# Generative Adversarial Imitation Learning

Jonathan Ho, Stefano Ermon - NIPS 2016

Stanford University
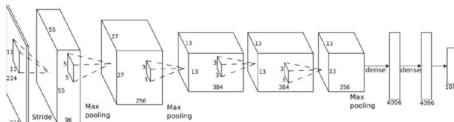
Presenter: Jinsung Yoon

# Problem: Imitation Learning

- Learning to perform a task from **demonstrations (Expert trajectories) without interaction** with the expert or access to reinforcement signal
- Applications: When the **rewards are hard to define**



$\mathbf{o}_t$ $\qquad$ $\pi_\theta(\mathbf{u}_t | \mathbf{o}_t)$ $\qquad$ $\mathbf{u}_t$

**Examples:**

- **Autonomous vehicles:**
  https://www.youtube.com/watch?v=cFtnflNe5fM
- **Medicine**

# Problem: Imitation Learning

Two main approaches
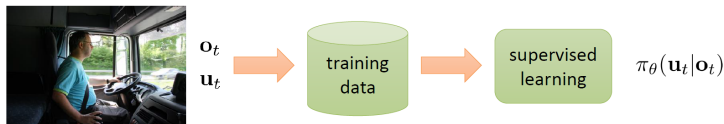
- **Behavioral cloning (BC)**
    - A supervised learning problem that maps state/action pairs to policy. **(State: feature, action: label)**
    - Requires a large number of expert trajectories **(high sample complexity)** - due to **compounding error** caused by covariate shift
    - Copies unnecessary actions as well
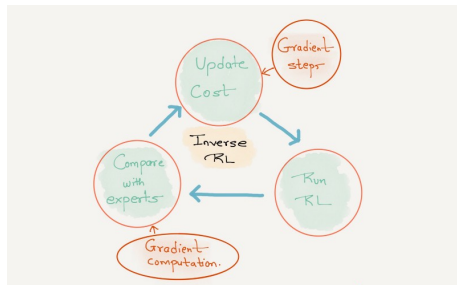- **Inverse Reinforcement Learning (IRL)**
    - Learns the **reward function** from expert trajectories that **prioritizes** entire trajectories over others, then derives the optimal policy
    - Expensive to run **(Inner loop has RL)**
    - Indirectly learns optimal policy from the reward function **(Using RL)**

# Problem: Imitation Learning

- **Behavioral cloning (BC)**



$$\mathbf{o}_t$$
$$\mathbf{u}_t$$ → training data → supervised learning → $\pi_\theta(\mathbf{u}_t | \mathbf{o}_t)$

- **Inverse Reinforcement Learning (IRL)**

## Proposed Model

**Generative Adversarial Imitation Learning** (GAIL)

- **Directly extracting a policy** from data as if it were obtained by RL following IRL.
- **Bypassing** any intermediate IRL step
- Draws an analogy between **imitation learning** and **generative adversarial networks (GAN)**
- Derive a model-free imitation learning algorithm with **significant performance improvement** with low sample and computational complexity.

# Background: Preliminaries on RL

- $\mathcal{S}$: Finite state space
- $\mathcal{A}$: Finite action space
- $\Pi$ : the set of all stochastic policies. Take action $\in \mathcal{A}$ given state $\in \mathcal{S}$
- $P(s'|s, a)$: Model dynamics
- $\pi \in \Pi$: A policy
- $\gamma$-discounted infinite horizon setting

$$\mathbb{E}_\pi[c(s, a)] = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t)]$$

- where $s_0 \sim p_0, a_t \sim \pi(\cdot|s_t), s_{t+1} \sim P(\cdot|s_t, a_t)$
- $\pi_E$: Expert policy

# Inverse Reinforcement Learning (IRL)

- **Maximum causal entropy IRL**

$$\max_{c \in \mathcal{C}} \left[ \min_{\pi \in \Pi} -H(\pi) + \mathbb{E}_\pi[c(s,a)] \right] - \mathbb{E}_{\pi_E}[c(s,a)]$$

