

PN-GAIL: LEVERAGING NON-OPTIMAL INFORMATION FROM IMPERFECT DEMONSTRATIONS

Qiang Liu, Huiqiao Fu, Kaiqiang Tang & Chunlin Chen*

School of Management and Engineering

Nanjing University

Nanjing, China

{qiangliu, hqfu, kqtang}@smail.nju.edu.cn, clchen@nju.edu.cn

Daoyi Dong

The Australian Artificial Intelligence Institute

University of Technology Sydney

Sydney, Australia

daoyidong@gmail.com

ABSTRACT

Imitation learning aims at constructing an optimal policy by emulating expert demonstrations. However, the prevailing approaches in this domain typically presume that the demonstrations are optimal, an assumption that seldom holds true in the complexities of real-world applications. The data collected in practical scenarios often contains imperfections, encompassing both optimal and non-optimal examples. In this study, we propose Positive-Negative Generative Adversarial Imitation Learning (PN-GAIL), a novel approach that falls within the framework of Generative Adversarial Imitation Learning (GAIL). PN-GAIL innovatively leverages non-optimal information from imperfect demonstrations, allowing the discriminator to comprehensively assess the positive and negative risks associated with these demonstrations. Furthermore, it requires only a small subset of labeled confidence scores. Theoretical analysis indicates that PN-GAIL deviates from the non-optimal data while mimicking imperfect demonstrations. Experimental results demonstrate that PN-GAIL surpasses conventional baseline methods in dealing with imperfect demonstrations, thereby significantly augmenting the practical utility of imitation learning in real-world contexts. Our codes are available at <https://github.com/QiangLiuT/PN-GAIL>.

1 INTRODUCTION

In recent years, Reinforcement Learning (RL) has achieved significant success in addressing sequential decision-making problems (Sutton & Barto, 2018; Xia et al., 2020; Zha et al., 2021). Its primary goal is to optimize policies to maximize cumulative rewards. However, designing an appropriate reward function can be quite challenging; a poorly designed reward function can lead to suboptimal performance of RL agents. In contrast, Imitation Learning (IL) presents a more practical approach, as it learns solely from demonstrations, eliminating the need for explicitly defined reward functions. Generative Adversarial Imitation Learning (GAIL) (Ho & Ermon, 2016), which employs the framework of Generative Adversarial Networks (GANs) (Goodfellow et al., 2014), directly learns a policy from demonstrations. Following the development of GAIL, many variants have been proposed to enhance algorithmic performance across different problem domains (Li et al., 2017; Fu et al., 2018; Dadashi et al., 2020; Fu et al., 2024).

The imitation learning methods mentioned above can learn an optimal policy given optimal demonstrations. However, most imitation learning methods tend to fail when faced with data filled with imperfect demonstrations. Especially in the real world, the assumption that the provided demonstrations are of high quality may not always be valid (Yang et al., 2024). For instance, due to factors

*Corresponding author

such as fatigue and distractions, decisions made by human experts may not always be optimal. In such cases, simply assigning equal weight to all data can lead to a decrease in the quality of the learned policy. Therefore, we need a method that can extract useful information from imperfect demonstrations to learn an optimal policy.

Existing methods for imitation learning from imperfect demonstrations can be broadly divided into two categories: weighting-based methods (Wu et al., 2019; Wang et al., 2021b;a; Tangkaratt et al., 2020; Zhang et al., 2021; Wang et al., 2023) and ranking-based methods (Brown et al., 2019; 2020; Chen et al., 2021; Huo et al., 2023; Taranovic et al., 2022). Weighting-based methods achieve imitation of optimal demonstrations through reweighting different demonstrations, while ranking-based methods aim to guide the recovery of the reward function with additional ranking information, thereby learning an optimal policy based on the rewards. In contrast, weighting-based methods are more computationally efficient since they do not require trajectory sorting. Additionally, they are more flexible to use, as they do not necessitate demonstrations to be in a trajectory form.

In order to solve the problem of learning from imperfect demonstrations using GAIL, Wu et al. (2019) proposed two methods: two-step importance weighting IL (2IWIL) and generative adversarial IL with imperfect demonstration and confidence (IC-GAIL). The former trains a classifier to forecast confidence scores and subsequently proceeds with weighted imitation learning, employing a two-step learning approach. The latter introduces an end-to-end learning method but at a slower pace of learning. However, as discussed in Section 4.1, 2IWIL is susceptible to the influence of preferences inherent in imperfect demonstrations during training. In the learning process of the discriminator, 2IWIL tends to assign a higher “reward” to the state-action pair with a greater probability of occurrence in imperfect demonstrations. This discrepancy in “rewards” diverges from our intended objectives, potentially resulting in the acquisition of a suboptimal policy.

To tackle the aforementioned challenge, we propose a new method, Positive-Negative Generative Adversarial Imitation Learning (PN-GAIL), building upon the framework of GAIL. Different from 2IWIL, we leverage non-optimal information from imperfect demonstrations, enabling the discriminator to weigh both positive and negative risks of imperfect demonstrations comprehensively and requiring only a small subset of labeled confidence scores. In this way, it can provide more accurate rewards for subsequent RL methods. Theoretical analysis reveals that PN-GAIL not only mimics imperfect demonstrations but also avoids imitating non-optimal ones, illustrating the ability of PN-GAIL to learn an optimal policy. Additionally, to get more accurate confidence scores, we propose an improved semi-supervised confidence classifier. Experiments on six control tasks are conducted to show the efficiency of our method in dealing with imperfect demonstrations compared to baseline methods. In particular, the main contributions of this work are threefold:

1. We propose a new method called PN-GAIL, which can leverage non-optimal information to learn an optimal policy from imperfect demonstrations.
2. We theoretically analyze the output of the optimal discriminator in PN-GAIL, demonstrating that PN-GAIL learns an optimal policy by deviating from the non-optimal demonstrations.
3. We demonstrate the efficiency of our method across six control tasks, with results showing superior performance compared to other baseline methods.

