

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

Anonymous authors

Paper under double-blind review

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://anonymous.4open.science/r/PN-GAIL-3828>.

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 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). However, 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:

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

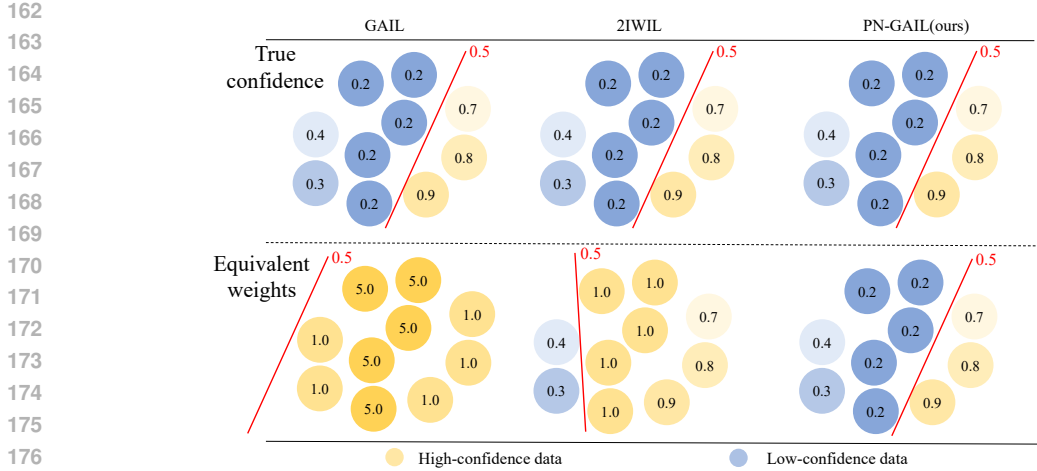


Figure 1: Schematic diagram of the difference between PN-GAIL, 2IWIL and GAIL. The top half of the graph is the actual confidence score, the bottom half is the equivalent weight when training the discriminator, and the red line is distinguished by a threshold of 0.5.

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 A.

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 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 $\mathcal{D} = \mathcal{D}_{opt} + \mathcal{D}_{non}$. 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:

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

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:

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

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

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

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 B.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 $p_1 = (p_\theta + (1-\eta)p_{\text{non}})/(2-\eta)$, $C = \eta\mathbb{E}_{x\sim p_{\text{opt}}}[\log \frac{\eta p_{\text{opt}}}{p}] + (1-\eta)\mathbb{E}_{x\sim p_{\text{non}}}[\log \frac{(1-\eta)p_{\text{non}}}{p}] + \log(2-\eta) - (1-\eta)\log(1-\eta)/(2-\eta) - 2\log 2$, 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 is shown in Algorithm 1.

Algorithm 1 PN-GAIL

```

1: Input: Imperfect demonstrations and confidence  $\mathcal{D}_c = \{(x_{c,i}, r_i)\}_{i=1}^{n_c}$ ,  $\mathcal{D}_u = \{x_{u,i}\}_{i=1}^{n_u}$ 
2: Train a probabilistic classifier by minimizing Eq. (13) with  $\alpha = \frac{n_u}{n_c+n_u}$ ,  $\beta = \frac{n_c}{n_c+n_u}$ 
3: Predict confidence scores  $\{\hat{r}_{u,i}\}_{i=1}^{n_u}$  for  $\{x_{u,i}\}_{i=1}^{n_u}$ 
4: for  $i = 0, 1, 2, \dots$  do
5:   Sample trajectories  $\tau_\theta \sim \pi_\theta$ ,  $\tau_e \sim \{\mathcal{D}_c, \mathcal{D}_u\}$ 
6:   Update  $D_w$  by maximizing Eq. (11)
7:   Update  $\pi_\theta$  by TRPO with reward  $-\log D_w(x)$ 
8: end for

```

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 B.1, B.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 B.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 B.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 demonstra-

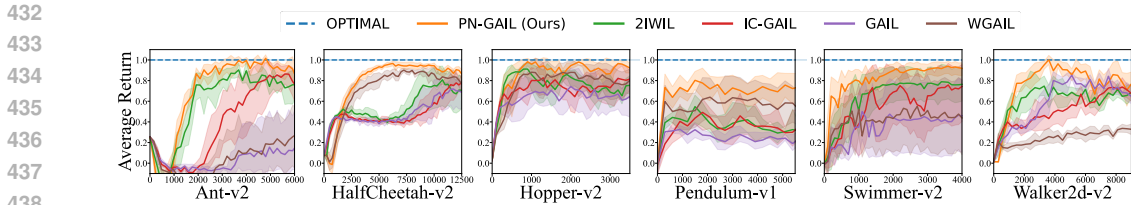


Figure 2: Normalized average returns of PN-GAIL and baseline methods during training. The x-axis is the number of training steps.

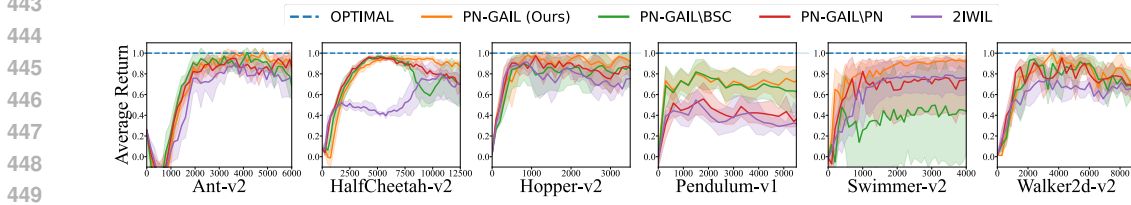


Figure 3: Normalized average returns of ablation experiments during training. The x-axis is the number of training steps.

tions. 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 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.

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 ± 0.023	0.184 ± 0.011	0.307 ± 0.025	0.126 ± 0.014	0.362 ± 0.049	0.132 ± 0.015
	RMSE	0.345 ± 0.033	0.272 ± 0.009	0.519 ± 0.022	0.164 ± 0.013	0.595 ± 0.040	0.246 ± 0.032
BSC	MAE	0.056 ± 0.011	0.057 ± 0.012	0.169 ± 0.126	0.097 ± 0.006	0.286 ± 0.179	0.014 ± 0.002
	RMSE	0.212 ± 0.026	0.175 ± 0.013	0.371 ± 0.138	0.138 ± 0.005	0.472 ± 0.188	0.101 ± 0.010

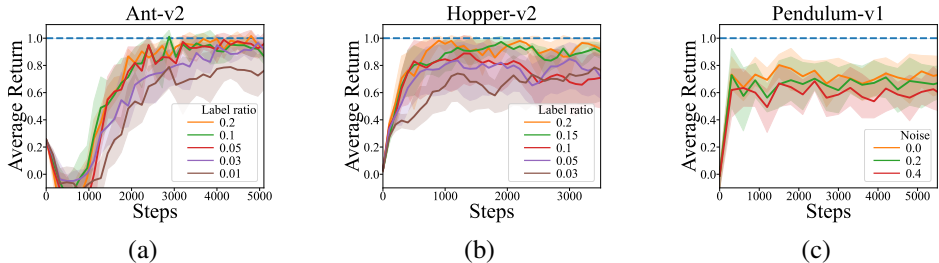


Figure 4: (a) Ant-v2 experiments with different label ratios. (b) Hopper-v2 experiments with different label ratios. (c) Pendulum-v1 experiments with different standard deviations of Gaussian noise.

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 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 B.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.