- where $H(\pi) = [E]_\pi[-\log \pi(a|s)]$
- Try to find a cost function $c \in \mathcal{C}$ that assigns **low cost** to the expert policy $\pi_E$ and **high cost** to other poilces $(\pi)$
- Using **RL procedure**, we can find the expert policy based on the cost $c$

$$RL(c) = \arg\min_{\pi \in \Pi} -H(\pi) + \mathbb{E}_\pi[c(s,a)]$$

- Inner loop has **RL**; thus, slow

## Proposed Framework

- Use the **largest possible set of cost functions**
  $\mathcal{C} = \mathbb{R}^{\mathcal{S} \times \mathcal{A}} = \{c : \mathcal{S} \times \mathcal{A} \to \mathbb{R}\}$

- Use Gaussian processes or **Neural networks** to find the best cost function $c$ among large cost function class $\mathcal{C}$

- Avoid overfitting, we use a **"convex" cost function regularizer**
  $\psi : \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \to \mathbb{R} \cup \infty$

- With $\psi$, IRL procedure can be written as

$$\mathsf{IRL}_\psi(\pi_E) = \arg \max_{c \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}} -\psi(c) + \left[ \min_{\pi \in \Pi} -H(\pi) + \mathbb{E}_\pi[c(s,a)] \right] - \mathbb{E}_{\pi_E}[c(s,a)]$$

- Let $\tilde{c} \in \mathsf{IRL}_\psi(\pi_E)$, we are interested in $\pi = \mathsf{RL}(\tilde{c})$

# Occupancy Measure

- For a policy $\pi \in \Pi$, define its occupancy measure $\rho_\pi : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ as

$$\rho_\pi(s, a) = \pi(a|s) \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \pi)$$

- In words, occupancy measure is the **distribution** of state-action pairs with policy $\pi$

$$\mathbb{E}_\pi[c(s, a)] = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t)] = \sum_{s,a} \rho_\pi(s, a) c(s, a)$$

- The set of valid occupancy measures $\mathcal{D} = \{\rho_\pi : \pi \in \Pi\}$ can be written as

$$\mathcal{D} = \{\rho : \rho \geq 0 \,\&\, \sum_a \rho(s, a) = p_0(s) + \gamma \sum_{s',a} P(s|s', a) \rho(s', a) \forall s \in \mathcal{S}\}$$

- Note that there is **1-1 correspondence** between $\Pi$ and $\mathcal{D}$

# Occupancy Measure

- $\pi_\rho$ to denote the **unique policy for an occupancy measure** $\rho$
- **Convex conjugate:** for a function $f : \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \to \mathbb{R} \cup \infty$, its convex conjugate $f^* : \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \to \mathbb{R} \cup \infty$ is

$$f^*(x) = \sup_{y \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}} x^T y - f(y)$$

- Then, $\mathrm{RL}(\tilde{c})$ can be written as

$$\mathrm{RL} \odot \mathrm{IRL}_\psi(\pi_E) = \arg \min_{\pi \in \Pi} -H(\pi) + \psi^*(\rho_\pi - \rho_{\pi_E})$$

- Proof:

$$\arg \min_{\pi \in \Pi} -H(\pi) + \psi^*(\rho_\pi - \rho_{\pi_E})$$

$$= \arg \min_{\pi \in \Pi} \max_c -H(\pi) - \psi(c) + \sum_{s,a} \rho(s,a)c(s,a) - \sum_{s,a} \rho_{\pi_E}(s,a)c(s,a)$$

# Occupancy Measure

- The above theorem said that $\psi$-regularized IRL, implicitly seeks a **policy whose occupancy measure is close to the expert's** as measured by the convex function $\psi^*$

