

2024.09.11 (Reference: Slide from Abbeel, Levine, Sadigh, & Salakhutdinov)

- BC vs IRL
 - Learn r s.t. $\pi^* = \arg \max_{\theta} \mathbb{E}_{s \sim p(s|s)} r(s, \pi_\theta(s))$
 - $\arg \min_{\theta} \mathbb{E}_{(s, a) \sim p^*} L(a^*, \pi_\theta(s))$

- <pros>
- No reasoning about outcomes or dynamics
 - No notion of intentions
 - Expert can be suboptimal
 - :

- Problem Setting: MDP/r (Same as BC)

→ State Space : S

Action Space : A

An expert policy: $\pi^*: S \rightarrow A$

Transition model: $P_{sa}(S_{t+1}|S_t, A_t)$

(Given) (Sometimes)

} Goal: Learn a reward function $r(S, a)$ or $r(S)$ assuming $\pi^* = \text{optimal}$

Describes desirability of being in a state.

Dynamics Model $T = P_{sa}$

Probability distribution over next states given current state and action

Reward Function R

Reinforcement Learning / Optimal Control

(Given Demonstration)

* Recover *

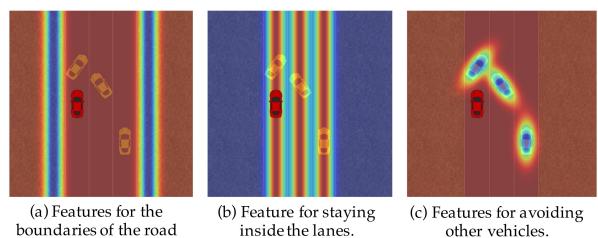
$$\arg \max_{\pi} \mathbb{E}[\sum_t \gamma^t R(s_t)|\pi]$$

Prescribes action to take for each state

- Big Assumption (1): Reward func. is a linear combination of features.

$$R(s) = w^\top \varphi(s)$$

↓ ↓
weights features of state \Rightarrow
 $\varphi: S \rightarrow \mathbb{R}^n$
 $w \in \mathbb{R}^n$



- Big Assumption (2): Expert is optimal, R^* explains π^* .

$$\rightarrow \text{Find } R^* \text{ s.t. } \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R^*(s_t) | \pi^*\right] \geq \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R^*(s_t) | \pi\right] \quad \forall \pi.$$

- Combine two assumptions.

Challenges.

$$\Rightarrow R(s) = w^\top \varphi(s)$$

$$\underbrace{\mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R^*(s_t) | \pi^*\right]}_{V^{\pi}(s) \text{ in RL!}} \geq \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R^*(s_t) | \pi\right] \quad \forall \pi$$

↳ how to compute for π^* ?
"Requires scale"

$$\begin{aligned} \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) | \pi\right] &= \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t w^\top \varphi(s_t) | \pi\right] \\ &= w^\top \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t \varphi(s_t) | \pi\right] \\ &= w^\top M(\pi) \end{aligned}$$

Depend on state visitation distribution.
feature expectations.

$$\Rightarrow \text{Find } w^* \text{ s.t. } w^{*\top} M(\pi^*) \geq w^{*\top} M(\pi), \quad \forall \pi. \quad (\text{If } \pi^* \text{ is suboptimal, Infeasible!})$$

$$(V^{\pi^*}(s) \geq V^{\pi}(s) \quad \forall \pi)$$

R=0 can be solution
<ambiguity>

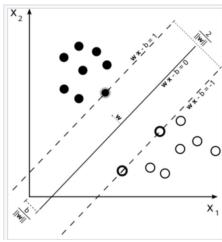
* Reward Ambiguity : Many R exist under π^* 's demonstrations.

↳ How to pick R ?

- Maximum Margin Planning
- Maximum Entropy IRL

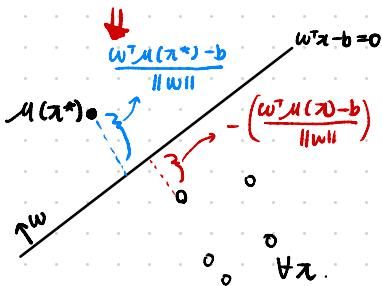
(MMP)

- Maximum Margin Planning : Find R that best separates π^* from π .