REFERENCES

- 540
541
542 Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. In
543 *International Conference on Machine Learning*, pp. 1, 2004.
- 544
545 Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and
546 structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- 547
548 Raunak Bhattacharyya, Blake Wulfe, Derek J Phillips, Alex Kuefler, Jeremy Morton, Ransalu
549 Senanayake, and Mykel J Kochenderfer. Modeling human driving behavior through generative
550 adversarial imitation learning. *IEEE Transactions on Intelligent Transportation Systems*, 24(3):
2874–2887, 2022.
- 551
552 Daniel Brown, Wonjoon Goo, Prabhat Nagarajan, and Scott Niekum. Extrapolating beyond sub-
553 optimal demonstrations via inverse reinforcement learning from observations. In *International
554 Conference on Machine Learning*, pp. 783–792. PMLR, 2019.
- 555
556 Daniel S Brown, Wonjoon Goo, and Scott Niekum. Better-than-demonstrator imitation learning via
557 automatically-ranked demonstrations. In *Conference on Robot Learning*, pp. 330–359. PMLR,
2020.
- 558
559 Letian Chen, Rohan Paleja, and Matthew Gombolay. Learning from suboptimal demonstration via
560 self-supervised reward regression. In *Conference on Robot Learning*, pp. 1262–1277. PMLR,
561 2021.
- 562
563 Robert Dadashi, Leonard Hussenot, Matthieu Geist, and Olivier Pietquin. Primal wasserstein imita-
564 tion learning. In *International Conference on Learning Representations*, 2020.
- 565
566 Marthinus C Du Plessis, Gang Niu, and Masashi Sugiyama. Analysis of learning from positive and
unlabeled data. *Advances in Neural Information Processing Systems*, 27, 2014.
- 567
568 Huiqiao Fu, Kaiqiang Tang, Yuanyang Lu, Yiming Qi, Guizhou Deng, Flood Sung, and Chunlin
569 Chen. Ess-infogail: Semi-supervised imitation learning from imbalanced demonstrations. *Ad-
570 vances in Neural Information Processing Systems*, 36, 2024.
- 571
572 Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse rein-
forcement learning. In *International Conference on Learning Representations*, 2018.
- 573
574 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
575 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in Neural Informa-
576 tion Processing Systems*, 27, 2014.
- 577
578 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
579 maximum entropy deep reinforcement learning with a stochastic actor. In *International Confer-
ence on Machine Learning*, pp. 1861–1870. PMLR, 2018.
- 580
581 Jeffrey Hawke, Richard Shen, Corina Gurau, Siddharth Sharma, Daniele Reda, Nikolay Nikolov,
582 Przemysław Mazur, Sean Micklethwaite, Nicolas Griffiths, Amar Shah, et al. Urban driving
583 with conditional imitation learning. In *2020 IEEE International Conference on Robotics and
584 Automation (ICRA)*, pp. 251–257. IEEE, 2020.
- 585
586 Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in Neural
Information Processing Systems*, 29, 2016.
- 587
588 Liangyu Huo, Zulin Wang, and Mai Xu. Learning noise-induced reward functions for surpassing
589 demonstrations in imitation learning. In *Proceedings of the AAAI Conference on Artificial Intel-
590 ligence*, volume 37, pp. 7953–7961, 2023.
- 591
592 Liyiming Ke, Sanjiban Choudhury, Matt Barnes, Wen Sun, Gilwoo Lee, and Siddhartha Srinivasa.
593 Imitation learning as f-divergence minimization. In *Algorithmic Foundations of Robotics XIV:
Proceedings of the Fourteenth Workshop on the Algorithmic Foundations of Robotics 14*, pp.
313–329. Springer, 2021.

- 594 Geon-Hyeong Kim, Seokin Seo, Jongmin Lee, Wonseok Jeon, HyeongJoo Hwang, Hongseok Yang,
595 and Kee-Eung Kim. Demodice: Offline imitation learning with supplementary imperfect demon-
596 strations. In *International Conference on Learning Representations*, 2021.
- 597 Yunzhu Li, Jiaming Song, and Stefano Ermon. Infogail: Interpretable imitation learning from visual
598 demonstrations. *Advances in Neural Information Processing Systems*, 30, 2017.
- 600 Ziniu Li, Tian Xu, Zeyu Qin, Yang Yu, and Zhi-Quan Luo. Imitation learning from imperfection:
601 Theoretical justifications and algorithms. *Advances in Neural Information Processing Systems*,
602 36, 2024.
- 603 Tianwei Ni, Harshit Sikchi, Yufei Wang, Tejus Gupta, Lisa Lee, and Ben Eysenbach. f-irl: Inverse
604 reinforcement learning via state marginal matching. In *Conference on Robot Learning*, pp. 529–
605 551. PMLR, 2021.
- 606 Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. Amp: Adversarial
607 motion priors for stylized physics-based character control. *ACM Transactions on Graphics (ToG)*,
608 40(4):1–20, 2021.
- 610 Dean A Pomerleau. Alvin: An autonomous land vehicle in a neural network. *Advances in Neural
611 Information Processing Systems*, 1, 1988.
- 612 Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John
613 Wiley & Sons, 2014.
- 614 Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. In *Proceedings of the
615 thirteenth international conference on artificial intelligence and statistics*, pp. 661–668. JMLR
616 Workshop and Conference Proceedings, 2010.
- 618 Kangrui Ruan and Xuan Di. Learning human driving behaviors with sequential causal imitation
619 learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 4583–
620 4592, 2022.
- 621 Fumihiko Sasaki and Ryota Yamashina. Behavioral cloning from noisy demonstrations. In *Internat-
622 ional Conference on Learning Representations*, 2020.
- 624 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region
625 policy optimization. In *International Conference on Machine Learning*, pp. 1889–1897. PMLR,
626 2015.
- 627 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
628 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 630 Jiaming Song, Hongyu Ren, Dorsa Sadigh, and Stefano Ermon. Multi-agent generative adversarial
631 imitation learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- 632 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 633 Voot Tangkaratt, Bo Han, Mohammad Emtiyaz Khan, and Masashi Sugiyama. Variational imitation
634 learning with diverse-quality demonstrations. In *International Conference on Machine Learning*,
635 pp. 9407–9417. PMLR, 2020.
- 637 Aleksandar Taranovic, Andras Gabor Kupcsik, Niklas Freymuth, and Gerhard Neumann. Adversar-
638 ial imitation learning with preferences. In *International Conference on Learning Representations*,
639 2022.
- 640 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control.
641 In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pp. 5026–5033.
642 IEEE, 2012.
- 643 Yunke Wang, Chang Xu, and Bo Du. Robust adversarial imitation learning via adaptively-selected
644 demonstrations. In *IJCAI*, pp. 3155–3161, 2021a.
- 645 Yunke Wang, Chang Xu, Bo Du, and Honglak Lee. Learning to weight imperfect demonstrations.
646 In *International Conference on Machine Learning*, pp. 10961–10970. PMLR, 2021b.