- It shows that various settings of $\psi$ lead to various imitation learning algorithms that **directly solve the optimization problem.** (without IRL and RL iteration)

- Therefore we can deduce the following things
  - IRL is a dual of an **occupancy measure matching problem**
  - The induced optimal policy is the **primal optimum.**

- Therefore, the traditional IRL definition (**finding a cost function that the expert poilicy is uniquely optimal**) changes to (**finding a policy that matches the expert's occupancy measure**)

# Generative Adversarial Imitation Learning (GAIL)

- Proposed $\psi$

$$\psi_{GA}(c) = \begin{cases} \mathbb{E}_{\pi_E}[g(c(s,a))], & \text{if } c < 0. \\ \infty, & \text{otherwise.} \end{cases}$$

where

$$g(x) = \begin{cases} -x - \log(1 - e^x), & \text{if } x < 0. \\ \infty, & \text{otherwise.} \end{cases}$$

- **Low penalty** when $x$ is far from 0. **High penalty** when $x$ is close to 0.

# $\psi_{GA}$

- The regularization function $\psi_{GA}$ is **motivated by the following fact.**

$$\psi_{GA}^*(\rho_\pi - \rho_{\pi_E}) = \max_{D \in (0,1)^{\mathcal{S} \times \mathcal{A}}} \mathbb{E}_\pi[\log(D(s,a))] + \mathbb{E}_{\pi_E}[\log(1 - D(s,a))]$$

- where the **maximum ranges over discriminative classifiers** is $D : \mathcal{S} \times \mathcal{A} \to (0,1)$

- The above equation is the **optimal negative log loss** of the binary classification problem of distinguishing between state-action pairs of $\pi$ and $\pi_E$.

- Therefore, the optimal loss is the **Jensen-Shannon divergence**

$$D_{JS}(\rho_\pi, \rho_{\pi_E}) = D_{KL}(\rho_\pi || (\rho_\pi + \rho_E)/2) + D_{KL}(\rho_{\pi_E} || (\rho_\pi + \rho_{\pi_E})/2)$$

# Proposed Optimization Problem

$$\min_{\pi} \psi_{GA}^*(\rho_\pi - \rho_{\pi_E}) - \lambda H(\pi) = D_{JS}(\rho_\pi, \rho_{\pi_E}) - \lambda H(\pi)$$

- It finds a policy whose **occupancy measure minimizes Jensen-Shannon divergence to the expert's.**

**Proposed Algorithm**

$$\min_{\pi} \max_{D \in (0,1)^{S \times A}} \mathbb{E}_\pi[\log(D(s,a))] + \mathbb{E}_{\pi_E}[\log(1 - D(s,a))] - \lambda H(\pi)$$

- We find a saddle point $(\pi, D)$

## Algorithm

$$\min_{\pi} \max_{D \in (0,1)^{\mathcal{S} \times \mathcal{A}}} \mathbb{E}_{\pi}[\log(D(s,a))] + \mathbb{E}_{\pi_E}[\log(1 - D(s,a))] - \lambda H(\pi)$$

- Initialize the **policy** $\pi_{\theta}$, and a **discriminator** $D_w : \mathcal{S} \times \mathcal{A} \to (0,1)$
- **Alternatively update** $\theta$ and $w$
    - **Adam for gradient step on** $w$ to increase
      $\mathbb{E}_{\pi}[\log(D(s,a))] + \mathbb{E}_{\pi_E}[\log(1 - D(s,a))]$
    - **TPRO step on** $\theta$ to decrease
      $\mathbb{E}_{\pi}[\log(D(s,a))] + \mathbb{E}_{\pi_E}[\log(1 - D(s,a))] - \lambda H(\pi)$
- Discriminator network is a **local cost function** providing learning signal to the policy.
- Taking a policy step that **decreases expected cost** w.r.t
  $c(s,a) = \log D(s,a)$

