

DIGNet: Learning Decomposed Patterns in Representation Balancing for Treatment Effect Estimation

Anonymous authors

Paper under double-blind review

Abstract

Estimating treatment effects from observational data is often subject to a covariate shift problem incurred by selection bias. Recent research has sought to mitigate this problem by leveraging representation balancing methods that aim to extract balancing patterns from observational data and utilize them for outcome prediction. The underlying theoretical rationale is that minimizing the unobserved counterfactual error can be achieved through two principles: (I) reducing the risk associated with predicting factual outcomes and (II) mitigating the distributional discrepancy between the treated and controlled samples. **However, an inherent trade-off between the two principles can lead to a potential loss of information useful for factual outcome predictions and, consequently, deteriorating treatment effect estimations.** In this paper, we propose a novel representation balancing model, DIGNet, for treatment effect estimation. DIGNet incorporates two key components, PDIG and PPBR, which effectively mitigate the trade-off problem by improving one aforementioned principle without sacrificing the other. Specifically, PDIG captures more effective balancing patterns (Principle II) without affecting factual outcome predictions (Principle I), while PPBR enhances factual outcome prediction (Principle I) without affecting the learning of balancing patterns (Principle II). Our comprehensive ablation studies confirm the effectiveness of PDIG and PPBR in improving treatment effect estimation, and experimental results on benchmark datasets demonstrate the superior performance of our DIGNet model compared to baseline models.

1 Introduction

In the context of the ubiquity of personalized decision-making, causal inference has sparked a surge of research exploring causal machine learning in many disciplines, including economics and statistics (Wager & Athey, 2018; Athey & Wager, 2019; Farrell, 2015; Chernozhukov et al., 2018; Huang et al., 2021), healthcare (Qian et al., 2021; Bica et al., 2021a;b), and commercial applications (Guo et al., 2020a;b; Chu et al., 2021). The core of causal inference is to estimate *treatment effects*, which is closely related to the *factual outcomes* (observed outcomes) and *counterfactual outcomes*. The concept of the counterfactual outcome is closely linked to a fundamental hypothetical question: What would the outcome be if an alternative treatment were received? Answering this question is challenging because counterfactual outcomes are unobservable in reality, making it impossible to directly access ground-truth treatment effects from observational data. Consequently, an increasing amount of recent research has focused on developing innovative machine learning models that aim to enhance the estimation of counterfactual outcomes to obtain more accurate treatment effect estimates.

One of the challenges in estimating counterfactual outcomes lies in the *covariate shift* problem. In observational data, the population can be typically divided into two groups: (i) individuals who received treatment ($T = 1$), referred to as *treated samples* or *treatment samples*, and (ii) individuals who did not receive treatment ($T = 0$), referred to as *controlled samples* or *control samples*. The covariate shift problem indicates the difference between the distribution of covariate in the treated group and that in the controlled group, meaning $P(X|T = 1) \neq P(X|T = 0)$. This phenomenon is a result of the non-random treatment assignment mechanism, where the decision to receive treatment (e.g., heart medicine) is often determined by the

covariate (e.g., age). For example, people receiving heart medicine treatment tend to be much older compared to those who do not receive such treatment, because the doctor’s decision-making regarding whether to undergo heart medicine treatment highly depends on the patients’ age. Such a non-random treatment assignment is known as the selection bias phenomenon in the causal inference literature. Although the covariate shift arises from the association between covariate and treatment, this issue can significantly exacerbate the difficulty in inferring counterfactual outcomes, as traditional machine learning models can be invalid in estimating potential outcomes when a covariate shift is present (Yao et al., 2018; Hassanpour & Greiner, 2019a). Specifically, to infer the potential outcome Y^0 for treated ($T = 1$) samples, the conventional approach is to first train a model $\hat{\tau}^0(X)$ using controlled ($T = 0$) samples, and then utilize $\hat{\tau}^0(X)$ to predict Y^0 for treated ($T = 1$) samples. This approach, known as the T-learner in the causal inference literature (Curth & Van Der Schaar, 2023; Mahajan et al., 2024), becomes problematic because the training data (control samples) used for model training do not have the same distribution as the test data (treated samples), i.e., $P(X|T = 1) \neq P(X|T = 0)$. This violates the assumption in machine learning that training data and test data should be independent and identically distributed.