- 648 Yunke Wang, Bo Du, and Chang Xu. Unlabeled imperfect demonstrations in adversarial imitation
649 learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 10262–
650 10270, 2023.
- 651 Hua Wei, Dongkuan Xu, Junjie Liang, and Zhenhui Jessie Li. How do we move: Modeling human
652 movement with system dynamics. In *Proceedings of the AAAI Conference on Artificial Intelli-*
653 *gence*, volume 35, pp. 4445–4452, 2021.
- 654 Yueh-Hua Wu, Nontawat Charoenphakdee, Han Bao, Voot Tangkaratt, and Masashi Sugiyama. Imitation
655 learning from imperfect demonstration. In *International Conference on Machine Learning*,
656 pp. 6818–6827. PMLR, 2019.
- 657 Fei Xia, Chengshu Li, Roberto Martín-Martín, Or Litany, Alexander Toshev, and Silvio Savarese.
658 Relmogen: Leveraging motion generation in reinforcement learning for mobile manipulation.
659 *arXiv preprint arXiv:2008.07792*, 2020.
- 660 Danfei Xu and Misha Denil. Positive-unlabeled reward learning. In *Conference on Robot Learning*,
661 pp. 205–219. PMLR, 2021.
- 662 Haoran Xu, Xianyuan Zhan, Honglei Yin, and Huiling Qin. Discriminator-weighted offline imitation
663 learning from suboptimal demonstrations. In *International Conference on Machine Learning*, pp.
664 24725–24742. PMLR, 2022.
- 665 Hanlin Yang, Chao Yu, Siji Chen, et al. Hybrid policy optimization from imperfect demonstrations.
666 *Advances in Neural Information Processing Systems*, 36, 2024.
- 667 Lantao Yu, Tianhe Yu, Jiaming Song, Willie Neiswanger, and Stefano Ermon. Offline imitation
668 learning with suboptimal demonstrations via relaxed distribution matching. In *Proceedings of the*
669 *AAAI conference on artificial intelligence*, volume 37, pp. 11016–11024, 2023.
- 670 Daochen Zha, Jingru Xie, Wenye Ma, Sheng Zhang, Xiangru Lian, Xia Hu, and Ji Liu. Douzero:
671 Mastering doudizhu with self-play deep reinforcement learning. In *International Conference on*
672 *Machine Learning*, pp. 12333–12344. PMLR, 2021.
- 673 Songyuan Zhang, Zhangjie Cao, Dorsa Sadigh, and Yanan Sui. Confidence-aware imitation learn-
674 ing from demonstrations with varying optimality. *Advances in Neural Information Processing*
675 *Systems*, 34:12340–12350, 2021.
- 676 Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy
677 inverse reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
678 volume 8, pp. 1433–1438. Chicago, IL, USA, 2008.
- 679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

Appendices

A DERIVATIONS AND PROOFS

A.1 PROOF OF THEOREM 4.1

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

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

Proof. Since $p_{\text{opt}}(x) = p(x|y=0) = \frac{r(x)p(x)}{\eta}$, $p_{\text{non}}(x) = p(x|y=1) = \frac{(1-r(x))p(x)}{1-\eta}$, we have

$$\begin{aligned} \eta \mathbb{E}_{x \sim p_{\text{opt}}} [\log(1 - D_w(x))] &= \eta \int p_{\text{opt}}(x) \log(1 - D_w(x)) dx \\ &= \eta \int \frac{r(x)}{\eta} p(x) \log(1 - D_w(x)) dx \\ &= \int p(x) r(x) \log(1 - D_w(x)) dx \\ &= \mathbb{E}_{x \sim p} [r(x) \log(1 - D_w(x))], \end{aligned}$$

$$\begin{aligned} (1 - \eta) \mathbb{E}_{x \sim p_{\text{non}}} [\log D_w(x)] &= (1 - \eta) \int p_{\text{non}}(x) \log D_w(x) dx \\ &= (1 - \eta) \int \frac{1 - r(x)}{1 - \eta} p(x) \log D_w(x) dx \\ &= \int p(x) (1 - r(x)) \log D_w(x) dx \\ &= \mathbb{E}_{x \sim p} [(1 - r(x)) \log D_w(x)], \end{aligned}$$

According to Eq. (8), the objective of the discriminator can be rewritten as

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

□

A.2 DERIVATION OF BALANCED SEMI-CONF (BSC) CLASSIFICATION

Recall that in 2IWIL:

$$\begin{aligned} R_{\text{SC},\ell}(g) &= \mathbb{E}_{x,r \sim q} [r\ell(g(x)) + (1-r)\ell(-g(x))] \\ &= \mathbb{E}_{x,r \sim q} \left[r\ell(g(x)) + (1-r)\ell(-g(x)) + \underbrace{\beta\ell(-g(x)) - \beta\ell(-g(x))}_{=0} \right] \end{aligned}$$

By introducing $\mathbb{E}_{x \sim p} [\alpha\ell(g(x))] - \mathbb{E}_{x \sim q} [\alpha\ell(g(x))]$ with theoretical values of 0, we have

$$\begin{aligned} R_{\text{BSC},\ell}(g) &= \mathbb{E}_{x,r \sim q} \left[r\ell(g(x)) + (1-r)\ell(-g(x)) + \underbrace{\alpha\ell(g(x)) - \alpha\ell(g(x))}_{=0} + \underbrace{\beta\ell(-g(x)) - \beta\ell(-g(x))}_{=0} \right] \\ &= \mathbb{E}_{x,r \sim q} [r\ell(g(x)) + (1-r)\ell(-g(x)) - \alpha\ell(g(x)) - \beta\ell(-g(x))] + \mathbb{E}_{x \sim p} [\alpha\ell(g(x)) + \beta\ell(-g(x))]. \end{aligned}$$

A.3 PROOF OF THEOREM 4.2

Theorem. Let d_1 denote $\text{Var}(\ell(-g(x)))$, d_2 denote $\text{Var}(\ell(g(x)))$, σ_{cov1} 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 $\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i}))$, σ_{cov2} 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

$$\alpha = \frac{n_u}{n_c + n_u} - \frac{d_1 \text{cov} - \text{cov}^2}{d_1 d_2 - \text{cov}^2} \frac{n_u}{n_c + n_u} + \frac{d_1 \sigma_{\text{cov2}} - \text{cov} \sigma_{\text{cov1}}}{d_1 d_2 - \text{cov}^2} \frac{n_c n_u}{n_c + n_u},$$

$$\beta = \frac{n_u}{n_c + n_u} - \frac{d_2 \text{cov} - \text{cov}^2}{d_1 d_2 - \text{cov}^2} \frac{n_u}{n_c + n_u} + \frac{d_2 \sigma_{\text{cov1}} - \text{cov} \sigma_{\text{cov2}}}{d_1 d_2 - \text{cov}^2} \frac{n_c n_u}{n_c + n_u}.$$

Proof. Denote

$$\begin{aligned} \mu &\triangleq \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [\widehat{R}_{\text{BSC},\ell}(g)], \\ \mu_1 &\triangleq \mathbb{E}_{\mathcal{D}_c} \left[\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right] = \mathbb{E}_{\mathcal{D}_u} \left[\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) \right] = \mathbb{E}_{x \sim p} [\ell(-g(x))], \\ \mu_2 &\triangleq \mathbb{E}_{\mathcal{D}_c} \left[\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(g(x_{c,i})) \right] = \mathbb{E}_{\mathcal{D}_u} \left[\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(g(x_{u,i})) \right] = \mathbb{E}_{x \sim p} [\ell(g(x))], \\ d_1 &\triangleq \text{Var}_{\mathcal{D}_c} [\ell(-g(x_c))] = \text{Var}_{\mathcal{D}_u} [\ell(-g(x_u))] = \text{Var} [\ell(-g(x))], \\ d_2 &\triangleq \text{Var}_{\mathcal{D}_c} [\ell(g(x_c))] = \text{Var}_{\mathcal{D}_u} [\ell(g(x_u))] = \text{Var} [\ell(g(x))], \\ \omega &\triangleq \mathbb{E}_{\mathcal{D}_c} \left[\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \ell(g(x_{c,i})) \right] = \mathbb{E}_{\mathcal{D}_u} \left[\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) \ell(g(x_{u,i})) \right] = \mathbb{E}_{x \sim p} [\ell(-g(x)) \ell(g(x))], \\ \text{cov} &\triangleq \text{Cov}_{\mathcal{D}_c} (\ell(-g(x_c)), \ell(g(x_c))) = \text{Cov}_{\mathcal{D}_u} (\ell(-g(x_u)), \ell(g(x_u))) = \text{Cov} (\ell(-g(x)), \ell(g(x))) = \omega - \mu_1 \mu_2, \\ \sigma_{\text{cov1}} &\triangleq \text{Cov} \left(\frac{1}{n_c} \sum_{i=1}^{n_c} r_i (\ell(g(x_{c,i})) - \ell(-g(x_{c,i}))), \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right), \\ \sigma_{\text{cov2}} &\triangleq \text{Cov} \left(\frac{1}{n_c} \sum_{i=1}^{n_c} (1 - r_i) (\ell(-g(x_{c,i})) - \ell(g(x_{c,i}))), \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(g(x_{c,i})) \right). \end{aligned}$$