Try to find the policy $\pi \in \Pi$ that minimizes the cost function $c(s, a)$

$$L(\theta) = \mathbb{E}_{(s,a) \sim p}[c(s, a)]$$

Therefore, its gradient is

$$
\begin{aligned}
\nabla_\theta L(\theta) &= \nabla_\theta \int_{s,a} c(s, a) p(s, a) \\
&= \nabla_\theta \int_{s,a} c(s, a) \pi_\theta(a|s) p(s) \\
&= \int_{s,a} c(s, a) \nabla_\theta \pi_\theta(a|s) p(s) \\
&= \int_{s,a} c(s, a) \frac{\nabla_\theta \pi_\theta(a|s)}{\pi_\theta(a|s)} \pi_\theta(a|s) p(s) \\
&= \int_{s,a} c(s, a) \nabla_\theta \log(\pi_\theta(a|s)) \pi_\theta(a|s) p(s) \\
&= \mathbb{E}_{(s,a) \sim p}[c(s, a) \nabla_\theta \log(\pi_\theta(a|s))]
\end{aligned}
$$

# Algorithm

---

**Algorithm 1** Generative adversarial imitation learning

---

1: **Input:** Expert trajectories $\tau_E \sim \pi_E$, initial policy and discriminator parameters $\theta_0, w_0$
2: **for** $i = 0, 1, 2, \ldots$ **do**
3:     Sample trajectories $\tau_i \sim \pi_{\theta_i}$
4:     Update the discriminator parameters from $w_i$ to $w_{i+1}$ with the gradient

$$\hat{\mathbb{E}}_{\tau_i}[\nabla_w \log(D_w(s,a))] + \hat{\mathbb{E}}_{\tau_E}[\nabla_w \log(1 - D_w(s,a))] \qquad (17)$$

5:     Take a policy step from $\theta_i$ to $\theta_{i+1}$, using the TRPO rule with cost function $\log(D_{w_{i+1}}(s,a))$. Specifically, take a KL-constrained natural gradient step with

$$\hat{\mathbb{E}}_{\tau_i}\left[\nabla_\theta \log \pi_\theta(a|s)Q(s,a)\right] - \lambda \nabla_\theta H(\pi_\theta),$$
$$\text{where } Q(\bar{s}, \bar{a}) = \hat{\mathbb{E}}_{\tau_i}[\log(D_{w_{i+1}}(s,a)) \,|\, s_0 = \bar{s}, a_0 = \bar{a}] \qquad (18)$$

6: **end for**

---

# Experiments

**Experimental Settings** (Run on OpenAI Gym)

- **Low-dimensional** control tasks: (e.g. Cartpole, Acrobot)
- **High-dimensional** tasks: (e.g. 3D humanoid locomotion)

**Procedures**

- Generate expert behavior for these tasks by **running TRPO** on the **true cost functions** to create expert policies.
- Run GAIL and other benchmarks on the **generated expert policies**.
- Evaluate imitation performance w.r.t **sample complexity of expert data**.

**Benchmarks**

- Behavior Cloning
- Feature expectation matching (FEM): with **linear** cost function
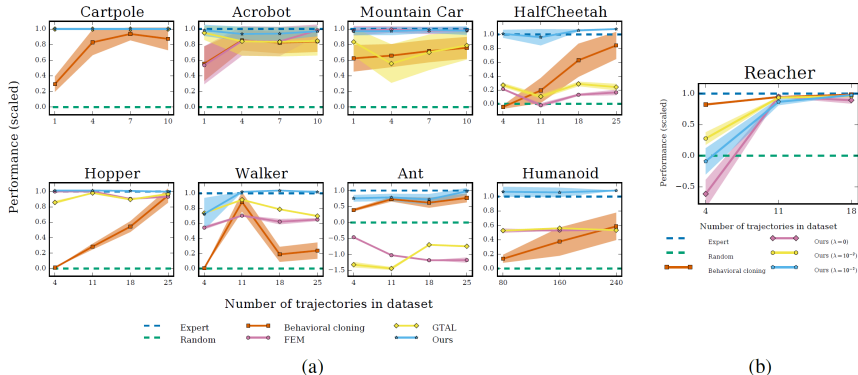- Game-theoretic apprenticeship learning (GTAL): with **convex** cost function

Figure 1: (a) Performance of learned policies. The $y$-axis is negative cost, scaled so that the expert achieves 1 and a random policy achieves 0. (b) Causal entropy regularization $\lambda$ on Reacher.

