

Imitation Learning from Imperfect Demonstration

Yueh-Hua Wu^{1,2}, Nontawat Charoenphakdee^{3,2}, Han Bao^{3,2},
Voot Tangkaratt², Masashi Sugiyama^{2,3}

¹National Taiwan University

²RIKEN Center for Advanced Intelligence Project

³The University of Tokyo

Poster #47

- Imitation learning
 - learning from **demonstration** instead of a reward function
- Demonstration
 - a set of **decision makings** (state-action pairs x)
- Collected demonstration may be **imperfect**
 - Driving: traffic violation
 - Playing basketball: technical foul

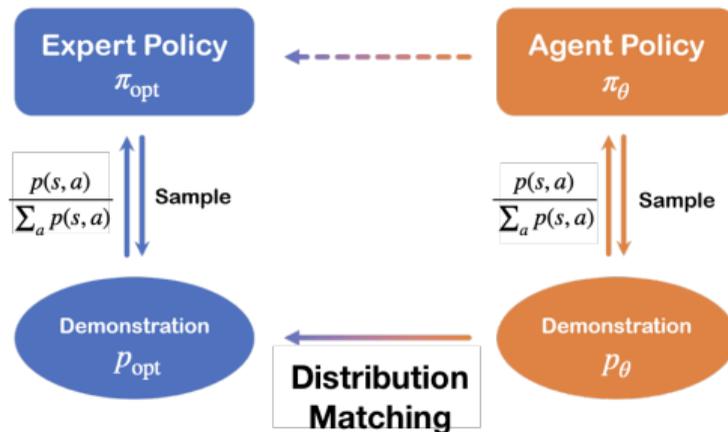
- **Confidence**: how optimal is state-action pair x (between 0 and 1)
- A **semi-supervised** setting: demonstration **partially** equipped with **confidence**
- How?
 - **crowdsourcing**: $N(1)/(N(1) + N(0))$.
 - **digitized score**: $0.0, 0.1, 0.2, \dots, 1.0$

Generative Adversarial Imitation Learning [1]

- **One-to-one correspondence** between **the policy** π and **the distribution of demonstration** [2]
- Utilize **generative adversarial training**

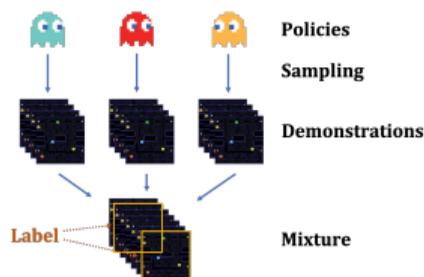
$$\min_{\theta} \max_w \mathbb{E}_{x \sim p_{\theta}} [\log D_w(x)] + \mathbb{E}_{x \sim p_{\text{opt}}} [\log(1 - D_w(x))]$$

D_w : discriminator, p_{opt} : demonstration distribution of π_{opt} , and p_{θ} : trajectory distribution of agent π_{θ}



Problem Setting

Human switches to **non-optimal policies** when they **make mistakes** or **are distracted**



$$p(x) = \underbrace{\alpha p(x|y = +1)}_{p_{\text{opt}}(x)} + (1 - \alpha) \underbrace{p(x|y = -1)}_{p_{\text{non}}(x)}$$

- **Confidence:** $r(x) \triangleq \Pr(y = +1|x)$
- **Unlabeled demonstration:** $\{x_i\}_{i=1}^{n_u} \sim p$
- **Demonstration with confidence:** $\{(x_j, r_j)\}_{j=1}^{n_c} \sim q$

Proposed Method 1: Two-Step Importance Weighting Imitation Learning

Step 1: estimate confidence by learning a confidence scoring function g

- Unbiased risk estimator (come to Poster #47 for details):

$$R_{\text{SC},\ell}(g) = \underbrace{\mathbb{E}_{x,r \sim q}[r \cdot (\ell(g(x)))]}_{\text{Risk for optimal}} + \underbrace{\mathbb{E}_{x,r \sim q}[(1-r)\ell(-g(x))]}_{\text{Risk for non-optimal}}$$

Theorem

For $\delta \in (0, 1)$, with probability at least $1 - \delta$ over repeated sampling of data for training \hat{g} ,

$$R_{\text{SC},\ell}(\hat{g}) - R_{\text{SC},\ell}(g^*) = \mathcal{O}_p\left(\underbrace{n_c^{-1/2}}_{\# \text{ of confidence}} + \underbrace{n_u^{-1/2}}_{\# \text{ of unlabeled}} \right)$$

Step 2: employ importance weighting to reweight GAIL objective

- Importance weighting

$$\min_{\theta} \max_w \mathbb{E}_{x \sim p_{\theta}}[\log D_w(x)] + \mathbb{E}_{x \sim p}\left[\frac{\hat{r}(x)}{\alpha} \log(1 - D_w(x))\right]$$