2 RELATED WORK

Imitation Learning Imitation learning methods can learn an optimal policy when given optimal demonstrations. Behavior cloning (BC) (Pomerleau, 1988) learns policies directly through a supervised learning paradigm and is mostly used in autonomous driving tasks (Hawke et al., 2020). While straightforward, it suffers from compounded errors due to covariate shift (Ross & Bagnell, 2010) and typically demands extensive data for effective training. Inverse Reinforcement Learning (IRL) (Abbeel & Ng, 2004; Ziebart et al., 2008) first seeks to recover the underlying reward function and then learns a policy through RL. On the other hand, GAIL views an imitation learning problem through the lens of occupancy measures (Puterman, 2014), and can learn a policy directly from the demonstrations. GAIL has demonstrated success across various imitation tasks, including multi-agent scenarios (Song et al., 2018), robot control (Peng et al., 2021), human motion simulation (Wei et al., 2021), and imitation of driver behavior (Bhattacharyya et al., 2022; Ruan & Di, 2022). How-

ever, these methods presuppose access to optimal demonstrations. When provided with imperfect demonstrations, they struggle to learn a good policy.

Weighting-based imitation learning from imperfect demonstrations Weighting-based imitation learning from imperfect demonstrations learns an optimal policy by reweighting different demonstrations and amplifying the significance of the optimal ones. 2IWIL and IC-GAIL (Wu et al., 2019) first propose to reweight imitation learning based on confidence. WGAIL (Wang et al., 2021b) connects confidence with the agent policy and discriminator without requiring additional prior information on confidence. However, it needs a high proportion of optimal demonstrations in imperfect demonstrations. VILD (Tangkaratt et al., 2020) employs a variational method to jointly estimate demonstration quality and reward, but it assumes that the quality of demonstrations be correlated with variance. CAIL (Zhang et al., 2021) guides confidence estimation by introducing trajectory ranking. UID (Wang et al., 2023) treats imperfect demonstrations as unlabeled data, based on the idea of PU Learning (Du Plessis et al., 2014), mitigating the impact of non-optimal demonstrations. Nevertheless, this relies on the assumption that non-optimal demonstrations within the imperfect demonstrations can well match agent demonstrations. Additionally, some studies address imperfect demonstrations in offline imitation learning (Sasaki & Yamashina, 2020; Xu et al., 2022; Kim et al., 2021; Yu et al., 2023; Li et al., 2024). However, these methods either similarly assume that the proportion of the optimal demonstrations is dominant, or require an additional set of optimal demonstrations.

Ranking-based imitation learning from imperfect demonstrations Ranking-based imitation learning from imperfect demonstrations utilizes additional ranking information to guide the recovery of the reward function, thereby learning a policy based on the rewards. T-REX (Brown et al., 2019) infers the reward function from the given ranking trajectories and expects the reward function to conform to the given ranking order. However, this approach demands a substantial quantity of ranking trajectories to enhance its generalization capacity. D-REX (Brown et al., 2020) automatically generates ranking trajectories by introducing varying degrees of noise. SSRR (Chen et al., 2021) revises the structure of the reward function in D-REX to accommodate different levels of noise influence better. LERP (Huo et al., 2023) views suboptimal demonstrations as additive noise on the reward function, establishing a quantifiable relationship between noise and reward based on D-REX. However, the automatic generation of the ranking trajectories requires the assumption that the trajectory will receive lower rewards with the addition of noise, which is not necessarily true in cases where random demonstrations exist. AILP (Taranovic et al., 2022) necessitates the teacher’s access to the true reward function, thereby providing real-time correct ranking between two trajectories. Nevertheless, this condition is challenging to meet in practice.

3 PRELIMINARIES

In this section, we provide a brief background on RL, GAIL, and 2IWIL.

Reinforcement learning We consider the standard Markov Decision Process (MDP) (Sutton & Barto, 2018). An MDP typically comprises six components, denoted as $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \rho_0, \gamma \rangle$, where \mathcal{S} is the state space, \mathcal{A} is the action space, $\mathcal{P}(s_{t+1}|s_t, a_t)$ is the transition probability from state s_t and action a_t at time step t to state s_{t+1} at time step $t + 1$, $\mathcal{R}(s, a)$ is the reward function, ρ_0 is the distribution of initial states, and $\gamma \in (0, 1)$ stands for the discount factor. In an RL process, the agent aims to learn a policy $\pi(a|s)$ to maximize its expected discounted rewards $\mathbb{E}_{s_0 \sim \rho_0, \pi} [\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(s_t, a_t)]$. For any given policy π , there exists a corresponding occupancy measure $\rho_\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, establishing a one-to-one relationship between them.

GAIL and 2IWIL GAIL integrates GANs framework into imitation learning, leading the following min-max optimization problem by minimizing the Jensen-Shannon divergence between p_θ and p_E (Ke et al., 2021):

$$\min_{\theta} \max_w \mathbb{E}_{(s,a) \sim p_\theta} [\log D_w(s, a)] + \mathbb{E}_{(s,a) \sim p_E} [\log(1 - D_w(s, a))], \quad (1)$$

where p_θ and p_E are the corresponding normalized occupancy measures for the agent policy π_θ and the expert policy π_E , respectively. The discriminator D_w attempts to discern these distributions from

π_E and π_θ , while π_θ aims to “trick” the discriminator, thereby minimizing $\mathbb{E}_{(s,a)\sim p_\theta}[\log D_w(s,a)]$. Ultimately, the output of the discriminator, $-\log D_w(s,a)$, serves as a reward, which can then be utilized to learn the policy π_θ through RL methods such as TRPO (Schulman et al., 2015), PPO (Schulman et al., 2017) and SAC (Haarnoja et al., 2018).