To alleviate the covariate shift problem, recent advancements in representation balancing research have explored the representation learning model, such as CounterFactual Regression Network (CFRNet) (Shalit et al., 2017), to estimate individual treatment effects (ITEs). These representation balancing models aim to extract balancing patterns from observational data and utilize these patterns to predict outcomes. The corresponding objective function is typically concerned with minimizing the empirical risk of factual outcomes while concurrently minimizing the distributional distance between the treatment and control groups in the representation space (Shalit et al., 2017; Johansson et al., 2022a). The underlying theoretical logic behind these studies is that minimizing counterfactual error can be achieved by two principles in the representation space: *(Principle I) minimizing the risk associated with factual outcome prediction*, and *(Principle II) reducing the distributional discrepancy between the treated and controlled samples*. The theoretical foundation and the classic CFRNet structure proposed in Shalit et al. (2017) have inspired many subsequent studies on representation balancing methods for treatment effect estimation, including Yao et al. (2018); Shi et al. (2019); Zhang et al. (2020); Hassanpour & Greiner (2019a); Assaad et al. (2021); Huang et al. (2022).

While the representation balancing framework provides a powerful tool to tackle the covariate shift issue, models based on the classic CFRNet structure (Figure 2(a)) still face a trade-off problem between the aforementioned two principles, because enforcing models to learn merely balancing patterns can undermine the predictive power of the outcome function (Zhang et al., 2020; Assaad et al., 2021; Huang et al., 2022). We will now discuss two cases to gain a deeper understanding of this trade-off phenomenon. (1) The case without representation balancing: In this case, the outcome functions are fitted by $Y^1 = \hat{\tau}^1(X^{treat})$ and $Y^0 = \hat{\tau}^0(X^{control})$ using treated and controlled samples, respectively. $\hat{\tau}^1(X^{treat})$ and $\hat{\tau}^0(X^{control})$ can be good estimates of factual outcomes based on the well-preserved pre-balancing information (group information). However, the estimated counterfactual outcomes $\hat{\tau}^0(X^{treat})$ and $\hat{\tau}^1(X^{control})$ can be problematic due to the presence of the covariate shift problem $P(X|T = 1) \neq P(X|T = 0)$, where the distribution of training data $P(X, Y|T = t)$ differs from that of the test data $P(X, Y|T = 1 - t)$ for $t \in \{0, 1\}$ ¹. (2) The case with representation balancing: In this case, the outcome functions are fitted by $Y^1 = \hat{h}^1(\Phi(X^{treat}))$ and $Y^0 = \hat{h}^0(\Phi(X^{control}))$ using treated and controlled samples, respectively. Using $\Phi(X)$ to fit factual outcomes can improve the accuracy of the counterfactual estimates $\hat{h}^0(\Phi(X^{treat}))$ and $\hat{h}^1(\Phi(X^{control}))$, because representation balancing enforces the distributions of treated and controlled samples to be as close as possible in the representation space. As a result, representation balancing effectively tackles the covariate shift issue, resulting in training data and test data following the same distribution². However, executing representation balancing can inevitably lead to a loss of outcome-predictive information in $\Phi(X)$. This occurs naturally as Φ becomes insensitive to the treatment variable, thereby sacrificing pre-balancing information (group information) that contributes to factual outcome predictions. To illustrate the negative impact of losing pre-balancing information in balanced representations, we present a motivating example below.

¹By unconfoundedness, we have $P(Y|X, T = t) = P(Y|X, T = 1 - t)$. Due to the covariate shift $P(X|T = t) \neq P(X|T = 1 - t)$, we have $P(Y|X, T = t)P(X|T = t) \neq P(Y|X, T = 1 - t)P(X|T = 1 - t)$, i.e., $P(X, Y|T = t) \neq P(X, Y|T = 1 - t)$.

