

# Apprenticeship Learning via Inverse Reinforcement Learning

- Key insight: Instead of learning the policy from the "expert", it learns a reward function

## Preliminaries:

$S$ : states,  $A$ : actions,  $T = \{P_{sa}\}$  transition probability,  $\gamma \in [0, 1)$  discount factor,

$D$ : initial-state distribution,  $R: S \mapsto A$  reward function  $\leq 1$ . MDP  $\mid R$  (without a reward function)

$\phi$ : features  $S \rightarrow [0, 1]^k$ ,  $R^*(s) = w^* \cdot \phi(s), w^* \in \mathbb{R}^k$  true reward function

$$\pi: \text{policy } S \mapsto D_A, E_{s \sim D} [V^\pi(s_0)] = E[\sum_{t=0}^{\infty} \gamma^t R(s_t) | \pi] = E[\sum_{t=0}^{\infty} \gamma^t w^* \cdot \phi(s_t) | \pi] = w^* \cdot E[\sum_{t=0}^{\infty} \gamma^t \phi(s_t) | \pi]$$

$$\mu(\pi): \text{feature expectations}, \mu(\pi) = E[\sum_{t=0}^{\infty} \gamma^t \phi(s_t) | \pi] \in \mathbb{R}^k = w^* \cdot \underbrace{\mu(\pi)}_{\mu(\pi)}$$

Note,  $R$  is a linear combination of  $\phi$ .

Let  $\pi_1, \pi_2 \in \Pi$ ,  $\pi_3$  be  $\pi_1 + \pi_2$  with  $P(\pi_1) = \lambda, P(\pi_2) = 1 - \lambda \Rightarrow \mu(\pi_3) = \lambda \mu(\pi_1) + (1 - \lambda) \mu(\pi_2)$

Generally,  $\pi_1, \dots, \pi_d \in \Pi \rightarrow$  convex combination  $\sum_{i=1}^d \lambda_i \mu(\pi_i)$  ( $\lambda_i \geq 0, \sum \lambda_i = 1$ )

$\pi_E$ : "expert" policy, can be viewed as optimal  $R^* = w^{*T} \phi$

Estimator:  $\hat{M}_E = \hat{M}(\pi_E)$ ,  $m$  trajectories  $\{S_0^{(i)}, S_1^{(i)}, \dots\}_{i=1}^m$  generated by expert

$$\hat{M}_E = \frac{1}{m} \sum_{i=1}^m \sum_{t=0}^{\infty} \gamma^t \phi(S_t^{(i)})$$

## Algorithm:

Goal: find  $\tilde{\pi}$  s.t.  $\|\mu(\tilde{\pi}) - M_E\|_2 \leq \epsilon$

$$|E[\sum_{t=0}^{\infty} \gamma^t R(s_t) | \pi_E] - E[\sum_{t=0}^{\infty} \gamma^t R(s_t) | \tilde{\pi}]|$$

$$= |w^* \mu(\pi_E) - w^* \hat{M}_E| \leq \|w\|_2 \|\mu(\tilde{\pi}) - M_E\|_2 \leq 1 \cdot \epsilon = \epsilon$$

Problem turns into finding a policy  $\tilde{\pi}$  that makes  $\mu(\tilde{\pi})$  close to  $M_E$

Steps: (QP-based)

1. Randomly pick  $\pi^{(0)}$ , compute  $\mu^{(0)} = \mu(\pi^{(0)})$ , set  $i = 1$ .

2. Compute  $t^{(i)} = \max_{\|w\|_2=1} \min_{\pi \in \Pi} w^T (M_E - \mu^{(i)})$ , let  $w^{(i)}$  be the max

3. If  $t^{(i)} \leq \epsilon$ , terminate

4. Compute the optimal policy  $\pi^{(i)}$  using  $R = (w^{(i)})^T \phi$

5. Compute  $\mu^{(i)} = \mu(\pi^{(i)})$

6. Set  $i = i+1$ , go to step 2

My comment, this is kinda a trial & error algorithm

Projection-based method: replace step (2) by  
 set  $\bar{\mu}^{(i-1)} = \bar{\mu}^{(i-2)} + \frac{(\mu^{(i-1)} - \bar{\mu}^{(i-2)})^T (\mu_E - \bar{\mu}^{(i-2)})}{(\mu^{(i-1)} - \bar{\mu}^{(i-2)})^T (\mu^{(i-1)} - \bar{\mu}^{(i-2)})} (\mu^{(i-1)} - \bar{\mu}^{(i-2)})$   
 set  $w^{(i)} = \mu_E - \bar{\mu}^{(i-1)}, \gamma^{(i)} = \|\mu_E - \bar{\mu}^{(i-1)}\|_2$

- Theoretical results

Theorem 1: Let an MDPVR, features  $\phi: S \mapsto [0,1]^k$ , and any  $\epsilon > 0$  be given. Then the apprenticeship learning algorithm will terminate with  $t^{(i)} \leq \epsilon$  after at most  
 $n = O\left(\frac{k}{(1-\gamma)^2 \epsilon} \log \frac{k}{(1-\gamma)\epsilon}\right)$  Iterations.

Theorem 2: lower bound for  $m$ , (samples required).

$$m \geq \frac{2k}{(\epsilon(1-\gamma))^2} \log \frac{2k}{\delta}$$

Proofs are skipped.