Since GAIL assigns the same weights to all demonstrations, if the given demonstrations are non-optimal, then the learned policy will also be non-optimal. To address this issue, 2IWIL considers the following setup:

$$\begin{aligned}\mathcal{D}_c &\triangleq \{(x_{c,i}, r_i)\}_{i=1}^{n_c} \stackrel{\text{i.i.d.}}{\sim} q(x,r), \\ \mathcal{D}_u &\triangleq \{x_{u,i}\}_{i=1}^{n_u} \stackrel{\text{i.i.d.}}{\sim} p(x),\end{aligned}$$

where x is the state-action pair, r denotes confidence score, indicating the probability that x belongs to the optimal demonstrations, $q(x,r) = p(x)p_r(r|x)$ and $p_r(r_i|x) = \delta(r_i - r(x))$ is Dirac delta function. \mathcal{D}_c and \mathcal{D}_u represent confidence data and unlabeled data, respectively.

2IWIL first trains a probabilistic classifier, which forecasts the confidence scores of demonstrations in \mathcal{D}_u through *semi-conf (SC) classification*, leveraging the knowledge of confidence scores in \mathcal{D}_c . The probabilistic classifier is trained with the loss function as follows:

$$R_{SC,\ell}(g) = \mathbb{E}_{x,r\sim q} [r\ell(g(x)) + (1-r)\ell(-g(x)) - \beta\ell(-g(x))] + \mathbb{E}_{x\sim p}[\beta\ell(-g(x))], \quad (2)$$

where g is a prediction function, ℓ is a loss function which uses logistic loss and $\beta = \frac{n_u}{n_c+n_u}$. After obtaining confidence scores for all demonstrations, 2IWIL uses Bayes’ rule to reweight the GAIL objective. The final objective becomes

$$\min_{\theta} \max_w \mathbb{E}_{x\sim p_\theta} [\log D_w(x)] + \mathbb{E}_{x\sim p} \left[\frac{r(x)}{\eta} \log(1 - D_w(x)) \right], \quad (3)$$

where η is a class-prior, denoting the proportion of optimal demonstrations within the imperfect demonstrations, and p is the corresponding normalized occupancy measures for \mathcal{D}_c and \mathcal{D}_u .

4 APPROACH

In this section, we begin by elucidating the motivation behind our method. We illustrate the problem of 2IWIL through an example and then introduce our method PN-GAIL with theoretical analysis. Details of derivations and proofs in this section can be found in Appendix B.

4.1 MOTIVATION

2IWIL aims to reweight demonstrations based on confidence, assigning greater weights to those with high confidence so that the discriminator can give higher rewards. However, it is worth noting that this weighting behavior can be influenced by the preferences inherent in imperfect demonstrations. As shown in Fig. 1, the top half of the graph represents the actual confidence scores, while the bottom half represents the equivalent weights during the discriminator’s training. When imperfect demonstrations favor a low-confidence state-action pair, we consider the goals of GAIL, 2IWIL:

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

GAIL assigns the same weights to all given demonstrations, which means $\frac{r(x)}{\eta} \equiv 1.0$. For the given expert demonstrations, we only need to consider the second term of the above equation, which can be expanded as: $\sum p(x) \frac{r(x)}{\eta} \log(1 - D_w(x))$. Here, the coefficient of $\log(1 - D_w(x))$ is $\frac{r(x)}{\eta}$. Therefore, if there is a higher probability of x_1 appearing in imperfect demonstrations, e.g., $p(x_1) = 5p(x_{other})$ (assuming that the probabilities of other state-action pairs are the same), then, for x_1 , the coefficient of $\log(1 - D_w(x))$ is $p(x_1) \frac{r(x_1)}{\eta} = p(x_{other}) \frac{5r(x_1)}{\eta}$. This can also be explained as that the probability of x_1 appearing is the same as for other demonstrations, but the confidence score is 5 times higher, since η is constant across all demonstrations. In the case of GAIL, since $\frac{r(x)}{\eta} \equiv 1.0$, the confidence score becomes 5 times of the original, which is calculated as $1.0 \times 5 = 5.0$. This

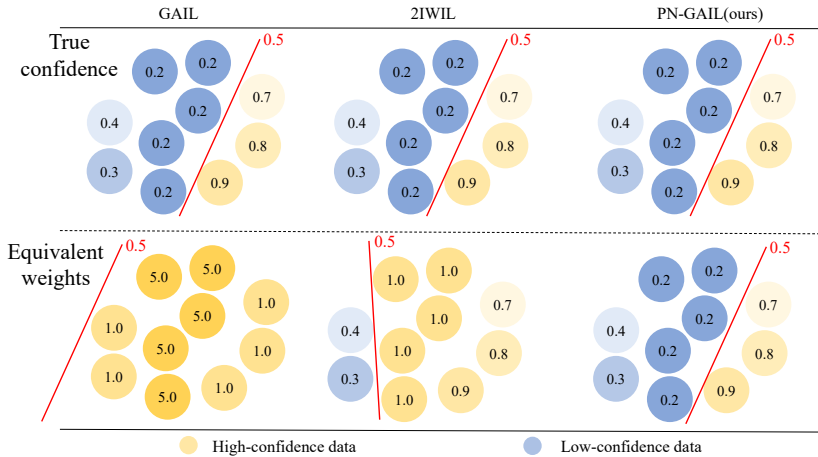


Figure 1: Schematic diagram of the difference between PN-GAIL, 2IWIL and GAIL. The top half of the graph is the actual confidence score, the bottom half is the equivalent weight when training the discriminator, and the red line is distinguished by a threshold of 0.5.

can lead to low-confidence data being treated as high-confidence data, not aligning with the actual situation. In addition, for a clearer explanation, we also provide a simple example.

