

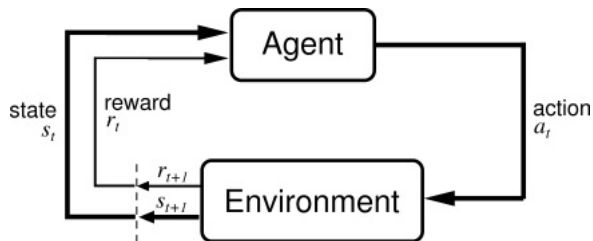
# End-to-End Differentiable Adversarial Imitation Learning

Shie Mannor  
EE @ Technion

Joint work with Nir Baram, Oron Anschel, Itai Caspi  
May 9th, ICRI

# Reinforcement Learning

(Agent interacts with an environment)



Applications:

1. Games (Go, Chess, Atari)
2. Robotics (self-driving cars)
3. Power grids, healthcare, ...

# Reinforcement Learning

(Different approaches)

## Value-based methods

$$Q^\pi(s, a) = \mathbb{E} \left[ \sum_{t=0}^N \gamma^t r_t \mid \pi, s, a \right]$$

$$\pi(s) \leftarrow \underset{a}{\operatorname{argmax}} Q(s, a)$$

## Policy-based Gradient

$$J(\theta) = \mathbb{E} \left[ \sum_{t=0}^N \gamma^t r_t \mid \pi_\theta \right]$$

Gradient Descent:  $\frac{\partial J}{\partial \theta}$

## Imitation Learning

$$D = \left\{ (s_0, a_0), (s_1, a_1), \dots \right\}$$

ERM:  $\mathcal{F} : s \rightarrow a$

# Why is imitation learning useful?

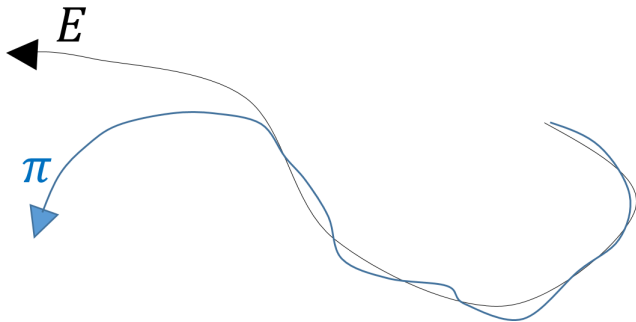
- ▶ Imitating is easy!
- ▶ Alleviate the temporal-credit assignment problem
- ▶ Learn complex behaviour quickly (dancing!)
- ▶ Combine simulator and data

But, need to use stochastic policies.

# Imitation Learning as a Supervised Problem

(Behavioral Cloning)

$$p_E(s) \neq p_\pi(s)$$



Impossible to avoid a stochastic policy.

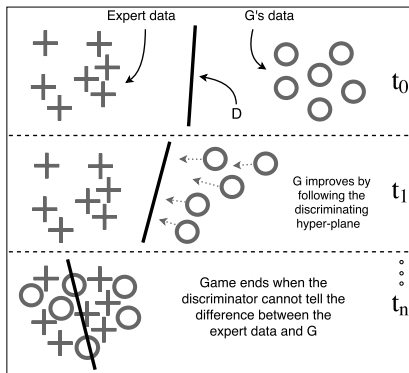
# Imitation Learning as a Reinforcement Problem

(Composition of two difficult problems)



# Imitation Learning using a GAN

(What are GANs?)



# Imitation Learning using GANs

( Generative Adversarial Imitation Learning)

## GAN

$$\operatorname{argmin}_G \operatorname{argmax}_{D \in (0,1)} \mathbb{E}_{x \sim p_E} [\log D(x)] + \mathbb{E}_{z \sim p_Z} [\log (1 - D(G(z)))]$$

A game between  $G$  (generator) and  $D$  (judge)

## GAIL

$$\operatorname{argmin}_\pi \operatorname{argmax}_{D \in (0,1)} \mathbb{E}_\pi [\log D(s, a)] + \mathbb{E}_{\pi_E} [\log(1 - D(s, a))] - \lambda H(\pi)$$

A game between  $\pi$  (policy) and  $D$  (judge)

Main technical issue:

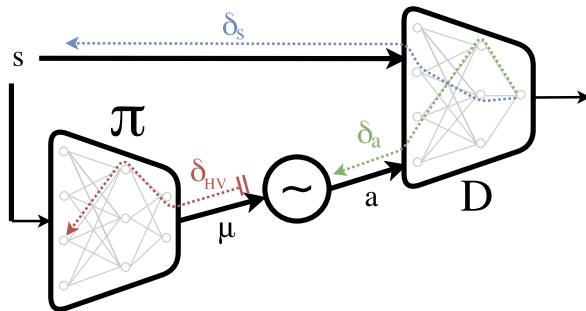
$$\mathbb{E}_\pi [\log D(s, a)] = \mathbb{E}_{s \sim \rho(\pi)} \mathbb{E}_{a \sim \pi(\cdot|s)} [\log D(s, a)].$$

How to differentiate w.r.t. (parameters of)  $\pi$ ?