Next, we adopt the symbols defined above to express several formulas that will be used:

$$\begin{aligned} \mathbb{E}_{\mathcal{D}_c} \left[\left(\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right)^2 \right] &= \frac{1}{n_c^2} \mathbb{E}_{\mathcal{D}_c} \left[\sum_{i=1}^{n_c} \ell(-g(x_{c,i}))^2 + 2 \sum_{i=1}^{n_c} \sum_{j=1}^{i-1} \ell(-g(x_{c,i})) \ell(-g(x_{c,j})) \right] \\ &= \frac{1}{n_c^2} \left(n_c \mathbb{E}_{x \sim p} [\ell(-g(x))^2] + n_c(n_c - 1) \mathbb{E}_{x \sim p} [\ell(-g(x))]^2 \right) \\ &= \frac{1}{n_c} \text{Var}(\ell(-g(x))) + \mu_1^2 \\ &= \frac{1}{n_c} d_1 + \mu_1^2. \end{aligned} \tag{17}$$

810 Similarly, we obtain

$$813 \mathbb{E}_{\mathcal{D}_u} \left[\left(\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) \right)^2 \right] = \frac{1}{n_u} d_1 + \mu_1^2, \quad (18)$$

$$817 \mathbb{E}_{\mathcal{D}_c} \left[\left(\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(g(x_{c,i})) \right)^2 \right] = \frac{1}{n_c} d_2 + \mu_2^2, \quad (19)$$

$$821 \mathbb{E}_{\mathcal{D}_u} \left[\left(\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(g(x_{u,i})) \right)^2 \right] = \frac{1}{n_u} d_2 + \mu_2^2, \quad (20)$$

$$824 \mathbb{E}_{\mathcal{D}_u} \left[\left(\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) \right) \left(\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(g(x_{u,i})) \right) \right] = \frac{1}{n_u} cov + \mu_1 \mu_2, \quad (21)$$

$$827 \mathbb{E}_{\mathcal{D}_c} \left[\left(\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right) \left(\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(g(x_{c,i})) \right) \right] = \frac{1}{n_c} cov + \mu_1 \mu_2. \quad (22)$$

831 We also have

$$834 \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} \left[\left(\frac{1}{n_c} \sum_{i=1}^{n_c} r(x_i) (\ell(g(x_{c,i})) - \ell(-g(x_{c,i}))) \right) \left(\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) - \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right) \right]$$

$$837 = \mathbb{E}_{\mathcal{D}_c} \left[\left(\frac{1}{n_c} \sum_{i=1}^{n_c} r(x_i) (\ell(g(x_{c,i})) - \ell(-g(x_{c,i}))) \right) \right] \mathbb{E}_{\mathcal{D}_u} \left[\left(\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) \right) \right]$$

$$841 - \mathbb{E}_{\mathcal{D}_c} \left[\left(\frac{1}{n_c} \sum_{i=1}^{n_c} r(x_i) (\ell(g(x_{c,i})) - \ell(-g(x_{c,i}))) \right) \left(\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right) \right]$$

$$843 = \mathbb{E}_{\mathcal{D}_c} \left[\left(\frac{1}{n_c} \sum_{i=1}^{n_c} r(x_i) (\ell(g(x_{c,i})) - \ell(-g(x_{c,i}))) \right) \right] \mathbb{E}_{\mathcal{D}_c} \left[\left(\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right) \right]$$

$$846 - \mathbb{E}_{\mathcal{D}_c} \left[\left(\frac{1}{n_c} \sum_{i=1}^{n_c} r(x_i) (\ell(g(x_{c,i})) - \ell(-g(x_{c,i}))) \right) \left(\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right) \right]$$

$$849 = \text{Cov} \left(\frac{1}{n_c} \sum_{i=1}^{n_c} r_i (\ell(g(x_{c,i})) - \ell(-g(x_{c,i}))), \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right)$$

$$851 = -\sigma_{\text{cov}1}. \quad (23)$$

854 Similarly, we obtain

$$857 \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} \left[\left(\frac{1}{n_c} \sum_{i=1}^{n_c} (1 - r(x_i)) (\ell(-g(x_{c,i})) - \ell(g(x_{c,i}))) \right) \left(\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(g(x_{u,i})) - \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(g(x_{c,i})) \right) \right]$$

$$860 = -\sigma_{\text{cov}2}. \quad (24)$$

863 Hence, we can express the estimator variance $\text{Var}(\widehat{R}_{\text{BSC}, \ell}(g))$ as follows:

$$\begin{aligned}
& \text{Var}(\widehat{R}_{\text{BSC},\ell}(g)) \\
&= \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} \left[\left(\widehat{R}_{\text{BSC},\ell}(g) \right)^2 \right] - \mu^2 \\
&= \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} \left[\left(\underbrace{\left(\frac{1}{n_c} \sum_{i=1}^{n_c} r_i \ell(g(x_{c,i})) + (1-r_i) \ell(-g(x_{c,i})) \right)}_A + \beta \underbrace{\left(\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) - \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right)}_B \right)}_C \right. \\
&\quad \left. + \alpha \left(\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(g(x_{u,i})) - \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(g(x_{c,i})) \right) \right)^2 \right] - \mu^2 \\
&= \mathbb{E}_{\mathcal{D}_c} [A^2] + \beta^2 \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2] + \alpha^2 \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [C^2] + 2\beta \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [AB] + 2\alpha \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [AC] + 2\beta\alpha \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [BC] - \mu^2 \\
&= \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2] \left(\beta + \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [AB]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2]} + \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [BC]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2]} \alpha \right)^2 + \left(\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [C^2] - \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}^2 [BC]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2]} \right) \alpha^2 \\
&\quad + 2 \left(\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [AC] - \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [AB] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [BC]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2]} \right) \alpha + \underbrace{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [A^2] - \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}^2 [AB]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2]} - \mu^2}_{\text{const.w.r.t.}\alpha, \beta} \\
&= \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2] \left(\beta + \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [AB]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2]} + \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [BC]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2]} \alpha \right)^2 \\
&\quad + \left(\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [C^2] - \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}^2 [BC]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2]} \right) \left(\alpha + \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [AC] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2] - \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [AB] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [BC]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [C^2] - \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}^2 [BC]} \right)^2 \\
&\quad - \underbrace{\left(\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [C^2] - \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}^2 [BC]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2]} \right) \left(\frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [AC] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2] - \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [AB] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [BC]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [C^2] - \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}^2 [BC]} \right)^2}_{\text{const.w.r.t.}\alpha, \beta} + \text{const} \\
&= \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2] \left(\beta + \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [AB]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2]} + \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [BC]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2]} \alpha \right)^2 \\
&\quad + \left(\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [C^2] - \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}^2 [BC]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2]} \right) \left(\alpha - \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [AB] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [BC] - \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [AC] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [B^2] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [C^2] - \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}^2 [BC]} \right)^2 + \text{const.} \\
&\hspace{15em} (25)
\end{aligned}$$

By Eq. (21) and Eq. (22), we can obtain

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} [BC] &= \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} \left[\left(\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) - \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right) \left(\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(g(x_{u,i})) - \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(g(x_{c,i})) \right) \right] \\
&= \frac{1}{n_u} \text{cov} + \mu_1 \mu_2 - \mu_1 \mu_2 - \mu_1 \mu_2 + \frac{1}{n_c} \text{cov} + \mu_1 \mu_2 \\
&= \left(\frac{1}{n_u} + \frac{1}{n_c} \right) \text{cov}.
\end{aligned}$$