Suppose a state s_1 has two actions $x_1(s_1, a_1)$ and $x_2(s_1, a_2)$. In Fig. 1, the circle with a confidence score of 0.8 represents x_2 , and the five circles with a confidence score of 0.2 all represent x_1 , which means $p(x_1) = 5p(x_2)$, indicating a higher probability of x_1 occurring in imperfect demonstrations. It is clear that x_2 is better than x_1 . However, in imperfect demonstrations where $p(x_1) = 5p(x_2)$ and the prior η is the same, according to Eq. (3), the equivalent weight of x_1 will be 1.0 compared to x_2 ($0.2 \times 5 = 1.0$). This means that the discriminator will consider x_1 to be more likely the optimal demonstration than x_2 , resulting in a poor policy.

4.2 POSITIVE-NEGATIVE GENERATIVE ADVERSARIAL IMITATION LEARNING

To tackle the problem above, we propose Positive-Negative Generative Adversarial Imitation Learning (PN-GAIL). This method leverages non-optimal information from imperfect demonstrations, allowing the discriminator to comprehensively assess the positive and negative risks associated with these demonstration. By doing so, it mitigates the influence of preferences inherent in imperfect demonstrations on the discriminator, thus ensuring that its evaluations better reflect actual conditions. This, in turn, provides more accurate rewards for the subsequent RL process, leading to the learning of a better policy.

We begin by focusing on the training of the discriminator, denoting optimal demonstrations as positive examples and non-optimal demonstrations as negative examples. In 2IWIL, the discriminator only considers the positive risk of imperfect demonstrations, while ignoring negative risk. Therefore, the discriminator will heavily prioritize the positive risk training for state-action pairs frequently appearing in imperfect demonstrations, leading to incorrect results. For this reason, we aim to incorporate the negative risk into the training of the discriminator when dealing with imperfect demonstrations. Specifically, following Xu & Denil (2021), let (X, Y) represent the input and output of a binary classification problem, where X denotes the state-action pair and $Y \in \{0, 1\}$. We label optimal data as 0 and non-optimal data as 1. The imperfect demonstrations is denoted as \mathcal{D} , comprising \mathcal{D}_{opt} (optimal demonstrations) and \mathcal{D}_{non} (non-optimal demonstrations), where $\mathcal{D} = \mathcal{D}_{\text{opt}} + \mathcal{D}_{\text{non}}$. We aim to train a discriminator D_w using a loss function $\phi : \mathbb{R} \times \{0, 1\} \rightarrow \mathbb{R}$. Utilizing the labeled risk operator as follows:

$$R_{D_w}^y(\mathcal{D}) = \mathbb{E}_{\mathcal{D}}[\phi(D_w(x), y)]. \tag{4}$$

We expect the discriminator to provide accurate evaluation scores for both the dataset generated by the agent policy and the given imperfect demonstrations. To achieve this, we consider the risk associated with the dataset generated by the agent policy and the risk associated with the imperfect

demonstrations, respectively. The overall risk of the discriminator is

$$R_{D_w}^{pn}(\mathcal{D}_{\pi_\theta}, \mathcal{D}) = R_{D_w}^1(\mathcal{D}_{\pi_\theta}) + R_{D_w}^{pn}(\mathcal{D}), \quad (5)$$

where \mathcal{D}_{π_θ} is the demonstrations generated by agent policy π_θ . We can write the risk associated with the imperfect demonstrations as the sum of positive and negative risks:

$$R_{D_w}^{pn}(\mathcal{D}) = R_{D_w}^{pn}(\mathcal{D}_{\text{opt}}, \mathcal{D}_{\text{non}}) = \eta R_{D_w}^0(\mathcal{D}_{\text{opt}}) + (1 - \eta) R_{D_w}^1(\mathcal{D}_{\text{non}}), \quad (6)$$

where $\eta = p(y = 0)$ is a class-prior, denoting the proportion of optimal demonstrations within the imperfect demonstrations.

Based on Eq. (6), the overall risk of the discriminator can be rewritten as

$$R_{D_w}^{pn}(\mathcal{D}, \mathcal{D}_{\pi_\theta}) = R_{D_w}^1(\mathcal{D}_{\pi_\theta}) + \eta R_{D_w}^0(\mathcal{D}_{\text{opt}}) + (1 - \eta) R_{D_w}^1(\mathcal{D}_{\text{non}}). \quad (7)$$

Replacing the loss function with the standard logistic loss and tidying up the statement, the objective of the discriminator becomes

$$\max_w \mathbb{E}_{x \sim p_\theta} [\log D_w(x)] + \eta \mathbb{E}_{x \sim p_{\text{opt}}} [\log(1 - D_w(x))] + (1 - \eta) \mathbb{E}_{x \sim p_{\text{non}}} [\log D_w(x)]. \quad (8)$$

Since $r(x)$ denotes the probability that x belongs to the optimal demonstrations, which means $r(x) = p(y = 0|x)$ and $1 - r(x) = p(y = 1|x)$, according to the Bayes' rule we have

$$p_{\text{opt}}(x) = p(x|y = 0) = \frac{r(x)p(x)}{\eta}, \quad p_{\text{non}}(x) = p(x|y = 1) = \frac{(1 - r(x))p(x)}{1 - \eta}. \quad (9)$$

Then we can rewrite the objective of the discriminator in the following theorem.