Table 3: Learned policy performance

| Task | Dataset size | Behavioral cloning | FEM | GTAL | Ours |
|---|---|---|---|---|---|
| Cartpole | 1 | $72.02 \pm 35.82$ | $200.00 \pm 0.00$ | $200.00 \pm 0.00$ | $200.00 \pm 0.00$ |
| | 4 | $169.18 \pm 59.81$ | $200.00 \pm 0.00$ | $200.00 \pm 0.00$ | $200.00 \pm 0.00$ |
| | 7 | $188.60 \pm 29.61$ | $200.00 \pm 0.00$ | $199.94 \pm 1.14$ | $200.00 \pm 0.00$ |
| | 10 | $177.19 \pm 52.83$ | $199.75 \pm 3.50$ | $200.00 \pm 0.00$ | $200.00 \pm 0.00$ |
| Acrobot | 1 | $-130.60 \pm 55.08$ | $-133.14 \pm 60.80$ | $-81.35 \pm 22.40$ | $-77.26 \pm 18.03$ |
| | 4 | $-93.20 \pm 32.58$ | $-94.21 \pm 47.20$ | $-94.80 \pm 46.08$ | $-83.12 \pm 23.31$ |
| | 7 | $-96.92 \pm 34.51$ | $-95.08 \pm 46.67$ | $-95.75 \pm 46.57$ | $-82.56 \pm 20.95$ |
| | 10 | $-95.09 \pm 33.33$ | $-77.22 \pm 18.51$ | $-94.32 \pm 46.51$ | $-78.91 \pm 15.76$ |
| Mountain Car | 1 | $-136.76 \pm 34.44$ | $-100.97 \pm 12.54$ | $-115.48 \pm 36.35$ | $-101.55 \pm 10.32$ |
| | 4 | $-133.25 \pm 29.97$ | $-99.29 \pm 8.33$ | $-143.58 \pm 50.08$ | $-101.35 \pm 10.63$ |
| | 7 | $-127.34 \pm 29.15$ | $-100.65 \pm 9.36$ | $-128.96 \pm 46.13$ | $-99.90 \pm 7.97$ |
| | 10 | $-123.14 \pm 28.26$ | $-100.48 \pm 8.14$ | $-120.05 \pm 36.66$ | $-100.83 \pm 11.40$ |
| HalfCheetah | 4 | $-493.62 \pm 246.58$ | $734.01 \pm 84.59$ | $1008.14 \pm 280.42$ | $4515.70 \pm 549.49$ |
| | 11 | $637.57 \pm 1708.10$ | $-375.22 \pm 291.13$ | $226.06 \pm 307.87$ | $4280.65 \pm 1119.93$ |
| | 18 | $2705.01 \pm 2273.00$ | $343.58 \pm 159.66$ | $1084.26 \pm 317.02$ | $4749.43 \pm 149.04$ |
| | 25 | $3718.58 \pm 1856.22$ | $502.29 \pm 375.78$ | $869.55 \pm 447.90$ | $4840.07 \pm 95.36$ |
| Hopper | 4 | $50.57 \pm 0.95$ | $3571.98 \pm 6.35$ | $3065.21 \pm 147.79$ | $3614.22 \pm 7.17$ |
| | 11 | $1025.84 \pm 266.86$ | $3572.30 \pm 12.03$ | $3502.71 \pm 14.54$ | $3615.00 \pm 4.32$ |
| | 18 | $1949.09 \pm 500.61$ | $3230.68 \pm 4.58$ | $3201.05 \pm 6.74$ | $3600.70 \pm 4.24$ |
| | 25 | $3383.96 \pm 657.61$ | $3331.05 \pm 3.55$ | $3458.82 \pm 5.40$ | $3560.85 \pm 3.09$ |
| Walker | 4 | $32.18 \pm 1.25$ | $3648.17 \pm 327.41$ | $4945.90 \pm 65.97$ | $4877.98 \pm 2848.37$ |
| | 11 | $5946.81 \pm 1733.73$ | $4723.44 \pm 117.18$ | $6139.29 \pm 91.48$ | $6850.27 \pm 39.19$ |
| | 18 | $1263.82 \pm 1347.74$ | $4184.34 \pm 485.54$ | $5288.68 \pm 37.29$ | $6964.68 \pm 46.30$ |
| | 25 | $1599.36 \pm 1456.59$ | $4368.15 \pm 267.17$ | $4687.80 \pm 186.22$ | $6832.01 \pm 254.64$ |
| Ant | 4 | $1611.75 \pm 359.54$ | $-2052.51 \pm 49.41$ | $-5743.81 \pm 723.48$ | $3186.80 \pm 903.57$ |
| | 11 | $3065.59 \pm 635.19$ | $-4462.70 \pm 53.84$ | $-6252.19 \pm 409.42$ | $3306.67 \pm 988.39$ |
| | 18 | $2597.22 \pm 1366.57$ | $-5148.62 \pm 37.80$ | $-3067.07 \pm 177.20$ | $3033.87 \pm 1460.96$ |
| | 25 | $3235.73 \pm 1186.38$ | $-5122.12 \pm 703.19$ | $-3271.37 \pm 226.66$ | $4132.90 \pm 878.67$ |
| Humanoid | 80 | $1397.06 \pm 1057.84$ | $5093.12 \pm 583.11$ | $5096.43 \pm 24.96$ | $10200.73 \pm 1324.47$ |
| | 160 | $3655.14 \pm 3714.28$ | $5120.52 \pm 17.07$ | $5412.47 \pm 19.53$ | $10119.80 \pm 1254.73$ |
| | 240 | $5660.53 \pm 3600.70$ | $5192.34 \pm 24.59$ | $5145.94 \pm 21.13$ | $10361.94 \pm 61.28$ |