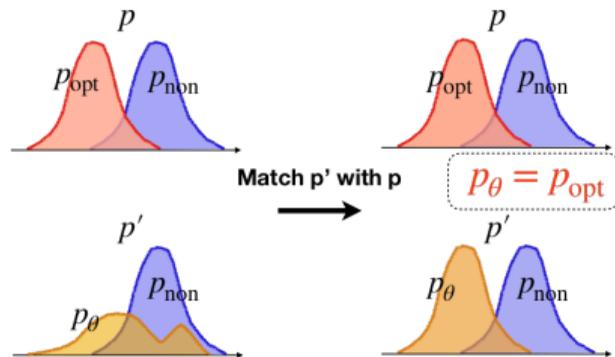
Proposed Method 2: GAIL with Imperfect Demonstration and Confidence

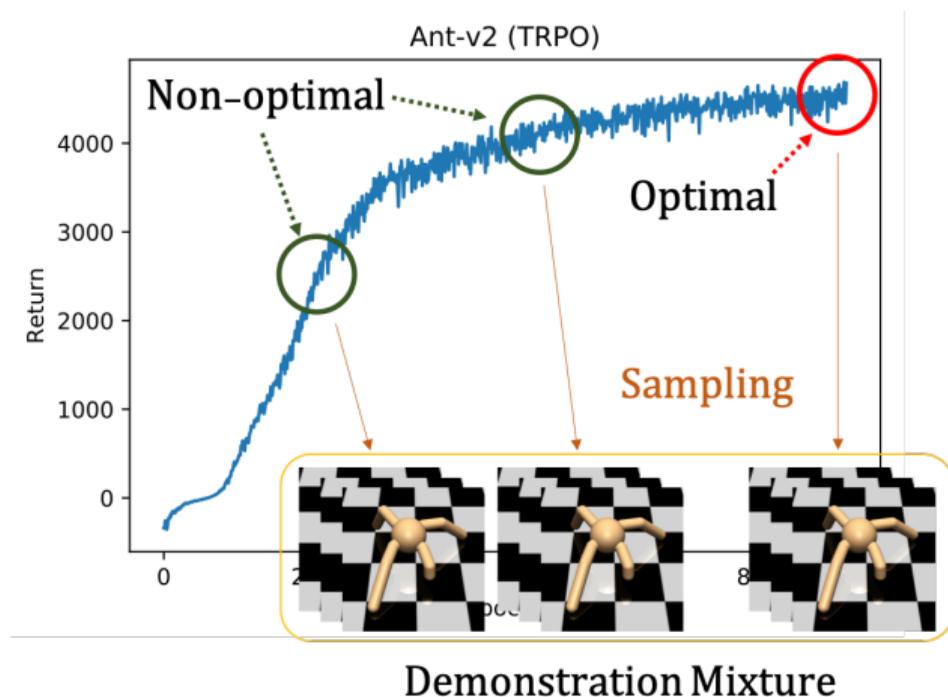
- Mix **the agent demonstration** with **the non-optimal one**

$$p' = \alpha p_{\theta} + (1 - \alpha) p_{\text{non}}$$

- Matching p' with p enables $p_{\theta} = p_{\text{opt}}$ and meanwhile **benefits from the large amount of unlabeled data.**
- Objective:

$$V(\theta, D_w) = \underbrace{\mathbb{E}_{x \sim p} [\log(1 - D_w(x))]}_{\text{Risk for P class}} + \underbrace{\alpha \mathbb{E}_{x \sim p_{\theta}} [\log D_w(x)] + \mathbb{E}_{x, r \sim q} [(1 - r) \log D_w(x)]}_{\text{Risk for N class}}$$





Confidence is given by a classifier trained with the demonstration mixture labeled as optimal ($y = +1$) and non-optimal ($y = -1$)