# Imitation Learning using GAN: Model free approach

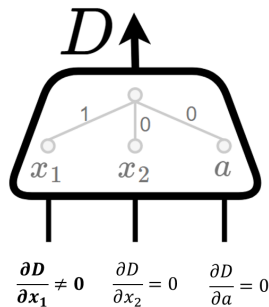
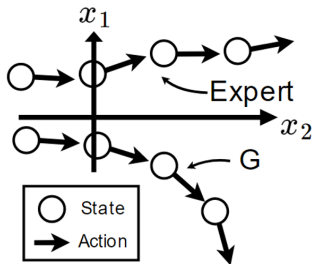
(The problem of training stochastic policies)



# Imitation Learning using GANs: Model free approach

(The importance of  $\frac{\partial D}{\partial s}$ )

Perfect discrimination based on  $x_1$ !



# Imitation Learning using GAN

(The role of  $D$  in imitation problems)

$D(s, a) = p(\pi|s, a)$  where  $\pi \in \{\pi_E, \pi\}$ .

$D(s, a)$  represents the likelihood ratio that the pair  $(s, a)$  is generated by  $\pi$  rather than by  $\pi_E$ .

Can show:

Policy likelihood ratio:  $\varphi(s, a) = \frac{p(a|s, \pi_E)}{p(a|s, \pi)}$ ,

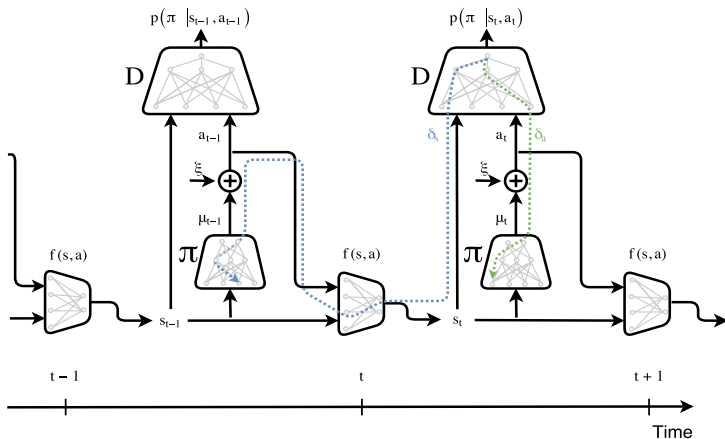
State distribution likelihood ratio:  $\psi(s) = \frac{p(s|\pi_E)}{p(s|\pi)}$ ,

$$D(s, a) = \frac{1}{1 + \varphi(s, a) \cdot \psi(s)}$$

Conclusion: Can easily compute  $\nabla_a D$  and  $\nabla_s D$ .

# model-based Imitation Learning using GAN

(How can we avoid throwing away  $\frac{\partial D}{\partial s}$ ,  $\frac{\partial D}{\partial a}$ )



# Main Conclusions

(The benefits of using a model)

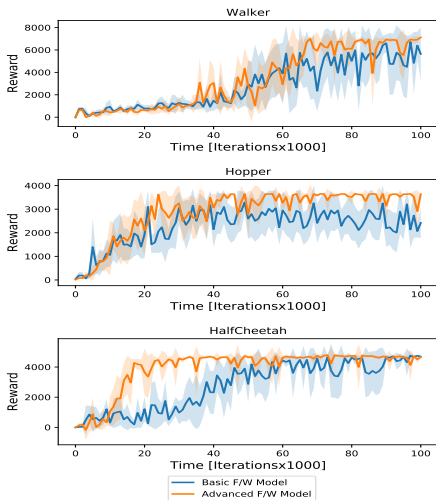
| Model Free (GAIL)  | Model Based (mGAIL)  |
|--|--|
| un-biased high-variance gradient<br>Discards $\frac{\partial D}{\partial s}$ | biased low-variance gradient<br>Uses $\frac{\partial D}{\partial s}$ |

# Some of the domains

Movie

# Results

(GAIL vs. mGAIL)



