

Last Time
Sequential

This Time

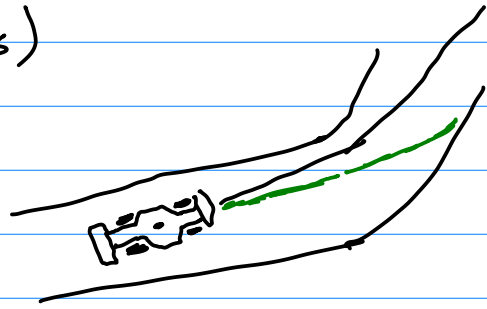
Imitation Learning
Inverse Reinforcement Learning

Behavioral Cloning

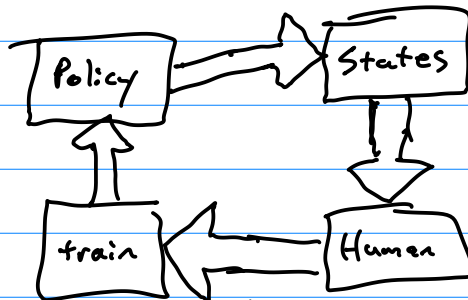
π_θ

$$\underset{\theta}{\text{maximize}} \prod_{(s,a) \in D} \pi_\theta(a|s)$$

Cascading Errors



Sequential Interactive Demonstration



Dagger "Dataset Aggregation"
SMILE

$$\pi^{(1)} = \pi^E$$

loop

execute $\pi^{(k)}$ to generate a dataset
query human response
train $\hat{\pi}^{(k)}$ with cloning
mix $\hat{\pi}^{(i)}$ with probability $\beta(1-\beta)^{i-1}$
to form $\hat{\pi}^{(k+1)}$

Apprenticeship Learning

Hard maneuver



Follow with controller

Expert Trajectories

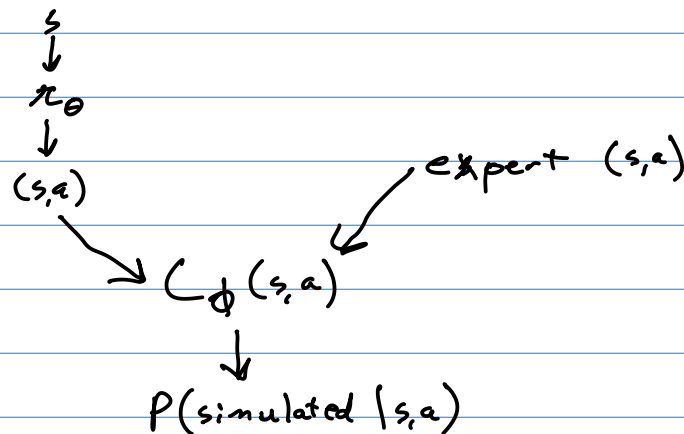
GAIL Generative Adversarial Imitation Learning

$\pi_\theta(a|s)$

$C_\phi(s,a)$

discriminator

$s,a \rightarrow C \rightarrow P(\text{simulated from } \pi_\theta)$



$$\max_{\phi} \min_{\theta} E_{(s,a) \sim \pi_\theta} [\log(C_\phi(s,a))] + E_{(s,a) \sim D} [\log(1 - C_\phi(s,a))]$$

\nwarrow from expert

Inverse Reinforcement Learning

forward DMU

Given (s, A, R, T)

Find π^*

Inverse RL

Given $s, A, T, \{\tau_i\}$

Find R

\nwarrow samples from π^*

Breakout Rooms

1	2	3
4	5	6
7	8	9

τ

1 →	1 →
2 →	2 ↓
3 ↓	5 →
6 ↓	6 ↓
9	9

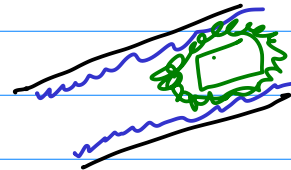
What is R ?

0	0	0
0	+10	0
0	-10	0

0	-1	1
-5	1	1
-5	-5	+10

Under specified

$$R_{\phi}(s, a) = \phi^T \beta(s, a)$$



Maximum Margin IRL

β binary

$$\|\phi\|_2 \leq 1$$

$$\begin{aligned}
 E[U(\tau)] &= E\left[\sum_{k=1}^{\tau} \gamma^{k-1} R_{\phi}(s^{(k)}, a^{(k)})\right] \\
 &= E\left[\sum \gamma^{k-1} \phi^T \beta(s^{(k)}, a^{(k)})\right] \\
 &= \phi^T E\left[\sum \gamma^{k-1} \beta(s^{(k)}, a^{(k)})\right] \\
 &= \phi^T \mu_{\pi} \quad \leftarrow \text{Feature Expectations}
 \end{aligned}$$

$\pi^{(1)} = \text{random}$

while $\epsilon \geq \epsilon$

$\mu^{(i)} \leftarrow$ feature expectations for $\pi^{(i)}$

$\phi^* \leftarrow$ maximize ϵ margin

subject to $\phi^T \mu_E \geq \phi^T \mu^{(i)} + \epsilon$ for $i=1, \dots, k-1$

$$\|\phi\| \leq 1$$

$\pi^{(i+1)} \leftarrow$ solve MDP with $R(s,a) = \phi^{*T} \bar{P}(s,a)$

minimize $\|\mu_E - \mu_\lambda\|_2$

$$\mu_\lambda = \sum \lambda_i \mu^{(i)}$$

subject to $\lambda \geq 0$

$$\|\lambda\|_1 = 1$$

Still Under-specified

Maximum Entropy IRL

Any policy introduces a distribution over trajectories
 $P_\pi(\tau)$

$$P_\phi(\tau) = \frac{1}{Z(\phi)} \exp(R_\phi(\tau))$$

$$R_\phi(\tau) = \sum_t \gamma^t \phi^T R(s(t), a(t))$$

$$Z(\phi) = \sum_{\tau} \exp(R_\phi(\tau))$$

normalization

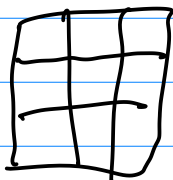
maximize $f(\phi)$

$$f(\phi) = \sum_{\tau \in D} \log P_\phi(\tau)$$

$$= \sum_{\tau \in D} \log \frac{1}{Z(\phi)} \exp(R_\phi(\tau))$$

$$= \left(\sum_{\tau \in D} R_\phi(\tau) \right) - |D| \log Z(\phi)$$

optimize this with gradient ascent



Takeaways

Imitation Learning matches actions of an expert

- Behavioral Cloning
- Dagger
- GAIL

Inverse Reinforcement Learning

- Learning R based on Expert Trajectories
- X Underspecified
- Maximum Entropy