Theorem 4.1. *Based on Eq. (9), the objective of the discriminator can be rewritten as*

$$\max_w \mathbb{E}_{x \sim p_\theta} [\log D_w(x)] + \mathbb{E}_{x \sim p} [r(x) \log(1 - D_w(x))] + \mathbb{E}_{x \sim p} [(1 - r(x)) \log D_w(x)]. \quad (10)$$

The agent receives a reward equivalent to $-\log D_w(x)$, and then the final objective to be optimized becomes

$$\min_\theta \max_w \mathbb{E}_{x \sim p_\theta} [\log D_w(x)] + \mathbb{E}_{x \sim p} [r(x) \log(1 - D_w(x))] + \mathbb{E}_{x \sim p} [(1 - r(x)) \log D_w(x)]. \quad (11)$$

Furthermore, recall that 2IWIL adopts a two-step learning approach, where \mathcal{D}_c and \mathcal{D}_u represent confidence data and unlabeled data, respectively. To get more accurate confidence scores, we refine the *semi-conf (SC) classification* proposed in 2IWIL, which is trained by minimizing the following risk:

$$R_{\text{SC}, \ell}(g) = \mathbb{E}_{x, r \sim q} [r \ell(g(x)) + (1 - r) \ell(-g(x)) - \beta \ell(-g(x))] + \mathbb{E}_{x \sim p} [\beta \ell(-g(x))]. \quad (12)$$

We note that for a state-action pair x occurring solely in \mathcal{D}_c , once $1 - r - \beta < 0$, where $r > 1 - \beta$, the coefficient of $\ell(-g(x))$ becomes negative. In order to minimize the risk, the classifier would then forecast $g(x)$ as positive infinity, leading to an excessively high estimation of confidence for demonstrations in \mathcal{D}_c . Concurrently, Eq. (12) tends to predict data in \mathcal{D}_u as negative, resulting in an underestimated confidence for demonstrations in \mathcal{D}_u . To balance this effect, we propose *balanced semi-conf (BSC) classification*. We introduce $\mathbb{E}_{x \sim p} [\alpha \ell(g(x))] - \mathbb{E}_{x \sim q} [\alpha \ell(g(x))]$, the theoretical value of which is 0 since \mathcal{D}_c and \mathcal{D}_u are drawn from the same distribution $p(x)$. The final risk is as follows:

$$R_{\text{BSC}, \ell}(g) = \mathbb{E}_{x, r \sim q} [r \ell(g(x)) + (1 - r) \ell(-g(x)) - \alpha \ell(g(x)) - \beta \ell(-g(x))] + \mathbb{E}_{x \sim p} [\alpha \ell(g(x)) + \beta \ell(-g(x))], \quad (13)$$

where the loss function ℓ uses the logistic loss. Next, similar to 2IWIL, we seek to derive the optimal values of α and β for minimizing the variance of the empirical unbiased estimator $\widehat{R}_{\text{BSC}, \ell}(g)$ through the following theorem.

Theorem 4.2. *Let d_1 denote $\text{Var}(\ell(-g(x)))$, d_2 denote $\text{Var}(\ell(g(x)))$, $\sigma_{\text{cov}1}$ denote the covariance between $\frac{1}{n_c} \sum_{i=1}^{n_c} r_i (\ell(g(x_{c,i})) - \ell(-g(x_{c,i})))$ and $\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i}))$, $\sigma_{\text{cov}2}$ denote the covariance between $\frac{1}{n_c} \sum_{i=1}^{n_c} (1 - r_i) (\ell(-g(x_{c,i})) - \ell(g(x_{c,i})))$ and $\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(g(x_{c,i}))$, cov denote $\text{Cov}(\ell(-g(x)), \ell(g(x)))$. The estimator $\widehat{R}_{\text{BSC}, \ell}(g)$ has the minimum variance when*

$$\alpha = \frac{n_u}{n_c + n_u} - \frac{d_1 \text{cov} - \text{cov}^2}{d_1 d_2 - \text{cov}^2} \frac{n_u}{n_c + n_u} + \frac{d_1 \sigma_{\text{cov}2} - \text{cov} \sigma_{\text{cov}1}}{d_1 d_2 - \text{cov}^2} \frac{n_c n_u}{n_c + n_u},$$

$$\beta = \frac{n_u}{n_c + n_u} - \frac{d_2 \text{cov} - \text{cov}^2}{d_1 d_2 - \text{cov}^2} \frac{n_u}{n_c + n_u} + \frac{d_2 \sigma_{\text{cov}1} - \text{cov} \sigma_{\text{cov}2}}{d_1 d_2 - \text{cov}^2} \frac{n_c n_u}{n_c + n_u}.$$

Since $d_1, d_2, \sigma_{\text{cov}1}, \sigma_{\text{cov}2}, \text{cov}$ are difficult to calculate, in practice, we assume that these covariances are sufficiently small for computational convenience. Consequently, we have $\alpha = \frac{n_u}{n_c + n_u}$ and $\beta = \frac{n_u}{n_c + n_u}$. During the training process, as we assume that the data from \mathcal{D}_c and \mathcal{D}_u are drawn from the same distribution $p(x)$, we guarantee this condition via the clip function (see more details in Appendix C.2).