²By one-to-one and invertible properties of Φ and unconfoundedness, we have $P(Y|\Phi(X), T = t) = P(Y|\Phi(X), T = 1 - t)$. Given an ideal representation balancing $P(\Phi(X)|T = t) = P(\Phi(X)|T = 1 - t)$, we have $P(Y|\Phi(X), T = t)P(\Phi(X)|T = t) = P(Y|\Phi(X), T = 1 - t)P(\Phi(X)|T = 1 - t)$, i.e., $P(\Phi(X), Y|T = t) = P(\Phi(X), Y|T = 1 - t)$.

Motivating example. Suppose there is a vaccine available to prevent a certain disease. We define X as the covariate, $T = 1$ as the treatment (receiving the vaccine), $T = 0$ as the control (not receiving the vaccine), and Y as the outcome (the level of specific antibodies). Assume that the outcome is determined by $Y = T \exp(X) + (1 - T) \cdot 0 = T \exp(X)$, which means that if an individual receives the treatment, the level of antibodies will be $y = \exp(x)$; otherwise, it will be $y = 0$. In observational data, the treatment is assigned based on the covariate of each individual. The left graph of Figure 1 illustrates the distributions of X in the treated and control groups. We observe that individuals with positive x values are more likely to receive the vaccine, resulting in a higher level of antibodies. Given some sample i , a well trained model first determines whether the sample is more likely to be in the treatment or control group based on its covariate $X_i = x$. If it determines the sample to be in the treatment group, the model then predicts $y = \exp(x)$; otherwise, it predicts $y = 0$. For example, if $x = 1$, the model would classify the sample as more likely to be in the treatment group and predict $y = e$. Therefore, in this case, the pre-balancing covariate remain informative for predicting the outcome. Now, let's consider a representation function Φ that achieves improper representation balancing between the treated and control samples. The right graph of Figure 1 shows the distributions of $\Phi(X)$ in the treated and control groups. In this case, it becomes challenging for a model to accurately predict Y using $\Phi(X)$ because the model may become confused about whether a sample is more likely to receive the treatment or the control. Consequently, improperly balanced representations can lead to a loss of outcome-predictive information.

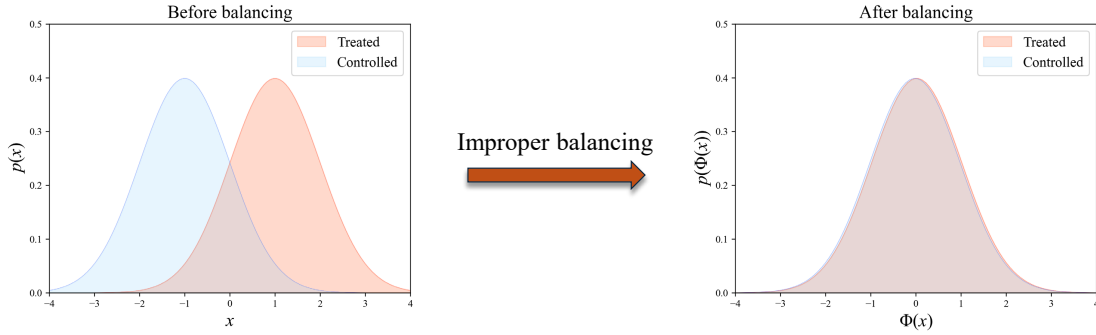


Figure 1: Motivating example for illustrating the trade-off between outcome prediction and representation balancing.

The above discussions and motivating example illustrate the inherent trade-off problem between outcome prediction (Principle I) and representation balancing (Principle II), which arises due to the fact that representation balancing models alleviate covariate shift at the expense of factual outcome prediction. This motivates us to ponder an interesting question: considering the inherent trade-off between the two principles, *is it possible to explore a scheme that enhances one principle without sacrificing the other?* More specifically, can we explore improving treatment effect estimation through the following two paths: *(Path I) learning more effective balancing patterns without sacrificing factual outcome prediction* and *(Path II) enhancing factual outcome prediction without sacrificing the learning of balancing patterns?*

