

687 2017-12-12

## Inverse RL

RL: Given samples  $S, A, P, R, d, \gamma$  find  $\pi^*$

IRL: " "  $S, A, P, d, \gamma, \pi^*$  find  $R$

↳ E.g., learning from demonstration (LfD)

- Abstraction selection

e.g.

$$R_E = W^T \phi(S_T)$$

## Multiagent RL

Stochastic game generalization of an MDP

$(\mathcal{S}, A_1, A_2, \dots, A_n, R_1, R_2, \dots, R_n, d, \gamma)$

$P(S, a_1, a_2, \dots, a_n, S') = \Pr(S_{t+1} = S' | S_t = s, A_1 = a_1, \dots, A_n = a_n)$

$R_i(S, a_1, a_2, \dots, a_n, S')$

- Fully cooperative if  $\forall i, j, R_i = R_j$

- Otherwise, need to employ Game Theory  $\rightarrow$  equilibria

## Pareto Frontier + multi-objective optimization

Let  $f_1, \dots, f_n$  be  $n$  objectives, where each  $f_i : \mathcal{X} \rightarrow \mathbb{R}$

The Pareto frontier is a set of solutions:

$$P \triangleq \{x \in \mathcal{X} : \forall x' \in \mathcal{X} : (\exists i : f_i(x') > f_i(x)) \Rightarrow (\exists j : f_j(x') < f_j(x))\}$$

making  $f_i$  better      makes  $f_j$  worse

## RL Theory

PAC-MDP : Probably approximately correct in MDPs

(To probability  $1-\delta$ , after a fixed number of timesteps less than a polynomial in  $|\mathcal{S}|, |\mathcal{A}|, 1/\epsilon, 1/\delta$ , and  $1-\gamma$ , it returns a policy whose expected return is within  $\epsilon$  of  $J(\pi^*)$ .)

- Sample complexity : the polynomial fn. [Kakade 2003]

## Hoeffding's Inequality

Let  $x_1, x_2, \dots, x_n$  be  $n$  i.i.d. random variables

s.t.  $\mu = \mathbb{E}[x_i]$  and  $x_i \in [a, b]$  (or  $\Pr(x_i \in [a, b]) = 1$ )

then

$$\Pr(\mu \geq \frac{1}{n} \sum_{i=1}^n x_i - (b-a)\sqrt{\frac{\ln 1/\delta}{2n}}) \geq 1-\delta$$



Theory:  $4 \times 4, 5 \times 5$  grid world Lattimore/Hutter 2015  
 $\hookrightarrow 10^{11}$  steps to guarantee within 10% to high prob.

Kearning (sp?) + Singh 1998

## Regret

The regret of an algorithm over  $T = K \cdot L$  time steps  
is  $\text{Regret}(K) \triangleq \sum_{k=1}^K v^*(s_k) - v^{\pi_k}(s_k)$

some constant

horizon

## UCB & Thompson Sampling

Bandit problem: Always choose bandits w/ higher upper bounds (e.g., via Hoeffding's)

$\hookrightarrow$  Sample according probability distributions that you update as you go. Ian Osband

Lower bound:  $\Omega\left(\sqrt{\frac{HSA}{L}}\right)$

Azar, Osband, Munos:  $\tilde{O}\left(\sqrt{HSAT} + H^2S^2A + H\sqrt{T}\right)$   
(upper bound, i.e., an algorithm)

## Off-policy (policy) evaluation

$$D = \{H_i\}_{i=1}^n \quad H: \pi_b \leftarrow \text{behavior policy}$$

Goal: estimate  $J(\pi_e)$  given  $D, \pi_b$

Method 1: model-based

Method 2: Sampling-based: Importance-sampling

$$\hat{J}(\pi_e, D) = \frac{1}{n} \sum_{i=1}^n \frac{\Pr(H_i | \pi_e)}{\Pr(H_i | \pi_b)} g(H_i) \quad \xrightarrow{\text{discounted sum of rewards}}$$

adjusts for likelihood of different histories  $H_i$  under different policies.

$$\begin{aligned} &= \frac{1}{n} \sum_{i=1}^n \frac{d_o(s_0) \pi_e(a_1 | s_1) P(s_1, a_1, s_2) P(R_1) \dots}{d_o(s_0) \pi_b(a_1 | s_1) P(s_1, a_1, s_2) P(R) \dots} g(H_i) \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\pi_e(a_1 | s_1) \pi_e(a_2 | s_2) \dots}{\pi_b(a_1 | s_1) \pi_b(a_2 | s_2) \dots} g(H_i) \end{aligned}$$

Batch RL - works from a batch of data from one policy  
 ↳ get a better policy

Safe RL - all policy moves are improvements

Average reward RL - average, not discounted, reward  
 Sridhar Mahadevan

## Deep Learning RL

DL → gives a function approximation space

RL → how to train a fu. approx. for

sequential decision problems (MDPs)

Q learning + Convolutional neural net + experience replay  
 ↳ tweaked to use target net

Dueling networks 199x Baird's Advantage Updating

→ Atari