller</u> when $\xi$ <u>closer</u> to $\xi_D$

↳ give everything but $\xi_D$ an advantage

e.g. $\ell[\xi, \xi_D] = \begin{cases} 0, & \text{if } \xi = \xi_D \\ 1, & \text{otherwise} \end{cases}$

in practice: $\ell$ smooth, eg $L_2$ norm

$$\max_{u} \min_{\varsigma} \left[ u[\varsigma] - \ell[\varsigma, \varsigma_D] \right] - u[\varsigma_D]$$

$$\min_{u} u[\varsigma_D] - \min_{\varsigma} \left[ u[\varsigma] - \ell[\varsigma, \varsigma_D] \right] +$$
$$+ \frac{\lambda}{2} \, \text{Regularization}(u)$$

Parametrize $u$, search its parameters

$$u[\varsigma] = w^T \int_\varsigma \qquad \int_\varsigma \text{ feature vector, e.g. } \begin{bmatrix} \text{length} \\ \text{obs dist} \\ \text{terrain} \end{bmatrix}$$

$$\min_{w} \underbrace{w^T \int_{\varsigma_D} - \min_{\varsigma} \left[ w^T \int_\varsigma - \ell(\varsigma, \varsigma_D) \right] + \frac{\lambda}{2} \|w\|^2}_{c(w)}$$

$$\min_{w} c(w) \qquad\qquad w_{i+1} = w_i - \alpha \nabla_{w_i} c$$

$$- \min_{\varsigma} \left[ w^T \int_\varsigma - \ell[\varsigma, \varsigma_D] \right] \quad \text{piecewise linear, convex:}$$



$$\text{Let} \quad \varsigma_w^* = \underset{\varsigma}{\text{argmin}} \left[ w^T \int_\varsigma - \ell(\varsigma, \varsigma_D) \right] \qquad w_i \to \varsigma_{w_i}^*$$

$$\nabla_{w_i} \left( - \min_{\varsigma} \left[ w^T \int_\varsigma - \ell[\varsigma, \varsigma_D] \right] = \nabla_{w_i} \left( - \left( w^T \int_{\varsigma_{w_i}^*} - \ell[\varsigma_w^*, \varsigma_D] \right) \right) \right)$$

"subgradient"

$$\nabla_{w_i} c = \int_{\varsigma_D} - \int_{\varsigma_{w_i}^*} + \lambda w_i$$

$$w_{i+1} = w_i - \underbrace{\alpha \lambda w_i}_{\text{shrink } w} - \alpha \underbrace{\left( \int_{\varsigma_D} - \int_{\varsigma_{w_i}^*} \right)}_{\text{gets you to match } \int_{\varsigma_D}}$$

example: $\xi_D$ goes on grass $\xi_{q_D} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ grass rock $\Bigg\} \Rightarrow$

$\xi_{wi}^*$ goes on rock $\xi_{\xi_{wi}^*} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$

$\Rightarrow w_{t+1} = w_i(1 - \alpha\gamma) - \alpha \begin{bmatrix} 1 \\ -1 \end{bmatrix}$

$w(0) - grass - \underline{decreases}$

$w(1) - rock - \underline{increases}$ $\Bigg\} \Rightarrow$

$\Rightarrow$ grass becomes $\underline{cheaper}$ ; rock becomes more expensive

## Maximum Entropy IRL

- assume # is not optimal, and nothing

else: $\frown$ entropy

$\max\limits_{P} H(P)$

$\text{s.t.} \; \mathbb{E}_{\xi \sim P} \big[ U[\xi_D] \big] = u^* + \varepsilon$

$\frown$ rationality coeff (dep. on $\varepsilon$)

$\rightarrow P(\xi | w) = \dfrac{e^{-\beta \, w^T \xi_\xi}}{\sum\limits_\xi e^{-\beta \, w^T \xi_\xi}}$ $\quad \beta = 0 \xrightarrow{\hspace{2cm}} \beta \to \infty$

$\frown$ uniform human

$\frown$ partition function

- $\varepsilon_D$ is not perfect :

$$P(\varsigma \mid w) \propto e^{-w^T \varsigma_\varsigma}$$

- MLE :

$$\max_w P(\varsigma_D \mid w)$$

$$(\Rightarrow) \max_w \log P(\varsigma_D \mid w)$$

$$(\Rightarrow) \max_w \log \frac{e^{-w^T \varsigma_{\varsigma_D}}}{\sum_\varsigma e^{-w^T \varsigma_\varsigma}}$$

$$(\Rightarrow) \max_w -w^T \varsigma_{\varsigma_D} - \log \sum_\varsigma e^{-w^T \varsigma_\varsigma}$$

$$\nabla_w : -\varsigma_{\varsigma_D} - \frac{1}{\sum_\varsigma e^{-w^T \varsigma_\varsigma}} \sum_\varsigma \left[ e^{-w^T \varsigma_\varsigma} \cdot (-\varsigma_\varsigma) \right]$$

$$\nabla_w : -\varsigma_{\varsigma_D} - \sum_\varsigma \frac{e^{-w^T \varsigma_\varsigma}}{\sum_{\bar\varsigma} e^{-w^T \varsigma_{\bar\varsigma}}} (-\varsigma_\varsigma)$$

$$\nabla_w : -\left( \varsigma_{\varsigma_D} - \sum_\varsigma P(\varsigma \mid w) \varsigma_\varsigma \right)$$

$$w_{t+1} = w_t - \alpha \left( \varsigma_{\varsigma_D} - \underbrace{\mathbb{E}_{\varsigma \sim P(\varsigma \mid w)} \varsigma_\varsigma}_{} \right)$$

ascent?

expected feature values
induced by current $w$

contrast to MMP : $\varsigma_{\varsigma_D} - \varsigma_{\varsigma \hat{w}_i}$

noisy demonstration version of MMP — $\varsigma$ not perfect

commonly presented as:

$$\max_{\omega} \underbrace{\min_{\Pi} \left( \lambda H(\Pi) - \underset{s,a \sim \Pi}{\mathbb{E}} [\omega^T f(s,a)] \right)}_{\text{cost of RL on } \omega \text{ (+entropy)}} - \underbrace{\underset{s,a \sim \Pi_D}{\mathbb{E}} [\omega^T f(s,a)]}_{\text{cost of } \Pi_D}$$

cost