In this paper, we propose a novel representation balancing model, **DIGNet** (Section 4.2.2), which is a neural Network that incorporates *Decomposed patterns* with *Individual propensity confusion* and *Group distance minimization*. The term of decomposed patterns denotes distinct components disentangled from some specific representations in DIGNet (Section 4.2). The individual propensity confusion aspect of DIGNet aims to learn representations that are difficult to utilize for characterizing the propensity of each individual being treated or controlled (Section 4.1.2), and the corresponding theoretical foundation is based on our derived \mathcal{H} -divergence guided counterfactual and ITE error bounds (Section 3.2). The group distance minimization aspect of DIGNet focuses on learning representations that minimize the distance between the treated and controlled groups (Section 4.1.1), and the corresponding theoretical foundation is supported by previous work (Shalit et al., 2017) on Wasserstein distance guided counterfactual and ITE error bounds (Section 3.1). To illustrate and explain these introduced concepts, we provide Figure 2 which visually depicts the proposed components and their relationships.

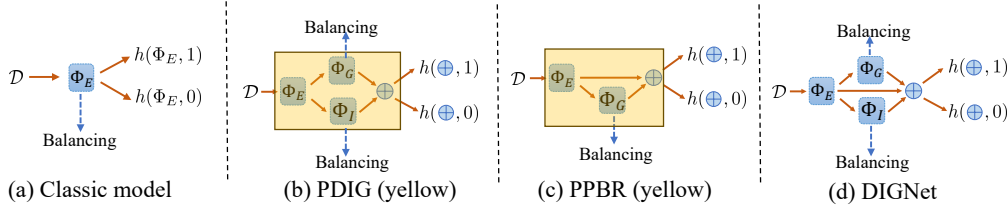


Figure 2: (a): The classic model (e.g., GNet in Section 4.1.1 and INet in Section 4.1.2) maps the original data \mathcal{D} into representations Φ_E to achieve representation balancing. The balanced representations are referred to as *balancing patterns*. These balancing patterns are also used for outcome prediction. (b): The PDIG (Section 4.2.1) is illustrated as the yellow part, where balancing patterns are decomposed into two distinct components, Φ_G and Φ_I . Φ_G serves for *group distance minimization* (Section 4.1.1) and Φ_I serves for *individual propensity confusion* (Section 4.1.2). The balancing patterns Φ_G and Φ_I are concatenated for predicting outcomes. (c): The PPBR (Section 4.2.1) is represented by the yellow section, where Φ_E is used for feature extraction and Φ_G is used for representation balancing. Here representations are decomposed into *pre-balancing patterns* Φ_E and balancing patterns Φ_G . Φ_E and Φ_G are concatenated for predicting outcomes. (d): The proposed model DIGNet (Section 4.2.2) integrates both PDIG and PPBR. Specifically, DIGNet decomposes balancing patterns into two distinct components, Φ_G and Φ_I . The outcome predictors are further formed by concatenating Φ_G , Φ_I , and pre-balancing patterns Φ_E .

Contributions. Our main contributions are summarized as follows:

1. We derive theoretical upper bounds for counterfactual error and ITE error based on \mathcal{H} -divergence (Section 3.2). In particular, this theoretical foundation highlights the important role of propensity score for representation balancing models, connecting the representation balancing with the concept of individual propensity confusion.
2. **We suggest learning decomposed patterns in representation balancing models (Section 4.2.1) to mitigate the trade-off problem rooted in classic causal representation balancing models.** First, we propose a **PDIG** method (Figure 2(b)), which aims to learn **P**atterns **D**ecomposed with **I**ndividual propensity confusion and **G**roup distance minimization to improve treatment effect estimation through Path I. Second, we propose a **PPBR** method (Figure 2(c)), which aims to learn **P**atterns of **P**re-balancing and **B**alancing **R**epresentations to improve treatment effect estimation through Path II.
3. Building upon PDIG and PPBR, we propose a novel representation balancing model, DIGNet (Figure 2(d)), for treatment effect estimation. In Section 5, ablation studies verify the efficacy of PDIG and PPBR in improving ITE estimation through Path I and Path II, respectively. Furthermore, experimental results on benchmark datasets demonstrate that DIGNet surpasses the performance of baseline models in terms of treatment effect estimation.