| Task | Dataset size | Behavioral cloning | Ours ($\lambda = 0$) | Ours ($\lambda = 10^{-3}$) | Ours ($\lambda = 10^{-2}$) |
|---|---|---|---|---|---|
| Reacher | 4 | $-10.97 \pm 7.07$ | $-67.23 \pm 88.99$ | $-32.37 \pm 39.81$ | $-46.72 \pm 82.88$ |
| | 11 | $-6.23 \pm 3.29$ | $-6.06 \pm 5.36$ | $-6.61 \pm 5.11$ | $-9.26 \pm 21.88$ |
| | 18 | $-4.76 \pm 2.31$ | $-8.25 \pm 21.99$ | $-5.66 \pm 3.15$ | $-5.04 \pm 2.22$ |

# References

- **Paper Link:** `https://arxiv.org/pdf/1606.03476.pdf` - **2016 NIPS paper**
- **Useful blog links:**
  - `https://medium.com/@sanketgujar95/generative-adversarial-imitation-learning-266f45634e60`
  - `https://hollygrimm.com/rl_gail`
- **Code links:**
  - `https://github.com/hollygrimm/gail-mujoco`
  - `https://github.com/andrewliao11/gail-tf`
- **Youtube links:**
  - **Imitation learning tutorial:** `https://www.youtube.com/watch?v=6rZTaboSY4k`
  - **Author presentation:** `https://www.youtube.com/watch?v=bcnCo9RxhB8`