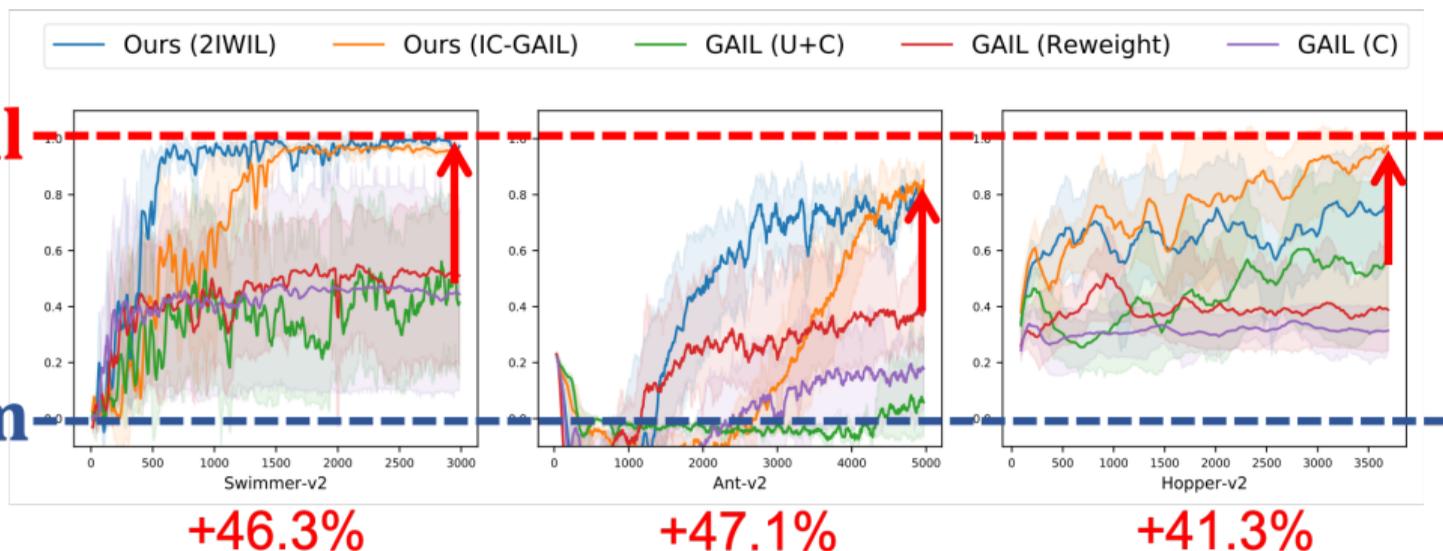
Results: Higher Average Return of the Proposed Methods

Environment: Mujoco

Proportion of labeled data: 20%

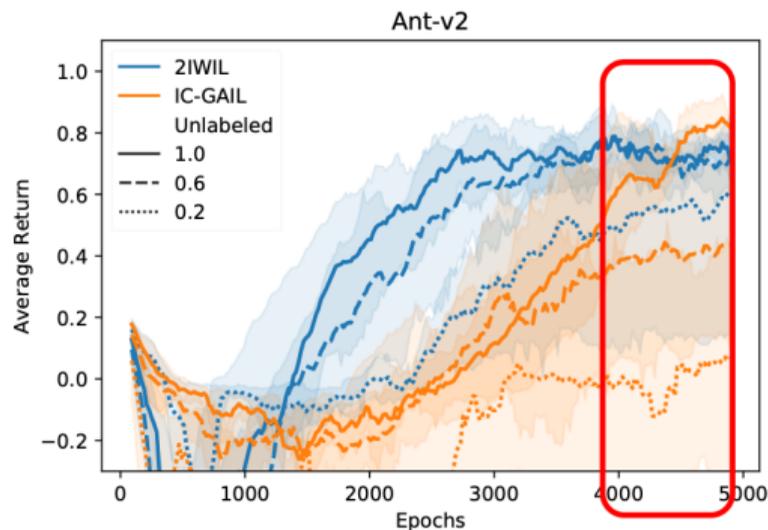
Optimal

Random

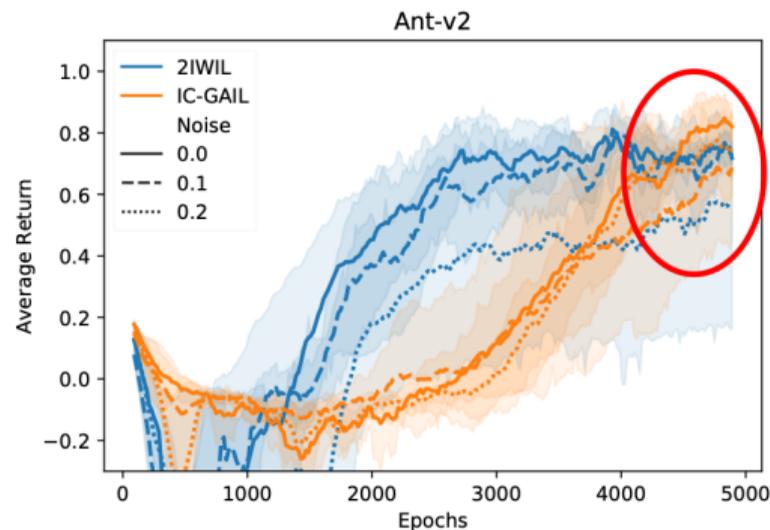


Results: Unlabeled Data Helps

- **More unlabeled data** results in **lower variance** and **better performance**
- proposed methods are **robust** to noise



(a) Number of unlabeled data. The number in the legend indicates **proportion** of original unlabeled data.



(b) Noise influence. The number in the legend indicates **standard deviation** of Gaussian noise.

- Two approaches that utilize **both unlabeled and confidence data** are proposed
- Our methods are **robust to labelers with noise**
- The proposed approaches can be generalized to other IL and IRL methods

Poster #47

- [1] Ho, Jonathan, and Stefano Ermon. "Generative adversarial imitation learning." Advances in Neural Information Processing Systems. 2016.
- [2] Syed, Umar, Michael Bowling, and Robert E. Schapire. "Apprenticeship learning using linear programming." Proceedings of the 25th international conference on Machine learning. ACM, 2008.