1.1 Related Work

The presence of a covariate shift problem stimulates the line of representation balancing works (Johansson et al., 2016; Shalit et al., 2017; Johansson et al., 2022a). These works aim to balance the distributions of representations between treated and controlled groups and simultaneously try to maintain representations predictive of factual outcomes. This idea is closely connected with domain adaptation. In particular, the ITE error bound based on Wasserstein distance is similar to the generalization bound in Ben-David et al. (2010); Long et al. (2014); Shen et al. (2018). In addition to Wasserstein distance based model, this paper derives a new ITE error bound based on \mathcal{H} -divergence (Ben-David et al., 2006; 2010; Ganin et al., 2016). **In addition to its connection to domain adaptation, causal representation learning is also linked to the field of fair representation learning, which aims to ensure that machine learning algorithms make fair decisions by learning fair representations. The main goal of these studies is to enforce a classification**

model to be less sensitive to certain sensitive variables when the representations of different groups are sufficiently similar (Zemel et al., 2013; Edwards & Storkey, 2015; Beutel et al., 2017; Madras et al., 2018; Zhang et al., 2018; Adel et al., 2019; Feng et al., 2019; Zhao et al., 2019a; Zhao & Gordon, 2022). Notably, the original idea of adversarial learned fair representations in Edwards & Storkey (2015) is also motivated by the domain adaptation work (Ben-David et al., 2006; 2010; Ganin et al., 2016), sharing a similar motivation to our utilization of INet, which relies on \mathcal{H} -divergence guided error bounds for ITE estimation. Moreover, Wasserstein distance has also been employed for learning fair representations in Jiang et al. (2020).

Another recent line of causal representation learning literature investigates efficient neural network structures for treatment effect estimation. Kuang et al. (2017); Hassanpour & Greiner (2019b) extract the original covariates into treatment-specific factors, outcome-specific factors, and confounding factors; X-learner (Künzel et al., 2019) and R-learner (Nie & Wager, 2021) are developed beyond the classic S-learner and T-learner; Curth & van der Schaar (2021) leverage structures for end-to-end learners to counteract the inductive bias towards treatment effect estimation, which is motivated by Makar et al. (2020). There are some other deep neural network models that have been employed in treatment effect estimation Louizos et al. (2017); Yao et al. (2018); Yoon et al. (2018); Shi et al. (2019); Du et al. (2021). To ensure comparability and consistency, we rigorously follow the same framework as these causal inference works. The causal graph in these studies satisfies the standard setup $T \leftarrow X \rightarrow Y$ and $T \rightarrow Y$. Additionally, it is also worth noting that there are many other causal inference works exploring treatment effect estimation under more complex causal graphs. For instance, studies such as Kallus et al. (2019); Jesson et al. (2021); Miao et al. (2023) specifically tackle the treatment effect estimation when unobserved confounders U present. In this case, the causal graph setup extends to $T \leftarrow X \rightarrow Y$, $T \rightarrow Y$, $T \leftarrow U \rightarrow Y$. A recent work (Cao et al., 2023) further expands this static causal graph to a dynamic setting. Moreover, some studies such as Angrist et al. (1996); Burgess et al. (2017); Wu et al. (2022); Yuan et al. (2023) estimate treatment effects with instrumental variables I involved. In this case, there are various causal graph setups such as $T \leftarrow X \rightarrow Y$, $I \rightarrow T \rightarrow Y$, and $T \leftarrow I \rightarrow Y$. More complex causal graph settings (Nogueira et al., 2021; Vowels et al., 2022; Zanga et al., 2022) have been studied with the development of Directed Graphical Models (Pearl, 2009), which represents another significant research direction known as causal discovery.