# Results

(GAIL vs. mGAIL)

| Task         | Dataset size | Behavioral cloning | GAIL                   | MGAIL                    |
|--------------|--------------|--------------------|------------------------|--------------------------|
| Cartpole     | 1            | 72.02 ± 35.82      | <b>200.00 ± 0.00</b>   | <b>200.00 ± 0.00</b>     |
|              | 4            | 169.18 ± 59.18     | <b>200.00 ± 0.00</b>   | <b>200.00 ± 0.00</b>     |
|              | 7            | 188.60 ± 29.61     | <b>200.00 ± 0.00</b>   | <b>200.00 ± 0.00</b>     |
|              | 10           | 177.19 ± 52.83     | <b>200.00 ± 0.00</b>   | <b>200.00 ± 0.00</b>     |
| Mountain Car | 1            | -136.76 ± 34.44    | <b>-101.55 ± 10.32</b> | <b>-107.4 ± 10.89</b>    |
|              | 4            | -133.25 ± 29.97    | <b>-101.35 ± 10.63</b> | <b>-100.23 ± 11.52</b>   |
|              | 7            | -127.34 ± 29.15    | <b>-99.90 ± 7.97</b>   | <b>-104.23 ± 14.31</b>   |
|              | 10           | -123.14 ± 28.26    | <b>-100.83 ± 11.40</b> | <b>-99.25 ± 8.74</b>     |
| Acrobot      | 1            | -130.60 ± 55.08    | <b>-77.26 ± 18.03</b>  | <b>-85.65 ± 23.74</b>    |
|              | 4            | -93.20 ± 35.58     | <b>-83.12 ± 23.31</b>  | <b>-81.91 ± 17.41</b>    |
|              | 7            | -96.92 ± 34.51     | <b>-82.56 ± 20.95</b>  | <b>-80.74 ± 14.02</b>    |
|              | 10           | -95.09 ± 33.33     | <b>-78.91 ± 15.76</b>  | <b>-77.93 ± 14.78</b>    |
| Hopper       | 4            | 50.57 ± 0.95       | 3614.22 ± 7.17         | <b>3669.53 ± 6.09</b>    |
|              | 11           | 1025.84 ± 266.86   | 3615.00 ± 4.32         | <b>3649.98 ± 12.36</b>   |
|              | 18           | 1949.09 ± 500.61   | 3600.70 ± 4.24         | <b>3661.78 ± 11.52</b>   |
|              | 25           | 3383.96 ± 657.61   | 3560.85 ± 3.09         | <b>3673.41 ± 7.73</b>    |
| Walker       | 4            | 32.18 ± 1.25       | 4877.98 ± 2848.37      | <b>6916.34 ± 115.20</b>  |
|              | 11           | 5946.81 ± 1733.73  | 6850.27 ± 91.48        | <b>7197.63 ± 38.34</b>   |
|              | 18           | 1263.82 ± 1347.74  | 6964.68 ± 46.30        | <b>7128.87 ± 141.98</b>  |
|              | 25           | 1599.36 ± 1456.59  | 6832.01 ± 254.64       | <b>7070.45 ± 30.68</b>   |
| Half-Cheetah | 4            | -493.62 ± 246.58   | 4515.70 ± 549.49       | <b>4891.56 ± 654.43</b>  |
|              | 11           | 637.57 ± 1708.10   | 4280.65 ± 1119.93      | <b>4844.61 ± 138.78</b>  |
|              | 18           | 2705.01 ± 2273.00  | 4749.43 ± 149.04       | <b>4876.34 ± 85.74</b>   |
|              | 25           | 3718.58 ± 1856.22  | 4840.07 ± 95.36        | <b>4989.95 ± 351.14</b>  |
| Ant          | 4            | 1611.75 ± 359.54   | 3186.80 ± 903.57       | <b>4645.12 ± 179.29</b>  |
|              | 11           | 3065.59 ± 635.19   | 3306.67 ± 988.39       | <b>4657.92 ± 94.27</b>   |
|              | 18           | 2579.22 ± 1366.57  | 3033.87 ± 1460.96      | <b>4664.44 ± 183.11</b>  |
|              | 25           | 3235.73 ± 1186.38  | 4132.90 ± 878.67       | <b>4637.52 ± 45.66</b>   |
| Humanoid     | 80           | 1397.06 ± 1057.84  | 10200.73 ± 1324.47     | <b>10312.34 ± 388.54</b> |
|              | 160          | 3655.14 ± 3714.28  | 10119.80 ± 1254.73     | <b>10428.39 ± 46.12</b>  |
|              | 240          | 5660.53 ± 3600.70  | 10361.94 ± 61.28       | <b>10470.94 ± 54.35</b>  |



# End-to-End Differentiable Adversarial Imitation Learning

- ▶ Imitation learning is the key to practical RL in many domains where expert trajectories are available
- ▶ A new architecture that allows differentiable backprop
- ▶ Main issue: the need to learn a forward model
- ▶ Transfer learning and multi-task are next