"Minimize $\|w\|$ subject to $y_i(w \cdot x_i - b) \geq 1$, for $i = 1, \dots, n$ "

Given a training dataset of $(x_1, y_1), \dots, (x_n, y_n)$, where y_i is either 1 or -1 identifying the class x_i is in. We want to find the maximum margin hyperplane that divides the points so the distance between the hyperplane and the nearest point from each class is maximized.



$$w^\top \mu(\pi^*) \geq w^\top \mu(\pi)$$

$$\Rightarrow w^\top \mu(\pi^*) - b \geq 0 \quad \& \quad w^\top \mu(\pi) - b \leq 0$$

OUR goal

$$\max_w \min_{\pi} \frac{(w^\top \mu(\pi^*) - b)}{\|w\|} + \frac{(-w^\top \mu(\pi) - b)}{\|w\|}$$

$$= \max_w \min_{\pi} \frac{w^\top \mu(\pi^*) - w^\top \mu(\pi)}{\|w\|}$$

$$= \max_w \frac{1}{\|w\|} \min_{\pi} w^\top \mu(\pi^*) - w^\top \mu(\pi)$$

$$\min_w \|w\|$$

$$\text{s.t., } w^\top \mu(\pi^*) - w^\top \mu(\pi) \geq 1$$

If student's less know linear algebra...

$$\begin{aligned} x_i &= x_i - d \\ d &= \alpha^\top w \\ w^\top x_i - b &= 0 \\ \Rightarrow w^\top (x_i - \alpha w) - b &= 0 \\ &= w^\top (x_i - \alpha w) - b = 0 \\ w^\top x_i - b &= 0, \quad \alpha = \frac{w^\top x_i - b}{w^\top w} \\ \Rightarrow \|w\|_2 &= \sqrt{d^\top d} = \sqrt{\alpha^\top w^\top w} = \frac{\|w^\top x_i - b\|}{\|w\|_2}. \end{aligned}$$

scale is invariant. Let's be clever, choose margin

MMP recap

$$\rightarrow \min_{\mathbf{w}} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } \mathbf{w}^\top \mathbf{u}(\pi^*) \geq \mathbf{w}^\top \mathbf{u}(\pi) + 1 \quad \forall \pi.$$

can be $m(\pi^*, \pi)$ (adaptive margin)

Max Entropy IAL.

$$\rightarrow \Sigma = \{(s_1, a_1), \dots, (s_T, a_T)\} : \text{trajectory}, \Sigma \in \Xi$$

$$\rightarrow D = \{\Sigma_1, \dots, \Sigma_{101}\} : \text{set of } \pi^* \text{'s demonstration}$$

$$\rightarrow f: \Xi \mapsto \mathbb{R}^n : \text{traj. feature.}, f_D = \frac{1}{|D|} \sum_{\Sigma \in D} f(\Sigma) \text{ is empirical feature expectation.}$$

$$\text{Then, } \mathbb{E}_{\Sigma \sim P(\Sigma)} \left[\sum_{t=1}^T \delta^t R(s_t) \right] = \mathbb{E}_{\Sigma \sim P(\Sigma)} \left[\sum_{t=1}^T \delta^t w^\top \mathbf{u}(s_t) \right] = \underbrace{w^\top \mathbb{E}_{\Sigma \sim P(\Sigma)} [f(\Sigma)]}_{\text{"weighted" feature expectation}}$$

Let's make noisy assumption (yet rational)

$$P(\Sigma) \propto \exp(R(\Sigma))$$

$$\text{Goal: } \max_P \int -P(\Sigma) \log P(\Sigma) d\Sigma$$

$$\text{s.t. } \mathbb{E}_{\Sigma \sim P(\Sigma)} [f(\Sigma)] = \int P(\Sigma) f(\Sigma) d\Sigma = f_D$$

$$\int P(\Sigma) d\Sigma = 1$$

$$P(\Sigma) \geq 0, \forall \Sigma \in \Xi$$

Find P from D
that matches f_D ,

& maximize Entropy.

strong preference for low-cost paths,
equal cost paths be equally probable.

$$\Rightarrow \text{Lagrangian: } L(P, \lambda, v) \Rightarrow \int -P(\Sigma) \log P(\Sigma) d\Sigma + \lambda^\top (\int P(\Sigma) f(\Sigma) d\Sigma - f_D) + v (\int P(\Sigma) d\Sigma - 1)$$

$$\Rightarrow \int (-P(\Sigma) \log P(\Sigma) + \lambda^\top P(\Sigma) f(\Sigma) + v P(\Sigma)) d\Sigma - v - \lambda^\top f_D$$

$\triangleright F(P(\Sigma))$

$$\Rightarrow \int F(P(\Sigma)) d\Sigma - v - \lambda^\top f_D \quad \triangleright \nabla_P L(P, \lambda, v) = 0$$

set gradient of each term
in the integral.