Our method is highly motivated by the trade-off problem between outcome prediction and representation balancing. In the causal representation learning literature, a similar trade-off phenomenon has been noticed by Zhang et al. (2020); Assaad et al. (2021); Huang et al. (2022), where the researchers argue that highly-balanced representations can have adverse effects on outcome modeling. However, the explanations for this phenomenon and its connections with other related literature are not extensively provided in their work. We highlight that the trade-off between outcome prediction and representation balancing is also connected with trade-offs observed in other research domains. In representation balancing models, representation balancing helps improve the model’s ability to generalize to counterfactual estimates. However, representation balancing can potentially sacrifice information necessary for predicting factual outcomes. In supervised machine learning, penalizing model complexity during model training helps the model to learn simpler patterns, thereby promoting generalization ability (reducing its variance) to unseen data. However, a bias-variance trade-off occurs because less flexible models tend to exhibit higher bias in training data (Geman et al., 1992; Domingos, 2000; Valentini & Dietterich, 2004; Yang et al., 2020). In the literature of domain adaptation (Shen et al., 2018; Zhao et al., 2019b), transfer learning (Long et al., 2015; 2017), out-of-distribution detection (Kumar et al., 2021; 2022), and fair representation learning (Zliobaite, 2015; Hardt et al., 2016), enforcing a model to capture proxy features that are domain-invariant helps the model to generalize well to unseen target (also known as out-of-distribution) data. However, a trade-off between classification accuracy and domain-invariance (or fairness in fair representation learning literature) occurs because the pursuit of domain-invariant features may lead to a loss of classification accuracy on the source (also known as in-distribution) data (Zhao et al., 2019a; Zhao & Gordon, 2022; Zhao et al., 2022).

2 Preliminaries

Notations. Suppose there are N i.i.d. random variables $\mathcal{D} = \{(\mathbf{X}_i, T_i, Y_i)\}_{i=1}^N$ with observed realizations $\{(\mathbf{x}_i, t_i, y_i)\}_{i=1}^N$, where there are N_1 treated units and N_0 controlled units. For each unit i , $\mathbf{X}_i \in \mathcal{X} \subset \mathbb{R}^d$ denotes d -dimensional covariates and $T_i \in \{0, 1\}$ denotes the binary treatment, with $e(\mathbf{x}_i) := p(T_i = 1 \mid \mathbf{X}_i =$

\mathbf{x}_i) defined as the propensity score (Rosenbaum & Rubin, 1983). Potential outcome framework (Rubin, 2005) defines the potential outcomes $Y^1, Y^0 \in \mathcal{Y} \subset \mathbb{R}$ for treatment $T = 1$ and $T = 0$, respectively. We let the observed outcome (factual outcome) be $Y = T \cdot Y^1 + (1 - T) \cdot Y^0$, and the unobserved outcome (counterfactual outcome) be $Y = T \cdot Y^0 + (1 - T) \cdot Y^1$. For $t \in \{0, 1\}$, let $\tau^t(\mathbf{x}) := \mathbb{E}[Y^t | \mathbf{X} = \mathbf{x}]$ be a function of Y^t w.r.t. \mathbf{X} , then our goal is to estimate the individual treatment effect (ITE) $\tau(\mathbf{x}) := \mathbb{E}[Y^1 - Y^0 | \mathbf{X} = \mathbf{x}] = \tau^1(\mathbf{x}) - \tau^0(\mathbf{x})$ ¹, and the average treatment effect (ATE) $\tau_{ATE} := \mathbb{E}[Y^1 - Y^0] = \int_{\mathcal{X}} \tau(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$. The introduced concepts PPBR and PDIG are illustrated in Figure 2, and the necessary representation functions Φ_E , Φ_G and Φ_I , as well as different model structures, are illustrated in Figure 3. Throughout the paper, we refer to patterns as meaningful representations. For instance, decomposed patterns are distinct components disentangled from some specific representations.

