

Lecture 27: Max Entropy IRL

Recall the IRL setting: finite horizon MDP

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

where reward is unobserved and we have a dataset $\mathcal{D} = \{s_i^*, a_i^*\} \sim d_{\mathcal{M}}^{\pi^*}$

Further assume linear reward function

$$r(s, a) = \underbrace{\Theta_*^T}_{\text{unknown}} \underbrace{\phi(s, a)}_{\text{known}}$$

Notice that $\mathbb{E}_{s, a \sim d_1} \phi(s, a) = \mathbb{E}_{s, a \sim d_2} \phi(s, a) \Rightarrow \mathbb{E}_{s, a \sim d_1} r(s, a) = \mathbb{E}_{s, a \sim d_2} r(s, a)$

(why? by linearity of expectation)

state/action distributions:

$P_h^{\pi}(s, a; \mathcal{M})$: probability of visiting (s, a) at step h using π

$d_{\mathcal{M}}^{\pi}(s, a) = \sum_{h=0}^{H-1} P_h^{\pi}(s, a; \mathcal{M}) / H$ average state-action distribution

$d_{\mathcal{M}}^{\pi}(s) = \sum_{a \in \mathcal{A}} d_{\mathcal{M}}^{\pi}(s, a)$ average state distribution.

Max-Ent IRL Problem

last lecture we arrived at the constrained optimization problem:

$$\min_{\pi} \mathbb{E}_{s, a \sim \pi} [\log(\pi(a|s))] \leftarrow \begin{array}{l} \text{equivalent to} \\ \text{maximizing} \\ \text{entropy} \end{array}$$

ensuring consistency w/ expert π^*

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

estimate from \mathcal{D}

d dimensional expression, $\phi: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$

Then using the Lagrange formulation:

$$\max_{w \in \mathbb{R}^d} \min_{\pi} \underbrace{\mathbb{E}_{s, a \sim \pi} [\log \pi(a|s)] + w^T (\mathbb{E}_{s, a \sim \pi} [\phi(s, a)] - \mathbb{E}_{s, a \sim \pi^*} [\phi(s, a)])}_{\mathcal{L}(\pi, w)}$$

Notice that we can write

$$\mathcal{L}(\pi, w) = \mathbb{E}_{s, a \sim \pi} [\log \pi(a|s) - w^T \phi(s, a)] + w^T \mathbb{E}_{s, a \sim \pi^*} [\phi(s, a)]$$

2) Iterative Max-Ent IRL

Initialize $w_0 \in \mathbb{R}^d$

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

$$\pi^t = \underset{\pi}{\operatorname{argmax}} \mathbb{E}_{s, \operatorname{and}_{\eta}^{\pi}} [w_t^T \phi(s, a) - \log \pi(a|s)] \quad (\text{best response})$$

$$w_{t+1} = w_t + \eta \left(\mathbb{E}_{s, \operatorname{and}_{\eta}^{\pi^t}} \phi(s, a) - \mathbb{E}_{s, \operatorname{and}_{\eta}^{\pi^t}} \phi(s, a) \right) \quad (\text{gradient update})$$

Return $\bar{\pi} = \operatorname{Uniform}(\pi^0, \dots, \pi^{T-1})$

$$\bar{\pi}(a|s) = \frac{1}{T} \sum_{t=0}^{T-1} \pi^t(a|s)$$

Best Response step is like RL problem with reward $w_t^T \phi(s, a)$ and addition policy dependent term $-\log \pi(a|s)$

3) Soft Value Iteration

We will solve this minimization problem with dynamic programming!

$$\arg\max_{\pi} \mathbb{E}[r(s,a) - \log \pi(a|s)] = \arg\max_{\pi} \mathbb{E} \left[\underbrace{\sum_{t=0}^{T-1} r(s_t, a_t) - \log \pi_t(a_t|s_t)}_{\substack{s_{t+1} \sim P(s_t, a_t) \\ a_t \sim \pi_t(s_t)}} \right]$$

Initialize $V_H^*(s) = 0$

For $h = H-1, \dots, 0,$

$$1) Q_h^*(s, a) = r(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} V_{h+1}^*(s')$$

$$2) \pi_h^*(\cdot|s) = \underset{p \in \Delta(\mathcal{A})}{\operatorname{argmin}} \sum_{a \in \mathcal{A}} p(a) Q_h^*(s, a) + \sum_{a \in \mathcal{A}} p(a) \underbrace{\log(p(a))}_{\text{entropy}}$$

constraint $\sum_a p(a) = 1$

(derivation below)

$$= \frac{\exp(Q_h^*(s, a))}{\sum_{a' \in \mathcal{A}} \exp(Q_h^*(s, a'))}$$

$$3) V_h^*(s) = \mathbb{E}_{a \sim \pi_h^*(\cdot|s)} [-\log \pi_h^*(a|s) + Q_h^*(s, a)]$$

(derivation below)

$$= \log \left(\sum_{a \in \mathcal{A}} \exp(Q_h^*(s, a)) \right)$$

Contrast π_h^*, V_h^* with the classic RL solution!

$$\operatorname{softmax}_a Q_h^*(s, a) \quad \text{vs.} \quad \max_a Q_h^*(s, a)$$

Deriving the softmax policy:

$$\max_{\rho} \min_w \sum_a \rho_a (\underbrace{Q_n^*(s,a) + w}_{\mathcal{L}(\rho, w)} - \rho_a \log \rho_a) + w (\sum_a \rho_a - 1)$$

$$\nabla_{\rho} \mathcal{L} = \nabla_{\rho} \left[\sum_a \rho_a (Q_n^*(s,a) + w) - \rho_a \log \rho_a \right]$$

$$\stackrel{=0}{\Rightarrow} Q_n^*(s,a) + w - \log \rho_a - 1 = 0 \quad \forall a$$

$$\rho_a = \exp(Q_n^*(s,a)) \exp(-1+w)$$

$$\nabla_w \mathcal{L} = \sum_a \rho_a - 1 \Rightarrow \sum_a \rho_a = 1$$

Combining the equations,

$$\sum_a \exp(Q_n^*(s,a)) (\exp(-1+w)) = 1$$

$$\Rightarrow \exp(-1+w) = \frac{1}{\sum_a \exp(Q_n^*(s,a))} \quad \checkmark$$

Deriving the value:

$$V_n^*(s) = \mathbb{E}_{a \sim \pi_n^*(\cdot|s)} [-\log \pi_n^*(a|s) + Q_n^*(s,a)]$$

$$= \mathbb{E}_{a \sim \pi_n^*(\cdot|s)} \left[\cancel{-Q_n^*(s,a)} + \log \left(\sum_{a' \in \mathcal{A}} \exp(Q_n^*(s,a')) \right) + \cancel{Q_n^*(s,a)} \right]$$

$$\text{(independent of } a) = \log \left(\sum_{a' \in \mathcal{A}} \exp(Q_n^*(s,a')) \right) \mathbb{E}_{a \sim \pi_n^*(\cdot|s)} [1] \quad \checkmark$$