By Eq. (17) and Eq. (18), we can obtain

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[B^2] &= \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} \left[\left(\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) - \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right)^2 \right] \\
&= \mathbb{E}_{\mathcal{D}_u} \left[\left(\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) \right)^2 \right] + \mathbb{E}_{\mathcal{D}_c} \left[\left(\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right)^2 \right] \\
&\quad - 2\mathbb{E}_{\mathcal{D}_u} \left[\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) \right] \mathbb{E}_{\mathcal{D}_c} \left[\frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right] \\
&= \frac{1}{n_u} d_1 + \mu_1^2 + \frac{1}{n_c} d_1 + \mu_1^2 - 2\mu^2 \\
&= \left(\frac{1}{n_u} + \frac{1}{n_c} \right) d_1.
\end{aligned}$$

Similarly, we can obtain $\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[C^2] = \left(\frac{1}{n_u} + \frac{1}{n_c} \right) d_2$. Hence, we have $\left(\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[C^2] - \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[BC]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[B^2]} \right) \geq 0$ since $d_1 d_2 \geq cov^2$. By Eq. (17) and Eq. (23), we can obtain:

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[AB] &= \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u} \left[\left(\frac{1}{n_c} \sum_{i=1}^{n_c} r(x_i) (\ell(g(x_{c,i})) - \ell(-g(x_{c,i}))) + \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right) \right. \\
&\quad \left. \left(\frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) - \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) \right) \right] \\
&= -\sigma_{cov1} + \mu_1^2 - \left(\frac{1}{n_c} d_1 + \mu_1^2 \right) \\
&= -\frac{1}{n_c} d_1 - \sigma_{cov1}.
\end{aligned}$$

Similarly, $\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[AC] = -\frac{1}{n_c} d_2 - \sigma_{cov2}$. Since $\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[B^2] \geq 0$, $\left(\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[C^2] - \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[BC]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[B^2]} \right) \geq 0$, and $\alpha, \beta \in \mathcal{R}$, according to Eq. (25), $\text{Var}(\widehat{R}_{\text{BSC}, \ell}(g))$ is minimized when

$$\begin{aligned}
\beta &= -\frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[AB]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[B^2]} - \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[BC]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[B^2]} \alpha, \\
\alpha &= \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[AB] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[BC] - \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[AC] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[B^2]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[B^2] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[C^2] - \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}^2[BC]}.
\end{aligned}$$

Substitute α , we obtain

$$\beta = \frac{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[AC] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[BC] - \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[AB] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[C^2]}{\mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[B^2] \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}[C^2] - \mathbb{E}_{\mathcal{D}_c, \mathcal{D}_u}^2[BC]}.$$

Through plugging in the above formula, we have

$$\begin{aligned}
\alpha &= \frac{n_u}{n_c + n_u} - \frac{d_1 cov - cov^2}{d_1 d_2 - cov^2} \frac{n_u}{n_c + n_u} + \frac{d_1 \sigma_{cov2} - cov \sigma_{cov1}}{d_1 d_2 - cov^2} \frac{n_c n_u}{n_c + n_u}, \\
\beta &= \frac{n_u}{n_c + n_u} - \frac{d_2 cov - cov^2}{d_1 d_2 - cov^2} \frac{n_u}{n_c + n_u} + \frac{d_2 \sigma_{cov1} - cov \sigma_{cov2}}{d_1 d_2 - cov^2} \frac{n_c n_u}{n_c + n_u}.
\end{aligned}$$

□

A.4 PROOF OF THEOREM 4.3

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

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

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

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

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 .

Proof. Denote

$$V(\pi_\theta, D_w) = \mathbb{E}_{x\sim p_\theta}[\log D_w(x)] + \mathbb{E}_{x\sim p}[r(x)\log(1-D_w(x))] + \mathbb{E}_{x\sim p}[(1-r(x))\log D_w(x)].$$

We denote $D_w^*(x) = \arg\min_D V(\pi_\theta, D_w)$ and have

$$\frac{\partial V(\pi_\theta, D_w)}{\partial D} = \frac{p_\theta}{D} - \frac{rp}{1-D} + \frac{(1-r)p}{D}.$$

The maximum value of $V(\pi_\theta, D_w)$ occurs when its partial derivative with respect to D is zero, given by $D_w = \frac{(1-r)p + p_\theta}{p + p_\theta}$. Consequently, we obtain $D_w^*(x) = \frac{(1-r)p + p_\theta}{p + p_\theta}$.

When the optimal discriminator $D_w^*(x)$ is given, $V(\pi_\theta, D_w)$ can be rewritten as

$$V(\pi_\theta, D_w) = \mathbb{E}_{x\sim p_\theta}\left[\log\frac{(1-r)p + p_\theta}{p + p_\theta}\right] + \mathbb{E}_{x\sim p}\left[r\log\frac{rp}{p + p_\theta}\right] + \mathbb{E}_{x\sim p}\left[(1-r)\log\frac{(1-r)p + p_\theta}{p + p_\theta}\right].$$

According to Eq. (9) and denoting $p_1 = (p_\theta + (1-\eta)p_{\text{non}})/(2-\eta)$, we obtain

$$\begin{aligned} V(\pi_\theta, D_w) &= \mathbb{E}_{x\sim p_\theta}\left[\log\frac{(1-\eta)p_{\text{non}} + p_\theta}{p + p_\theta}\right] + \mathbb{E}_{x\sim p_{\text{opt}}}\left[\eta\log\frac{\eta p_{\text{opt}}}{p + p_\theta}\right] + \mathbb{E}_{x\sim p_{\text{non}}}\left[(1-\eta)\log\frac{(1-\eta)p_{\text{non}} + p_\theta}{p + p_\theta}\right] \\ &= \mathbb{E}_{x\sim p_\theta}\left[-\log\frac{p_\theta}{(1-\eta)p_{\text{non}} + p_\theta} + \log\frac{p_\theta}{p + p_\theta}\right] + \eta\mathbb{E}_{x\sim p_{\text{opt}}}\left[\log\frac{\eta p_{\text{opt}}}{p + p_\theta}\right] \\ &\quad + (1-\eta)\mathbb{E}_{x\sim p_{\text{non}}}\left[-\log\frac{(1-\eta)p_{\text{non}}}{(1-\eta)p_{\text{non}} + p_\theta} + \log\frac{(1-\eta)p_{\text{non}}}{p + p_\theta}\right] \\ &= \mathbb{E}_{x\sim p_\theta}\left[\log\frac{p_\theta}{p + p_\theta}\right] + \eta\mathbb{E}_{x\sim p_{\text{opt}}}\left[\log\frac{\eta p_{\text{opt}}}{p + p_\theta}\right] + (1-\eta)\mathbb{E}_{x\sim p_{\text{non}}}\left[\log\frac{(1-\eta)p_{\text{non}}}{p + p_\theta}\right] \\ &\quad - \left(\mathbb{E}_{x\sim p_\theta}\left[\log\frac{p_\theta}{p_1}\right] - \underbrace{\log(2-\eta)}_{\text{const.w.r.t.}\pi_\theta} + (1-\eta)\mathbb{E}_{x\sim p_{\text{non}}}\left[\log\frac{p_{\text{non}}}{p_1}\right] + (1-\eta)\log\frac{1-\eta}{2-\eta}\right) \\ &\quad \underbrace{\hspace{10em}}_{\text{const.w.r.t.}\pi_\theta} \\ &= \mathbb{E}_{x\sim p_\theta}\left[\log\frac{p_\theta}{p + p_\theta}\right] + \eta\mathbb{E}_{x\sim p_{\text{opt}}}\left[\log\frac{\eta p_{\text{opt}}}{p} + \log\frac{p}{p + p_\theta}\right] + (1-\eta)\mathbb{E}_{x\sim p_{\text{non}}}\left[\log\frac{(1-\eta)p_{\text{non}}}{p}\right] \\ &\quad + \log\frac{p}{p + p_\theta} - \text{KL}(p_\theta||p_1) - (1-\eta)\text{KL}(p_{\text{non}}||p_1) + C_1 \\ &= \mathbb{E}_{x\sim p_\theta}\left[\log\frac{p_\theta}{p + p_\theta}\right] + \mathbb{E}_{x\sim p}\left[\log\frac{p}{p + p_\theta}\right] + \underbrace{\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]}_{\text{const.w.r.t.}\pi_\theta} \\ &\quad - \text{KL}(p_\theta||p_1) - (1-\eta)\text{KL}(p_{\text{non}}||p_1) + C_1 \\ &= \mathbb{E}_{x\sim p_\theta}\left[\log\frac{p_\theta}{(p + p_\theta)/2}\right] + \mathbb{E}_{x\sim p}\left[\log\frac{p}{(p + p_\theta)/2}\right] - 2\log 2 - \text{KL}(p_\theta||p_1) - (1-\eta)\text{KL}(p_{\text{non}}||p_1) + C_2 \\ &= 2\text{JSD}(p_\theta||p) - \text{KL}(p_\theta||p_1) - (1-\eta)\text{KL}(p_{\text{non}}||p_1) + C, \end{aligned}$$