$$\rightarrow \frac{\partial F(P)}{\partial P} = -\log P(\Sigma) - 1 + \lambda^\top f(\Sigma) + v = 0$$

$$\rightarrow P^*(\Sigma) = \exp(\lambda^\top f(\Sigma) + v - 1)$$

$$\rightarrow L(P^*, \lambda, v) = \int (-\exp(\lambda^\top f(\Sigma) + v - 1) \cdot (\lambda^\top f(\Sigma) + v - 1) + \lambda^\top \exp(\lambda^\top f(\Sigma) + v - 1) + v \exp(\lambda^\top f(\Sigma) + v - 1)) d\Sigma - v - \lambda^\top f_D$$

$$\Rightarrow \int \exp(\lambda^\top f(\Sigma) + v - 1) d\Sigma - \lambda^\top f_D - v$$

$$\rightarrow v^* = -\log(\int \exp(\lambda^\top f(\Sigma) - 1) d\Sigma)$$

$$\rightarrow v^* ? : \frac{\partial L(P^*, \lambda, v)}{\partial v} = 0 \Rightarrow e^{v-1} \int \exp(\lambda^\top f(\Sigma)) d\Sigma - 1 = 0$$

$$\Rightarrow e^{-v} = \int \exp(\lambda^\top f(\Sigma) - 1) d\Sigma$$

Recap

$$\rightarrow P^*(\xi) = \exp(\lambda^\top f(\xi) + v - 1)$$

$$v^* = -\log(\int \exp(\lambda^\top f(\xi) - 1) d\xi)$$

$$\Rightarrow \text{Put } v^* \text{ to } P^* \Rightarrow P^* = \exp(\lambda^\top f(\xi) - \log(\int \exp(\lambda^\top f(\xi) - 1) d\xi) - 1)$$

$$\Rightarrow P^* = \frac{\exp(\lambda^\top f(\xi) - 1)}{\int \exp(\lambda^\top f(\xi) - 1) d\xi} = \frac{\exp(\lambda^\top f(\xi))}{\int \exp(\lambda^\top f(\xi)) d\xi}$$

match to $P(\xi) \propto \exp(R(\xi))$

if we set $R(\xi) = \underline{\lambda^\top f(\xi)}$

Reward weights are dual variable λ .

$$\text{Finally, } \lambda^* = \arg \max_{\lambda} P(\xi_{1:D} | \lambda)$$

$$= \arg \max_{\lambda} \lambda^\top f_D - \log(\int \exp(\lambda^\top f(\xi)) d\xi)$$

Then, we do gradient ascent! $\nabla_{\lambda} M = f_D - \mathbb{E}_{\xi \sim P(\xi | \lambda)} [f(\xi)]$

$$\lambda_{i+1} \leftarrow \lambda_i + \alpha \nabla_{\lambda} M.$$

- 1) Initialize λ and collect expert demonstrations D .
- 2) Solve for the optimal policy $\pi_{\lambda}(a|s)$ with respect to λ . (RL Loop)
- 3) Solve for state visitation frequencies $P_{\lambda}(s)$.
- 4) Compute the gradient $\nabla_{\lambda} M$.
- 5) Update λ with one gradient step.

This assumes access to the dynamics (transition function) and having low dimensional systems to be able to solve for the policy using RL.

$$\nabla_{\lambda} M ? \Rightarrow \frac{\partial (\log \int \exp(\lambda^\top f(\xi)) d\xi)}{\partial \lambda} ?$$

$$\log g(\lambda), \quad g(\lambda) = \int \exp(\lambda^\top f(\xi)) d\xi$$

$$\begin{aligned} \frac{\partial \log g(\lambda)}{\partial \lambda} &= \frac{1}{g(\lambda)} \cdot \frac{\partial g(\lambda)}{\partial \lambda}, \quad \frac{\partial g(\lambda)}{\partial \lambda} = \int \frac{d}{d\lambda} \exp(\lambda^\top f(\xi)) d\xi \\ &= \int f(\xi) \exp(\lambda^\top f(\xi)) d\xi \end{aligned}$$

$$\Rightarrow \frac{\int f(\xi) \exp(\lambda^\top f(\xi)) d\xi}{\int \exp(\lambda^\top f(\xi)) d\xi} = \mathbb{E}[f(\xi)]$$

↳ Lagrangian?

$$\begin{aligned} \rightarrow \min_x f(x) \\ \text{s.t. } g(x) = 0 \end{aligned}$$

$$\begin{aligned} \xrightarrow{\text{Lag. multiplier}} L(x, \lambda) &= f(x) + \lambda \cdot g(x) \\ &\xrightarrow{\text{Find partial derivative w.r.t. } x \text{ & set to 0, find which } x \text{ is "optimal"}} \end{aligned}$$

$$\text{Ex. } \min_{x,y} f(x,y) = xy \quad \rightarrow L(x,y,\lambda) = xy - \lambda(x+y-1) \\ \text{s.t. } x+y=1$$

$$\begin{cases} \frac{\partial L}{\partial x} = y - \lambda = 0 \\ \frac{\partial L}{\partial y} = x - \lambda = 0 \\ \frac{\partial L}{\partial \lambda} = -(x+y-1) = 0 \end{cases} \quad \begin{cases} \lambda = \frac{1}{2} \\ x = \lambda \\ y = \lambda \end{cases}$$