

Last Time

- Multi Obj
- Constrained
- Linear Weighted

- Risk Aware

- Pure Info Gathering

This time: VOI, IRL

o, o' ← value of observing
↑ already

$$VOI(o' | o) = \sum_{o'} (P(o' | o) EU^*(o, o')) - EU^*(o)$$

Example



$$a \sim U([-r, r])$$

$$U(c) = \begin{cases} 1 & \text{if } |a-c| \leq r/4 \\ 0 & \text{otherwise} \end{cases}$$

VOI(A)?

$$VOI(A) = \int_a P(a) EU^*(a) da - EU^*(c)$$

$\uparrow \frac{1}{2\pi}$ $\uparrow 1$ $\uparrow \frac{1}{4}$

$$= \frac{3}{4}$$

$$VOI(o' | b) = \sum_{o'} P(o' | b) V^*(\tau(b, o')) - V^*(b)$$

Inverse RL IRL in IRL

forward RL

given S, A, T, R

find π^*

inverse

given $S, A, T, \{\tau_i\}$

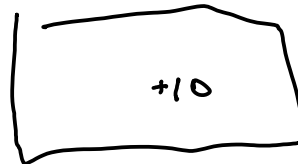
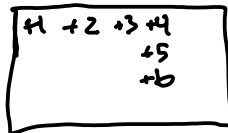
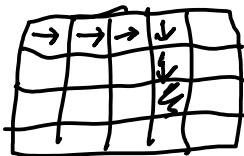
find R

↑ trajectories
sampled with
 π^*



Imitation Learning

Reinforcement Learning



Underspecification

$$\text{Linear } R_{\psi}(s, a) = \sum_i \psi_i f(s, a)$$

$$\begin{matrix} s = 0 \\ a = 0 \\ 0 \end{matrix} \begin{matrix} \times \\ \times \\ \times \end{matrix} \begin{matrix} 0 \\ 0 \\ 0 \end{matrix} R_{\psi}(s, a)$$

↑ parameters ψ

4 approaches

Feature matching IRL

π^{ψ} = optimal policy for r_{ψ}

$$\text{pick } \psi \text{ s.t. } E_{\pi^{\psi}}[f(s, a)] = E_{\pi^*}[f(s, a)]$$

↑ state-action
marginal distribution

↑ from samples

still ambiguous

Maximum margin

maximize m
 ψ, m

subject to $\psi^T E_{\pi^*} [f(s, a)] \geq \max_{\pi} \psi^T E_{\pi} [f(s, a)] + m$

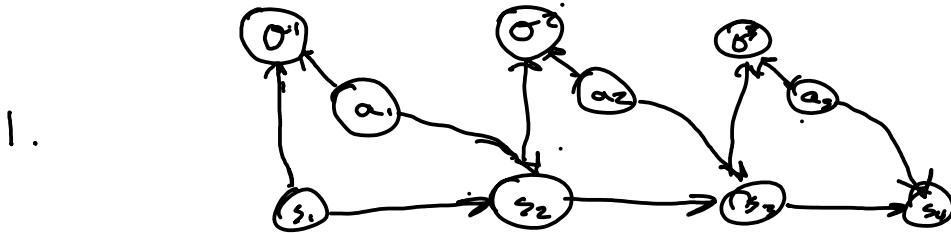


minimize $\frac{1}{2} \|\psi\|^2$
 ψ

s.t $\psi^T E_{\pi^*} [f(s, a)] \geq \max_{\pi} \psi^T E_{\pi} [f(s, a)] +$

$D(\pi, \pi^*)$

- kind of arbitrary
- No model of expert subopt.
- Messy Constrained Opt.



$$p(\sigma_t) \propto \exp(r(s_t, a_t))$$

$$\pi(a_t | s_t) \propto e^{Q_t(s_t, a_t) - V(s_t)} = e^{A_t(s_t, a_t)}$$

$$p(\tau, \sigma_{1:T}) = \frac{p(\tau, \sigma_{1:T})}{p(\sigma_{1:T})}$$

$$\propto p(\tau) \prod_t \exp(r(s_t, a_t)) = p(\tau) \exp(\sum_t r(s_t, a_t))$$

$$p(\tau | \sigma_{1:T}, \psi) \propto p(\tau) \exp(\sum_t r_\psi(s_t, a_t))$$

$$\max_{\psi} \frac{1}{N} \sum_{i=1}^N \log p(\tau_i | \sigma_{1:T}, \psi)$$

$$= \max_{\psi} \frac{1}{N} \sum_i r_{\psi}(\tau_i) - \log Z$$

$$Z = \int p(\tau) \exp(r_{\psi}(\tau)) d\tau$$

↑ normalization
 "partition function"
 hard part

$$\nabla_{\psi} \mathcal{L} = \frac{1}{N} \sum \nabla_{\psi} r_{\psi}(\tau_i) - \frac{1}{Z} \int p(\tau) \exp(r_{\psi}(\tau)) \nabla_{\psi} r_{\psi}(\tau) d\tau$$