2.1 Problem setup

In causal representation balancing works, we denote representation space by $\mathcal{R} \subset \mathbb{R}^d$, and $\Phi : \mathcal{X} \rightarrow \mathcal{R}$ is assumed to be a twice-differentiable, one-to-one and invertible function with its inverse $\Psi : \mathcal{R} \rightarrow \mathcal{X}$ such that $\Psi(\Phi(\mathbf{x})) = \mathbf{x}$ ³. The densities of the treated and controlled covariates are denoted by $p_{\mathbf{x}}^{T=1} = p^{T=1}(\mathbf{x}) := p(\mathbf{x} | T = 1)$ and $p_{\mathbf{x}}^{T=0} = p^{T=0}(\mathbf{x}) := p(\mathbf{x} | T = 0)$, respectively. Correspondingly, the densities of the treated and controlled covariates in the representation space are denoted by $p_{\Phi}^{T=1} = p_{\Phi}^{T=1}(\mathbf{r}) := p_{\Phi}(\mathbf{r} | T = 1)$ and $p_{\Phi}^{T=0} = p_{\Phi}^{T=0}(\mathbf{r}) := p_{\Phi}(\mathbf{r} | T = 0)$, respectively.

Our study is based on the potential outcome framework (Rubin, 2005). Assumption 1 states standard and necessary assumptions to ensure treatment effects are identifiable. Before proceeding with theoretical analysis, we also present some necessary terms and definitions in Definition 1.

Assumption 1 (Consistency, Overlap, and Unconfoundedness). *Consistency: If the treatment is t , then the observed outcome equals Y^t . Overlap: The propensity score is bounded away from 0 to 1, i.e., $0 < e(\mathbf{x}) < 1$. Unconfoundedness: $Y^t \perp\!\!\!\perp T | \mathbf{X}$, $\forall t \in \{0, 1\}$.*

Definition 1. Let $h : \mathcal{R} \times \{0, 1\} \rightarrow \mathcal{Y}$ be an hypothesis defined over the representation space \mathcal{R} such that $h(\Phi(\mathbf{x}), t)$ estimates y^t , and $L : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ be a loss function (e.g., the squared loss $L(y, y') = (y - y')^2$ or the absolute loss $L(y, y') = |y - y'|$). If we define the expected loss for (\mathbf{x}, t) as $\ell_{h, \Phi}(\mathbf{x}, t) = \int_{\mathcal{Y}} L(y^t, h(\Phi(\mathbf{x}), t)) p(y^t | \mathbf{x}) dy^t$, we then have factual and counterfactual errors, as well as them on the treated and controlled:

$$\begin{aligned} \epsilon_F(h, \Phi) &= \int_{\mathcal{X} \times \{0, 1\}} \ell_{h, \Phi}(\mathbf{x}, t) p(\mathbf{x}, t) d\mathbf{x} dt, & \epsilon_{CF}(h, \Phi) &= \int_{\mathcal{X} \times \{0, 1\}} \ell_{h, \Phi}(\mathbf{x}, t) p(\mathbf{x}, 1 - t) d\mathbf{x} dt, \\ \epsilon_F^{T=1}(h, \Phi) &= \int_{\mathcal{X}} \ell_{h, \Phi}(\mathbf{x}, 1) p^{T=1}(\mathbf{x}) d\mathbf{x}, & \epsilon_F^{T=0}(h, \Phi) &= \int_{\mathcal{X}} \ell_{h, \Phi}(\mathbf{x}, 0) p^{T=0}(\mathbf{x}) d\mathbf{x}, \\ \epsilon_{CF}^{T=1}(h, \Phi) &= \int_{\mathcal{X}} \ell_{h, \Phi}(\mathbf{x}, 1) p^{T=0}(\mathbf{x}) d\mathbf{x}, & \epsilon_{CF}^{T=0}(h, \Phi) &= \int_{\mathcal{X}} \ell_{h, \Phi}(\mathbf{x}, 0) p^{T=1}(\mathbf{x}) d\mathbf{x}. \end{aligned}$$

