

# 1) Max Entropy IRL

Finite horizon MDP

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

unknown & unobserved

expert dataset:  $\mathcal{D} = \{s_i^*, a_i^*\} \sim d_{\mathcal{M}}^{\pi^*}$

## State/action distributions

$P_h^{\pi}(s, a; \mathcal{M})$  = probability of visiting  $(s, a)$  at timestep  $h$  following  $P, \pi$

$$d_{\mathcal{M}}^{\pi}(s, a) = \frac{1}{H} \sum_{h=0}^{H-1} P_h^{\pi}(s, a; \mathcal{M})$$

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

## Linear Rewards

$$r(s, a) = \underbrace{\Theta_*^T}_{\text{unknown}} \underbrace{\phi(s, a)}_{\text{known}} \quad \Theta_* \in \mathbb{R}^d \quad \phi: S \times \mathcal{A} \rightarrow \mathbb{R}^d$$

$$\mathbb{E}_{s, a \sim d_{\mathcal{M}}^{\pi}} \phi(s, a) = \mathbb{E}_{s, a \sim d_{\mathcal{M}}^{\pi^*}} \phi(s, a) \Rightarrow \mathbb{E}_{s, a \sim d_{\mathcal{M}}^{\pi}} r(s, a) = \mathbb{E}_{s, a \sim d_{\mathcal{M}}^{\pi^*}} r(s, a)$$

## Max-Ent IRL Problem

$$\min_{\pi} \mathbb{E}_{s, a \sim d_{\mathcal{M}}^{\pi}} [\log(\pi(a|s))] \leftarrow \text{max ent} \hat{=} \text{min log likelihood}$$

s.t.  $\mathbb{E}_{s, a \sim d_{\mathcal{M}}^{\pi}} \phi(s, a) = \mathbb{E}_{s, a \sim d_{\mathcal{M}}^{\pi^*}} \phi(s, a) \quad \mathcal{D}$   $\leftarrow$  estimate from  $\mathcal{D}$

using the Lagrange Formulation

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

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

2) Iterative Max-Ent RL

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

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

$$\left[ \pi^t = \operatorname{argmax}_{\pi} \mathbb{E}_{s, a \sim \pi} [-\log \pi(a|s) + w_t^T \phi(s, a)] \right] \leftarrow \text{soft-VI}$$

$$w_{t+1} = w_t + \eta \left( \mathbb{E}_{s, a \sim \pi^*} \phi(s, a) - \mathbb{E}_{s, a \sim \pi^t} \phi(s, a) \right) \leftarrow \checkmark$$

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

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

Best response step: RL problem with reward  $w_t^T \phi(s, a)$  and a policy dependent term  $-\log \pi(a|s)$  (exercise)

$$\mathbb{E}_{s \sim \mu} [V^{\pi}(s)] = \mathbb{E} \left[ \sum_{t=0}^{H-1} r_t \mid r, p, \pi, \mu \right] = H \cdot \mathbb{E}_{s, a \sim \pi} [r(s, a)]$$

### 3) Soft Value Iteration

Use dynamic programming:

$$\operatorname{argmax}_{\pi} \mathbb{E} \left[ \sum_{t=0}^{H-1} \underbrace{r(s_t, a_t)} - \underbrace{\log \pi_t(a_t | s_t)} \right] \left[ \begin{array}{l} s_t \sim P(s_t, a_t) \\ a_t \sim \pi_t(s_t) \\ s_0 \sim \mu \end{array} \right]$$

Initialize  $V_H^*(s) = 0$

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

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

$$2) \pi_h^*(\cdot | s) = \operatorname{argmax}_{\pi} \mathbb{E}_{a \sim \pi(\cdot | s)} [Q_h^*(s, a) - \log \pi(a | s)]$$

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

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

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

Contrast  $\pi^* V^*$  with classic RL solution

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

$$\Pi(-/s) = \left[ \underset{\rho}{\operatorname{argmax}} \sum_{a \in \mathcal{A}} \rho(a) (Q_n^*(s, a) - \log \rho(a)) \right]$$

s.t.  $\left( \sum_{a \in \mathcal{A}} \rho(a) = 1 \right)$

$$\mathcal{J}(\rho, w) = \sum_{a \in \mathcal{A}} \rho(a) \underbrace{Q_n^*(s, a)}_{\times \log x} - \underbrace{\rho(a) \log \rho(a)}_{\times \log x} + w \left( \sum_a \rho(a) - 1 \right)$$

$\rho \in \mathbb{R}^{\mathcal{A}}, w \in \mathbb{R}$

$$\frac{\partial \mathcal{J}}{\partial \rho(a)} = Q_n^*(s, a) - \log \rho(a) - \frac{\rho(a)}{\rho(a)^2} + w = 0$$

$$\rho(a) = \exp(Q_n^*(s, a)) \cdot \exp(w-1) \quad \forall a$$

$\frac{\partial \mathcal{J}}{\partial w}$

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