4.3 THEORETICAL ANALYSIS

We consider the reward given by the optimal discriminator $D_w^*(x)$. In 2IWIL, when the discriminator is optimal, the reward is $-\log D_w^*(x) = \log((rp + \eta p_\theta) / (\eta p_\theta))$. Consequently, if imperfect demonstrations exhibit a pronounced preference for a certain state-action pair, it results in a significantly higher probability of p compared to other state-action pairs. The discriminator tends to provide an inflated reward, hindering the learning of an optimal policy. Conversely, in our method, we first give the following theorem:

Theorem 4.3. *Given a fixed agent policy π_θ , the optimal discriminator $D_w^*(x)$ of Eq. (11) can be written as*

$$D_w^*(x) = \frac{(1-r)p + p_\theta}{p + p_\theta}. \quad (14)$$

As a result, when the optimal discriminator $D_w^*(x)$ is given, the optimization of π_θ is equivalent to minimizing

$$2\text{JSD}(p_\theta||p) - \text{KL}(p_\theta||p_1) - (1-\eta)\text{KL}(p_{\text{non}}||p_1) + C, \quad (15)$$

where $p_1 = (p_\theta + (1-\eta)p_{\text{non}})/(2-\eta)$, $C = \eta\mathbb{E}_{x \sim p_{\text{opt}}} \left[\log \frac{\eta p_{\text{opt}}}{p} \right] + (1-\eta)\mathbb{E}_{x \sim p_{\text{non}}} \left[\log \frac{(1-\eta)p_{\text{non}}}{p} \right] + \log(2-\eta) - (1-\eta)\log(1-\eta)/(2-\eta) - 2\log 2$, which is a constant for π_θ .

According to Theorem 4.3, since p_1 is a weighted sum of p_θ and p_{non} , subtracting the second and third terms of the Kullback-Leibler (KL) divergence is equivalent to letting p_θ deviate from p_{non} . Thus, PN-GAIL aims to align p_θ with p and ensure that p_θ deviates from p_{non} . This illustrates that our method is able to avoid mimicking non-optimal data within imperfect demonstrations, thereby solely imitating the optimal ones. Additionally, the reward given by the optimal discriminator $D_w^*(x)$ in our method is $-\log D_w^*(x) = \log((p + p_\theta) / ((1-r)p + p_\theta))$. In cases where imperfect demonstrations exhibit a pronounced preference for a certain state-action pair, resulting in a significantly higher probability p compared to other state-action pairs, the presence of the term $(1-r)p$ in the denominator mitigates the impact of an excessively high p . Furthermore, even in extreme scenarios where p is much greater than p_θ , the maximum reward provided by the discriminator in our method is $-\log(1-r)$, rather than approaching positive infinity as in 2IWIL. As a result, in PN-GAIL, the discriminator can offer more accurate rewards, thereby facilitating the subsequent RL process to learn a better policy.

In the following theorem, we demonstrate that the estimation error of Eq. (13) is bounded, indicating that we can obtain a classifier by minimizing $\widehat{R}_{\text{BSC},\ell}$. We provide the estimation error bound with Rademacher complexity (Bartlett & Mendelson, 2002).

Theorem 4.4. *Denote \mathcal{G} as the hypothesis class being utilized and $\mathfrak{R}_n(\mathcal{G})$ as the Rademacher complexity of the function class \mathcal{G} with a sample size of n . Assume that the loss function ℓ is ρ_ℓ -Lipschitz continuous, and there exists a constant $C_\ell > 0$ such that for any $g \in \mathcal{G}$, $\sup_{x \in \mathcal{X}, y \in \{\pm 1\}} |\ell(yg(x))| \leq C_\ell$. Define \hat{g} as the minimizer of $\widehat{R}_{\text{BSC},\ell}(g)$ over $g \in \mathcal{G}$ and g^* as the minimizer of $R_{\text{BSC},\ell}(g)$ over $g \in \mathcal{G}$. For $\delta \in (0, 1)$, with probability at least $1 - \delta$ when repeatedly sampling data to train \hat{g} , we have*

$$R_{\text{BSC},\ell}(\hat{g}) - R_{\text{BSC},\ell}(g^*) \leq 16\rho_L((3 + \alpha - \beta)\mathfrak{R}_{n_c}(\mathcal{G}) + (\alpha + \beta)\mathfrak{R}_{n_u}(\mathcal{G})) + 4C_L \sqrt{\frac{\log(12/\delta)}{2}} \left((3 + \alpha - \beta)n_c^{-\frac{1}{2}} + (\alpha + \beta)n_u^{-\frac{1}{2}} \right). \quad (16)$$

4.4 OVERALL ALGORITHM

Through the aforementioned classifier, we can obtain the confidence scores for all demonstrations in the unlabeled data \mathcal{D}_u . Subsequently, we treat both \mathcal{D}_c and \mathcal{D}_u as imperfect demonstrations and optimize the discriminator D_w . Finally, we utilize Trust Region Policy Optimization (TRPO) (Schulman et al., 2015) to learn a policy π_θ based on the rewards provided by the discriminator. The pseudocode for the overall algorithm can be found in Appendix A.

5 EXPERIMENTS

In this section, we validate our method by conducting experiments on six control tasks, including Pendulum-v1 and five challenging MuJoCo (Todorov et al., 2012) environments. We aim to answer three questions: (1) Is 2IWIL influenced by the preferences inherent in imperfect demonstrations, and can our method alleviate such influence? (2) Does our proposed BSC outperform the SC proposed in 2IWIL? (3) How robust is our method?

Task setup We conduct experiments across six environments (Pendulum-v1, Ant-v2, Walker2d-v2, Hopper-v2, Swimmer-v2, and HalfCheetah-v2). Each experiment is conducted using five different random seeds. Additionally, to better showcase the performance of imitation, we normalize the cumulative rewards of the policies, where 1.0 represents the optimal policy and 0.0 represents the random policy. Due to space constraints, we place the details of the experiments, the performance of the optimal and the random policies and the uncropped figures of Ant-v2 in Appendix C.1, C.4.