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

A.5 PROOF OF THEOREM 4.4

Theorem. 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

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

Proof. Like 2IWIL, since \hat{g} and g^* are the minimizers of $\widehat{R}_{\text{BSC}, \ell}(g)$ and $R_{\text{BSC}, \ell}(g)$, respectively, we have

$$\begin{aligned} R_{\text{BSC}, \ell}(\hat{g}) - R_{\text{BSC}, \ell}(g^*) &= R_{\text{BSC}, \ell}(\hat{g}) - \widehat{R}_{\text{BSC}, \ell}(\hat{g}) + \widehat{R}_{\text{BSC}, \ell}(\hat{g}) - \widehat{R}_{\text{BSC}, \ell}(g^*) + \widehat{R}_{\text{BSC}, \ell}(g^*) - R_{\text{BSC}, \ell}(g^*) \\ &\leq \sup_{g \in \mathcal{G}} \left(R_{\text{BSC}, \ell}(g) - \widehat{R}_{\text{BSC}, \ell}(g) \right) + 0 + \sup_{g \in \mathcal{G}} \left(\widehat{R}_{\text{BSC}, \ell}(g) - R_{\text{BSC}, \ell}(g) \right) \\ &\leq 2 \sup_{g \in \mathcal{G}} \left| \widehat{R}_{\text{BSC}, \ell}(g) - R_{\text{BSC}, \ell}(g) \right|. \end{aligned}$$

Hence, we just need to bound the uniform deviation $\sup_{g \in \mathcal{G}} \left| \widehat{R}_{\text{BSC}, \ell}(g) - R_{\text{BSC}, \ell}(g) \right|$. We have

$$\begin{aligned} &\sup_{g \in \mathcal{G}} \left| \widehat{R}_{\text{BSC}, \ell}(g) - R_{\text{BSC}, \ell}(g) \right| \\ &\leq \sup_{g \in \mathcal{G}} \left| \frac{1}{n_c} \sum_{i=1}^{n_c} (r_i(\ell(g(x_{c,i})) - \ell(-g(x_{c,i}))) + (1 - \beta)\ell(-g(x_{c,i})) - \alpha\ell(g(x_{c,i}))) \right. \\ &\quad \left. - \mathbb{E}_{x, r \sim q} [r(\ell(g(x)) - \ell(-g(x))) + (1 - \beta)\ell(-g(x)) - \alpha\ell(g(x))] \right| \\ &\quad + \beta \sup_{g \in \mathcal{G}} \left| \frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) - \mathbb{E}_{x \sim p} [\ell(-g(x))] \right| + \alpha \sup_{g \in \mathcal{G}} \left| \frac{1}{n_u} \sum_{i=1}^{n_u} \ell(g(x_{u,i})) - \mathbb{E}_{x \sim p} [\ell(g(x))] \right| \\ &\leq \sup_{g \in \mathcal{G}} \left| \frac{1}{n_c} \sum_{i=1}^{n_c} r_i \ell(g(x_{c,i})) - \mathbb{E}_{x, r \sim q} [r \ell(g(x))] \right| + \sup_{g \in \mathcal{G}} \left| \frac{1}{n_c} \sum_{i=1}^{n_c} r_i \ell(-g(x_{c,i})) - \mathbb{E}_{x, r \sim q} [r \ell(-g(x))] \right| \\ &\quad + (1 - \beta) \sup_{g \in \mathcal{G}} \left| \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) - \mathbb{E}_{x, r \sim q} [\ell(-g(x))] \right| + \beta \sup_{g \in \mathcal{G}} \left| \frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) - \mathbb{E}_{x \sim p} [\ell(-g(x))] \right| \\ &\quad + \alpha \sup_{g \in \mathcal{G}} \left| \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(g(x_{c,i})) - \mathbb{E}_{x, r \sim q} [\ell(g(x))] \right| + \alpha \sup_{g \in \mathcal{G}} \left| \frac{1}{n_u} \sum_{i=1}^{n_u} \ell(g(x_{u,i})) - \mathbb{E}_{x \sim p} [\ell(g(x))] \right|. \end{aligned}$$

According to 2IWIL (Theorem 4.3 in Wu et al. (2019)), since the above six terms are the bounded differences with constants C_L/n_c , C_L/n_c , C_L/n_c , C_L/n_u , C_L/n_c , and C_L/n_u , respectively, we

can bound them with probability at least $1 - \delta/6$ as

$$\sup_{g \in \mathcal{G}} \left| \frac{1}{n_c} \sum_{i=1}^{n_c} r_i \ell(g(x_{c,i})) - \mathbb{E}_{x, r \sim q}[r \ell(g(x))] \right| \leq 8\rho_L \mathfrak{R}_{n_c}(\mathcal{G}) + 2C_L \sqrt{\frac{\log(12/\delta)}{2n_c}},$$

$$\sup_{g \in \mathcal{G}} \left| \frac{1}{n_c} \sum_{i=1}^{n_c} r_i \ell(-g(x_{c,i})) - \mathbb{E}_{x, r \sim q}[r \ell(-g(x))] \right| \leq 8\rho_L \mathfrak{R}_{n_c}(\mathcal{G}) + 2C_L \sqrt{\frac{\log(12/\delta)}{2n_c}},$$

$$\sup_{g \in \mathcal{G}} \left| \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(-g(x_{c,i})) - \mathbb{E}_{x, r \sim q}[\ell(-g(x))] \right| \leq 8\rho_L \mathfrak{R}_{n_c}(\mathcal{G}) + 2C_L \sqrt{\frac{\log(12/\delta)}{2n_c}},$$

$$\sup_{g \in \mathcal{G}} \left| \frac{1}{n_u} \sum_{i=1}^{n_u} \ell(-g(x_{u,i})) - \mathbb{E}_{x \sim p}[\ell(-g(x))] \right| \leq 8\rho_L \mathfrak{R}_{n_u}(\mathcal{G}) + 2C_L \sqrt{\frac{\log(12/\delta)}{2n_u}},$$

$$\sup_{g \in \mathcal{G}} \left| \frac{1}{n_c} \sum_{i=1}^{n_c} \ell(g(x_{c,i})) - \mathbb{E}_{x, r \sim q}[\ell(g(x))] \right| \leq 8\rho_L \mathfrak{R}_{n_c}(\mathcal{G}) + 2C_L \sqrt{\frac{\log(12/\delta)}{2n_c}},$$

$$\sup_{g \in \mathcal{G}} \left| \frac{1}{n_u} \sum_{i=1}^{n_u} \ell(g(x_{u,i})) - \mathbb{E}_{x \sim p}[\ell(g(x))] \right| \leq 8\rho_L \mathfrak{R}_{n_u}(\mathcal{G}) + 2C_L \sqrt{\frac{\log(12/\delta)}{2n_u}}.$$

In the end, we can bound the initial estimation error with probability of at least $1 - \delta$:

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

□

B DETAILS OF THE EXPERIMENTS AND ADDITIONAL EXPERIMENTS

B.1 HYPER-PARAMETERS SETTINGS AND TASK INFORMATION

All of our experiments are run on a single machine with 4 NVIDIA GeForce RTX 3080 GPUs. For the architectures of all neural networks, we employ two hidden layers of size 100 each, using Tanh as the activation function. Across all tasks, we utilize the same hyper-parameters as listed in Table 2. Table 3 shows the number of confidence data and unlabeled data used for each task, along with the cumulative rewards corresponding to the optimal and random policies.