$$= E_{\tau \sim \pi^*(\tau)} [\nabla_{\psi} r_{\psi}(\tau_i)] - E_{\tau \sim p(\tau | \sigma_{1:T}, \psi)} [\nabla_{\psi} r_{\psi}(\tau)]$$

\uparrow soft + optimal policy for current reward.

2.

$$\beta_t(s_t, a_t) = p(\sigma_{1:T} | s_t, a_t) \quad \text{"backward messages"}$$

$$\alpha_+(s_+) = p(s_+ | \sigma_{1:T-1}) \quad \text{can calculate with DP}$$

"forward messages"

Max Ent RL

3.

1. Given ψ , compute $\beta(s_t, a_t)$
2. Given ψ , β , compute $\alpha(s_+)$
3. $M_+(s_t, a_t) \propto \beta(s_t, a_t) \alpha(s_+)$
4. Eval $\nabla_{\psi} \mathcal{L} = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\psi} r_{\psi}(s_{i,t}, a_{i,t}) - \sum_t M_+(s_t, a_t) \nabla_{\psi} r_{\psi}(s_t, a_t)$
5. $\psi \leftarrow \psi + \eta \nabla_{\psi} \mathcal{L}$

Why "Max Ent"?

if $r_{\psi}(s_t, a_t) = \psi^T f(s_t, a_t)$

$$\underset{\psi}{\text{maximize}} \quad \underline{H(\pi^{\psi})} \quad \text{s.t.} \quad \underline{E_{\pi^{\psi}}[f]} = E_{\pi^*}[f]$$

3.

↑

3. Guided Cost Learning

Problem with Max Ent

1. solving for soft policy in inner loop
2. known dynamics

What if we can only sample?

$$\nabla_{\psi} \mathcal{L} = E_{\tau \sim \pi} [\nabla_{\psi} r_{\psi}(\tau)] - E_{\tau \sim p(\tau | \sigma_{\text{tr}}, \psi)} [\nabla_{\psi} r_{\psi}(\tau)]$$

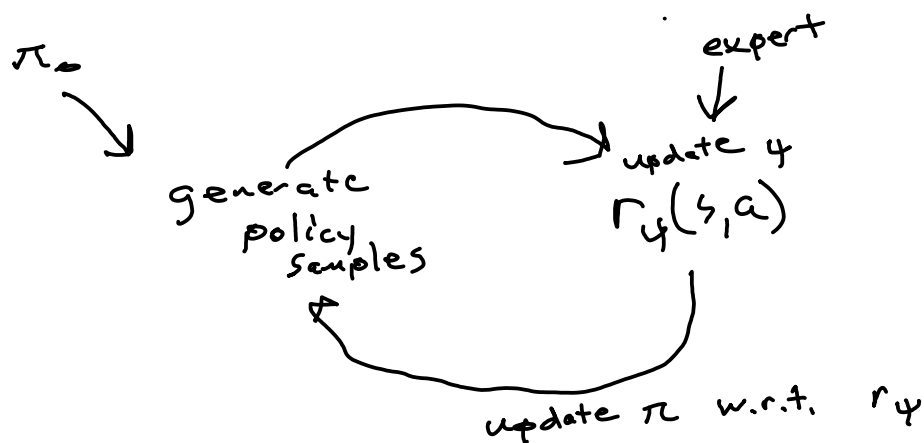
↑ solve with RL //

Do 1 step of RL

$$\approx - \frac{1}{M} \sum_{j=1}^M \nabla_{\psi} r_{\psi}(\tau_j)$$
$$\frac{1}{\sum w_j} \sum_j w_j \nabla_{\psi} r_{\psi}(\tau_j)$$

$$w_j = \frac{p(\tau) \exp(r_{\psi}(\tau_j))}{\pi(\tau_j)}$$

Guided Cost Learning



4 GAN

Looks like a game

policy tries to
make it harder
to distinguish
from Demos

reward learning tries
to make human policies
more likely

Generative Adversarial Network

Generator

policy

Discriminator

neural learner .