Demonstrations For the Pendulum-v1 environment, we train an optimal policy π_{opt} and an intermediate policy π_1 using TRPO. To highlight the preferences inherent in imperfect demonstrations, we aim for a higher proportion of samples to be drawn from π_1 . In that way, we ensure that the number of demonstrations generated by π_1 is four times that of π_{opt} , resulting in a final demonstrations ratio of $\pi_{\text{opt}} : \pi_1 = 1 : 4$, which are then merged together. Afterward, all demonstrations are annotated with confidence scores, utilizing normalized rewards. For the Ant-v2, Walker2d-v2, Hopper-v2, Swimmer-v2, and HalfCheetah-v2 environments, to maintain fairness, we directly utilize the demonstrations and confidence scores provided by the code of 2IWIL. During the practical experiments across all six environments, 20% of the given demonstrations are randomly selected to be assigned confidence scores, which means that the label ratio is 0.2.

Baselines We choose GAIL, 2IWIL, IC-GAIL, and WGAIL as our baseline methods. Among these methods, since GAIL and WGAIL do not require confidence information, we only provide them with demonstrations. Furthermore, we conduct ablation experiments, including **2IWIL**: Original 2IWIL. **PN-GAIL \ BSC**: PN-GAIL with *semi-conf (SC) classification*. **PN-GAIL \ PN**: 2IWIL with *balanced semi-conf (BSC) classification*. **PN-GAIL**: Our final method. All methods are trained jointly using both \mathcal{D}_c and \mathcal{D}_u . Meanwhile, we also test the performance of CAIL, ranking-based methods (T-REX, D-REX) and f-IRL (Ni et al., 2021) by constructing trajectory rankings from confidence scores in the Pendulum-v1 and Ant-v2 environments (due to the demonstrations provided by the 2IWIL’s code is not in trajectory form). The results can be seen in Appendix C.3.

5.1 PERFORMANCE

In our experiments, we use different numbers of $\mathcal{D}_c + \mathcal{D}_u$ for different tasks, and the specific values are shown in Appendix C.1. Fig. 2 and Fig. 3 show the normalized average returns during training. The results in Fig. 2 demonstrate that our method outperforms other baseline methods, achieving the highest returns in all six environments. Of particular note is its performance in Pendulum-v1. Here, imperfect demonstrations exhibit a preference for certain state-action pairs with lower confidence scores, adversely affecting the learning process of 2IWIL and leading to a poor policy. In contrast, our method addresses this issue by incorporating the negative risk of imperfect demonstrations. Experimental results demonstrate that our method is able to learn a near-optimal policy in the Pendulum-v1 environment while other baseline methods fail.

We observe that the performance of GAIL generally falls below that of other methods. This is because GAIL treats all demonstrations as optimal, unable to allocate distinct weights to different

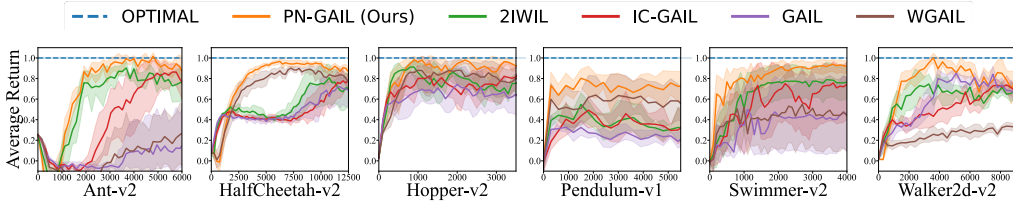


Figure 2: Normalized average returns of PN-GAIL and baseline methods during training. The x-axis is the number of training steps.

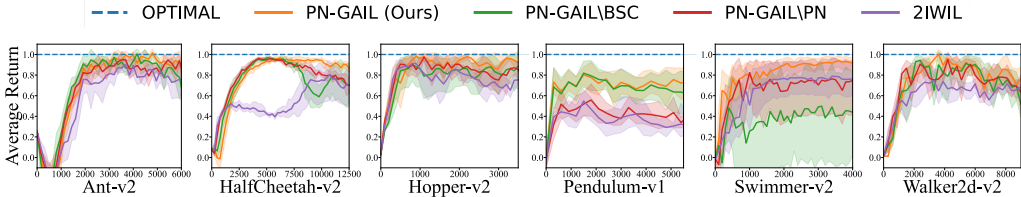


Figure 3: Normalized average returns of ablation experiments during training. The x-axis is the number of training steps.

demonstrations. However, in Walker2d-v2, neither 2IWIL nor IC-GAIL outperforms GAIL. We feel this might be due to the relatively low average confidence of demonstrations in Walker2d-v2. Meanwhile, we notice that WGAIL performs worse than GAIL in Walker2d-v2, which we attribute to its assumption of a higher proportion of optimal demonstrations within the imperfect demonstrations. Since the demonstrations provided in Walker2d-v2 do not align with this assumption, the confidence estimation of WGAIL would no longer be accurate.

Additionally, Fig. 3 shows the normalized average returns of the ablation experiments. In Fig. 3, the large difference between the performance of PN-GAIL and PN-GAIL\PN indicates that there is a preference in the imperfect demonstrations, resulting in the poor performance of the 2IWIL follow-up method. The large performance gap between the performance of PN-GAIL and PN-GAIL\BSC indicates that the prediction confidence of SC classification is not accurate enough, which affects the subsequent training. If the performance gap is not significant, it means that the above problems are not obvious or do not affect the final results. Our method outperforms other methods across all environments, thus confirming the performance enhancement brought by incorporating the negative risk of imperfect demonstrations and employing *balanced semi-conf (BSC) classification*.