Table 2: Hyper-parameters settings

Hyper-parameters	value
γ	0.995
τ (Generalized Advantage Estimation)	0.97
Batch size	5,000
Learning rate (Value network)	3×10^{-4}
Learning rate (Discriminator)	1×10^{-3}
Learning rate (Classifier)	3×10^{-4}
Optimizer	Adam

Table 3: Task information

Tasks	n_c	n_u	Optimal policy	Random policy
Ant-v2	120	480	4143.10	-72.30
HalfCheetah-v2	500	2000	3467.32	-288.44
Hopper-v2	20	80	3250.67	18.04
Pendulum-v1	200	800	-116.81	-1200.96
Swimmer-v2	5	20	348.99	2.31
Walker-v2	400	1600	3694.13	1.91

B.2 CLIP FUNCTION IN BSC

In the process of deriving Eq. (13), two terms with theoretical values of 0 are introduced: $\mathbb{E}_{x \sim p}[\alpha \ell(g(x))] - \mathbb{E}_{x \sim q}[\alpha \ell(g(x))]$ and $\mathbb{E}_{x \sim p}[\beta \ell(-g(x))] - \mathbb{E}_{x \sim q}[\beta \ell(-g(x))]$. If we directly minimize Eq. (13), the two terms may deviate far from 0. To mitigate this issue, we use the clip function to limit the sum of these two items to a neighborhood of 0. Specifically, we denote $R_1 = \mathbb{E}_{x \sim p}[\alpha \ell(g(x))] - \mathbb{E}_{x \sim q}[\alpha \ell(g(x))] + \mathbb{E}_{x \sim p}[\beta \ell(-g(x))] - \mathbb{E}_{x \sim q}[\beta \ell(-g(x))]$. According to Eq. (13), the final empirical risk can be written as

$$\widehat{R}_{\text{BSC},\ell}(g) = \widehat{R}_C(g) + \text{clip}_{[-\epsilon, \epsilon]} \widehat{R}_1,$$

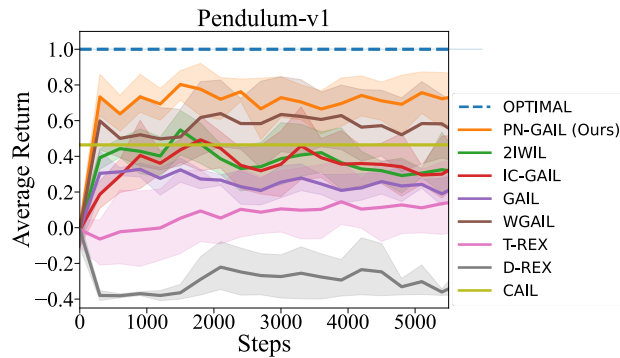
where

$$\widehat{R}_C(g) = \frac{1}{n_c} \sum_{i=1}^{n_c} [r_i \ell(g(x_{c,i})) + (1 - r_i) \ell(-g(x_{c,i}))].$$

In Walker2d-v2, the epoch for classifier training is set to 50000, with ϵ configured to 0.05. In other experiments, the epoch for classifier training is 25000, with ϵ configured to 0.01.

B.3 ADDITIONAL EXPERIMENTS

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199

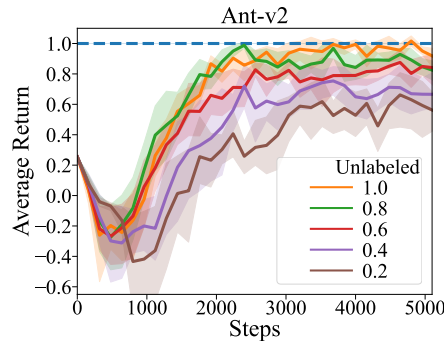


1200 Figure 5: Experiments with baselines including weighting-based and ranking-based methods.

1201
1202
1203
1204
1205

CAIL uses the average reward after training convergence, and the implementation is based on CAIL codebase. As shown in Fig. 5, PN-GAIL outperforms other baseline methods, achieving the highest return.

1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217

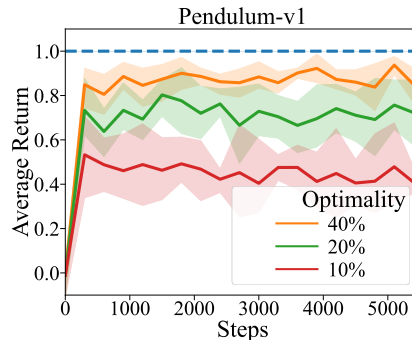


1218 Figure 6: Experiments on reducing the number of unlabeled confidence demonstrations.

1220
1221
1222
1223
1224

We gradually reduce the proportion of unlabeled confidence demonstrations. In Fig. 6, the numbers in the legend indicate the proportion of unlabeled data used as demonstrations. The performance of PN-GAIL improves as the amount of unlabeled data increases, illustrating how the use of unlabeled data can enhance the performance of our method.

1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237

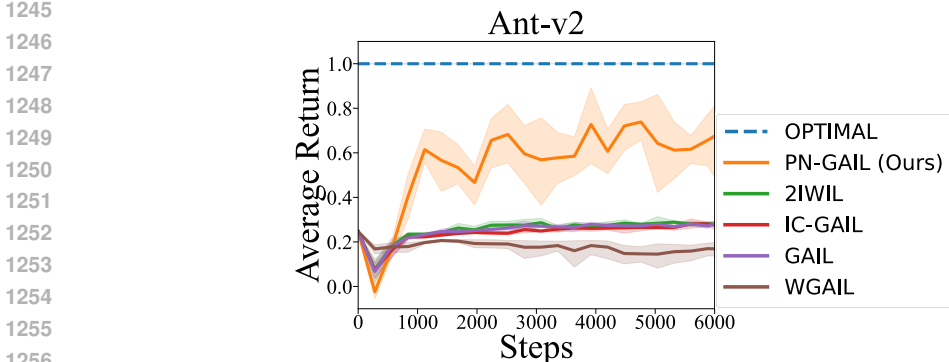


1238 Figure 7: Experiments under different optimalities of imperfect demonstrations.

1239
1240
1241

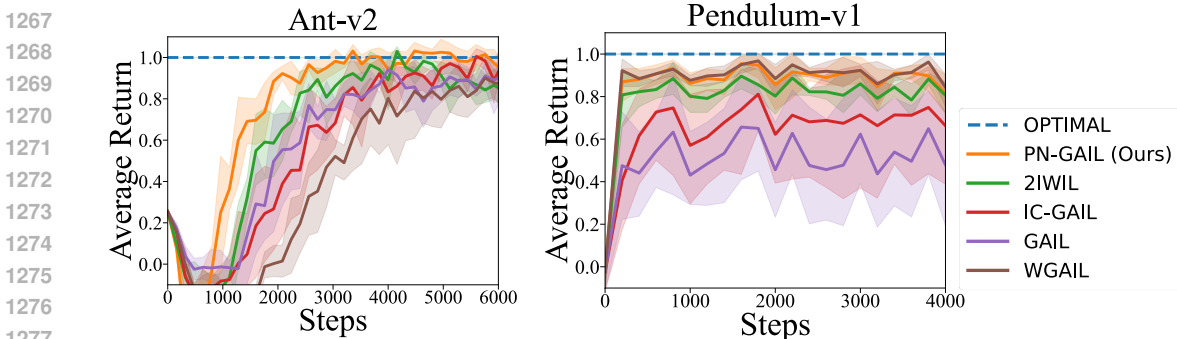
We evaluate the performance of PN-GAIL under different optimalities of imperfect demonstrations in the Pendulum-v1 environment. In Fig. 7, the numbers in the legend illustrate the proportion of demonstrations generated by the optimal policy when collecting datasets ($\pi_{\text{opt}} : \pi_1 = 2 : 3$,