5.2 ACCURACY OF CLASSIFIER

By comparing PN-GAIL with PN-GAIL\BSC as depicted in Fig. 3, it is clear that the performance of PN-GAIL can be improved by using the BSC classifier. This observation demonstrates the superior capability of the BSC classifier over the SC classifier in accurately predicting confidence scores. To illustrate the disparity between these two classifiers more clearly, we calculate the Mean Absolute Error (MAE) and Root Mean Square Error (RMSE) of the prediction confidence scores. Here, MAE represents the average of absolute errors, while RMSE denotes the square root of the average of squared differences between predicted and true values. As shown in Table 1, the MAE and RMSE of the BSC classifier are notably lower than those of the SC classifier, indicating that the predictions of the BSC classifier are closer to the ground truth. Consequently, BSC classifier provides more accurate confidence scores for subsequent imitation learning.

5.3 ROBUSTNESS OF PN-GAIL

To test the robustness of our method, We evaluate the performance of PN-GAIL at different label ratios in Ant-v2 and Hopper-v2 environments, the results are shown in Fig. 4 (a) and (b). As the label ratio decreases, PN-GAIL exhibits only a marginal decline in performance. This indicates that

Table 1: Accuracy of classifier measured by MAE and RMSE.

Classifier	Metrics	Ant-v2	HalfCheetah-v2	Hopper-v2	Pendulum-v1	Swimmer-v2	Walker2d-v2
SC	MAE	0.213 \pm 0.023	0.184 \pm 0.011	0.307 \pm 0.025	0.126 \pm 0.014	0.362 \pm 0.049	0.132 \pm 0.015
	RMSE	0.345 \pm 0.033	0.272 \pm 0.009	0.519 \pm 0.022	0.164 \pm 0.013	0.595 \pm 0.040	0.246 \pm 0.032
BSC	MAE	0.056 \pm 0.011	0.057 \pm 0.012	0.169 \pm 0.126	0.097 \pm 0.006	0.286 \pm 0.179	0.014 \pm 0.002
	RMSE	0.212 \pm 0.026	0.175 \pm 0.013	0.371 \pm 0.138	0.138 \pm 0.005	0.472 \pm 0.188	0.101 \pm 0.010

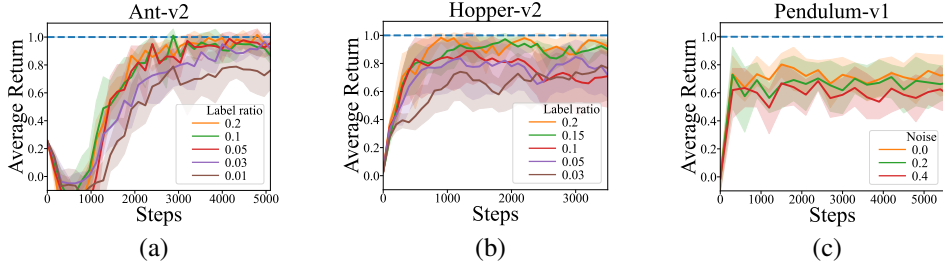


Figure 4: (a) Ant-v2 experiments with different label ratios. (b) Hopper-v2 experiments with different label ratios. (c) Pendulum-v1 experiments with different standard deviations of Gaussian noise.

PN-GAIL is not highly dependent on the label ratio, maintaining excellent performance even as the label ratio decreases.

In practice, considering that confidence scores are typically provided by human annotators, variations in their standards for labeling confidence may arise due to individual differences and factors such as fatigue. To assess the robustness of our method against noise in confidence scores, we conduct additional experiments. In Pendulum-v1, we introduce Gaussian noise to the confidence scores: $\hat{r}(x) = \text{clip}_{[0,1]}(r(x) + \epsilon)$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$, $\text{clip}_{[l,u]}(v) = \min\{\max\{v, l\}, u\}$. As shown in Fig. 4 (c), the numbers indicate the standard deviations of Gaussian noise. Even when confidence scores are subject to noise, our method still demonstrates satisfactory performance, indicating its robustness to noisy confidence scores.

We also test the performance of PN-GAIL in two scenarios: first, by reducing the number of unlabeled demonstrations; and second, by observing how PN-GAIL performs when the average optimality of imperfect demonstrations changes. Due to space constraints, we present the details of these experiments and the corresponding figures in Appendix C.3.

6 CONCLUSION

In this work, we proposed a novel algorithm termed PN-GAIL for imitation learning from imperfect demonstrations. PN-GAIL leverages non-optimal information embedded in these demonstrations, enabling the discriminator to weigh both positive and negative risks in a holistic manner. This approach facilitates the assignment of more refined reward signals. To enhance the precision of confidence estimation, we have integrated an advanced semi-supervised confidence classifier into our framework. Our theoretical investigations demonstrate that PN-GAIL is not merely capable of mimicking imperfect demonstrations but also adept at circumventing the imitation of suboptimal behaviors, thereby ensuring the acquisition of an optimal policy. Comprehensive experimental results indicate that our approach surpasses existing baselines in performance and exhibits remarkable robustness, thereby establishing a robust foundation for the practical deployment of imitation learning in real-world scenarios.

ACKNOWLEDGEMENT

This work was supported in part by the National Natural Science Foundation of China (Nos. 72394363 & 62073160), the Nanjing University Integrated Research Platform of the Ministry of Education-Top Talents Program and the Australian Research Council (FT220100656).