1242 $\pi_{\text{opt}} : \pi_1 = 1 : 4$, $\pi_{\text{opt}} : \pi_1 = 1 : 9$). The higher the proportion of demonstrations generated by the
 1243 optimal policy is, the better the performance of PN-GAIL.
 1244



1257
1258 Figure 8: Experiments with an extreme demonstration ratio.

1259
1260 In the Ant-v2 environment, we train an optimal policy π_{opt} and an intermediate policy π_1 using
 1261 TRPO. We collect demonstrations with ratio of $\pi_{\text{opt}} : \pi_1 = 1 : 10$ and all demonstrations use
 1262 normalized rewards to annotate confidence scores. We then evaluate the performance of PN-GAIL
 1263 and other baseline methods using the collected demonstrations. As shown in Fig. 8, under this
 1264 extreme demonstration ratio, all other methods fail, the final result is close to GAIL, and only PN-
 1265 GAIL learns valid information from the extreme demonstrations.
 1266



1278
1279 Figure 9: Experiments when optimal demonstrations are dominant.

1280
1281 In the Ant-v2 and Pendulum-v1 environments, we conduct the experiments at the demonstration
 1282 ratio of $\pi_{\text{opt}} : \pi_1 = 2 : 1$. As shown in Fig. 9, it can be seen that when the optimal demonstrations
 1283 are dominant, PN-GAIL still shows robust and excellent performance.

1284 We conduct experiments with different numbers of non-optimal demonstrations in the Ant-v2 en-
 1285 vironment. In Fig. 10, the numbers in the legend indicate the proportion of non-optimal demon-
 1286 strations used as demonstrations, and the blue dotted line represents the performance of the optimal
 1287 policy. When the number of non-optimal demonstrations decreases, the performance of PN-GAIL
 1288 does not decrease significantly. We speculate that PN-GAIL does not take advantage of additional
 1289 non-optimal demonstrations now because the number of optimal demonstrations in given expert
 1290 demonstrations is sufficient to learn an optimal policy.

1291 Since the ranking-based methods including CAIL and f-IRL all require demonstrations to be stored
 1292 in the form of a trajectory, we only evaluate them in the Pendulum-v1 and Ant-v2 environments. The
 1293 dataset used in the Pendulum-v1 environment is the same as in Fig. 2, and the dataset used in the
 1294 Ant-v2 environment is the same as in Fig. 8. The implementation of CAIL, T-REX and D-REX is
 1295 based on the CAIL codebase, and the implementation of f-IRL is based on the f-IRL codebase. All
 methods use the average reward after training convergence. We construct trajectory rankings based

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

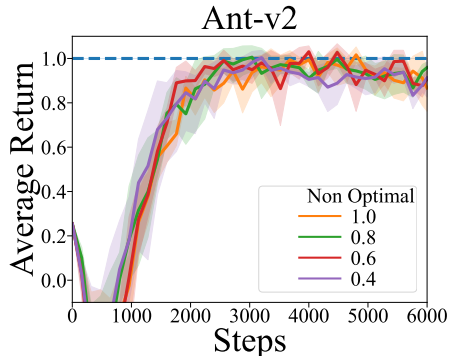


Figure 10: Experiments under different number of non-optimal demonstrations.

Table 4: Average returns of PN-GAIL, CAIL, ranking-based methods and f-IRL during training (only for evaluation).

Methods	Pendulum-v1	Ant-v2
PN-GAIL (Ours)	-465.978 ± 132.710	2960.458 ± 851.465
CAIL	-697.821 ± 91.459	2336.546 ± 735.378
T-REX	-1074.435 ± 258.078	-1858.025 ± 217.228
D-REX	-1532.958 ± 114.138	-2495.473 ± 220.703
FKL(f-IRL)	-698.064 ± 300.793	1563.796 ± 1372.482
RKL(f-IRL)	-603.760 ± 189.727	962.785 ± 816.547
JS(f-IRL)	-581.602 ± 185.096	882.026 ± 690.371
Optimal Policy	-116.81	4271.79

on confidence scores, and the average returns (only for evaluation) across all methods are shown in Table 4. PN-GAIL outperforms other methods, achieving the highest returns.

Table 5: The variance of the estimator $\hat{R}_{BSC,\ell}(g)$.

Var	Ant-v2	HalfCheetah-v2	Hopper-v2	Pendulum-v1	Swimmer-v2	Walker2d-v2
Origin	0.126±0.054	0.033±0.009	1.476±0.574	0.00064±0.00024	1.664±1.475	0.093±0.029
PN-GAIL(Ours)	0.100±0.050	0.025±0.006	1.412±0.566	0.00013±0.00004	0.007±0.015	0.084±0.028

Origin indicates the lack of α and β (i.e., $\alpha=0$ and $\beta=0$). The results in Table 5 show that the variance of PN-GAIL is consistently smaller, demonstrating the validity of the chosen values for α and β .

B.4 UNCROPPED FIGURES OF ANT-V2

In our experiments, we observe a decrease followed by an increase in performance within the Ant-v2 environment. For better comparison, we have cropped the figures of Ant-v2 in the main text, while the uncropped figures are presented below:

1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403

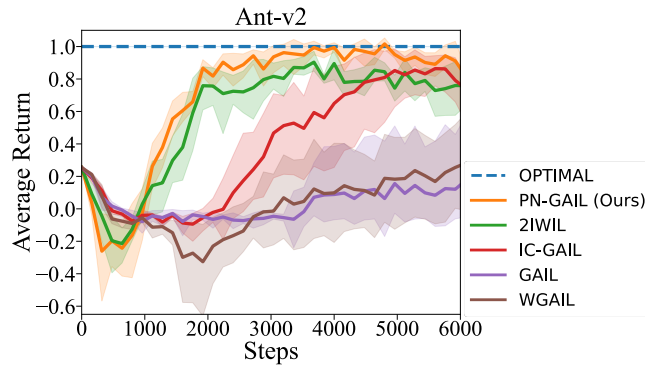


Figure 11: Normalized average returns of PN-GAIL and baseline methods during training.

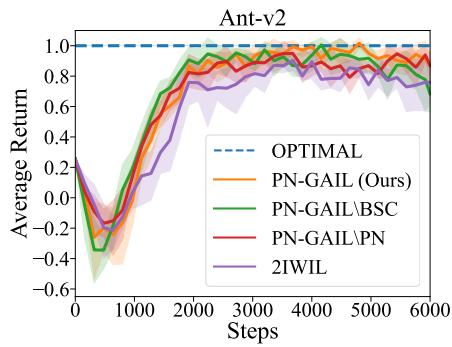


Figure 12: Normalized average returns of ablation experiments during training.