If we let $f(\mathbf{x}, t)$ be $h(\Phi(\mathbf{x}), t)$, where $f : \mathcal{X} \times \{0, 1\} \rightarrow \mathcal{Y}$ is a prediction function for outcome, then the estimated ITE over f is defined as $\hat{\tau}_f(\mathbf{x}) := f(\mathbf{x}, 1) - f(\mathbf{x}, 0)$. We can measure the error in ITE estimation with the metric, Precision in the expected Estimation of Heterogeneous Effect (PEHE):

$$\epsilon_{PEHE}(f) = \int_{\mathcal{X}} L(\hat{\tau}_f(\mathbf{x}), \tau(\mathbf{x})) p(\mathbf{x}) d\mathbf{x}. \quad (1)$$

Here, $\epsilon_{PEHE}(f)$ can also be denoted by $\epsilon_{PEHE}(h, \Phi)$ if we let $f(\mathbf{x}, t)$ be $h(\Phi(\mathbf{x}), t)$.

¹The term $\mathbb{E}[Y^1 - Y^0 | \mathbf{X} = \mathbf{x}]$ is commonly known as the Conditional Average Treatment Effect (CATE). In order to maintain consistency with the notion used in the existing causal representation balancing literature, e.g., Shalit et al. (2017), we refer to this term as ITE throughout this paper. Note that the original definition of ITE for the i -th individual is commonly expressed as the difference between their potential outcomes, represented as $Y_i^1 - Y_i^0$.

³Theoretically, the invertibility is necessary for deriving the upper bounds of ITE error, specifically for equation 39 and equation 47. However, the invertibility can be hard to verify in practice (Johansson et al., 2022b).

3 Theoretical Results

In this section, we first prove ϵ_{PEHE} is bounded by ϵ_F and ϵ_{CF} in Lemma 1. Next, we revisit the upper bound for Wasserstein distance guided representation balancing method in Section 3.1. Furthermore, we state the new theoretical results concerning \mathcal{H} -divergence guided representation balancing method in Section 3.2.

Lemma 1. *Let functions h and Φ be as defined in Definition 1. Recall that $\tau^t(\mathbf{x}) = \mathbb{E}[Y^t | \mathbf{X} = \mathbf{x}]$. Define $\sigma_y^2 = \min\{\sigma_{y^t}^2(p(\mathbf{x}, t)), \sigma_{y^t}^2(p(\mathbf{x}, 1-t))\}$ and $A_y = \max\{A_{y^t}(p(\mathbf{x}, t)), A_{y^t}(p(\mathbf{x}, 1-t))\} \forall t \in \{0, 1\}$, where $\sigma_{y^t}^2(p(\mathbf{x}, t)) = \int_{\mathcal{X} \times \{0,1\} \times \mathcal{Y}} (y^t - \tau^t(\mathbf{x}))^2 p(y^t | \mathbf{x}) p(\mathbf{x}, t) dy^t d\mathbf{x} dt$ and $A_{y^t}(p(\mathbf{x}, t)) = \int_{\mathcal{X} \times \{0,1\} \times \mathcal{Y}} |y^t - \tau^t(\mathbf{x})| p(y^t | \mathbf{x}) p(\mathbf{x}, t) dy^t d\mathbf{x} dt \forall t \in \{0, 1\}$.*

Let loss function L be the squared loss. Then we have:

$$\epsilon_{PEHE}(h, \Phi) \leq 2(\epsilon_{CF}(h, \Phi) + \epsilon_F(h, \Phi) - 2\sigma_y^2). \quad (2)$$

Let loss function L be the absolute loss. Then we have:

$$\epsilon_{PEHE}(h, \Phi) \leq \epsilon_{CF}(h, \Phi) + \epsilon_F(h, \Phi) + 2A_y. \quad (3)$$

Lemma 1 reveals that the ITE error ϵ_{PEHE} is closely connected with the factual error ϵ_F and counterfactual ϵ_{CF} , as well as a constant σ_y^2 (or A_y) that is unrelated with functions h and Φ . Here, σ_y^2 is the smaller value of the variance in Y^t w.r.t. the distribution $p(\mathbf{x}, t)$ and the variance in Y^{1-t} w.r.t. $p(\mathbf{x}, 1-t)$, and A_y is the larger value of the absolute deviation in Y^t w.r.t. the distribution $p(\mathbf{x}, t)$ and the absolute deviation in Y^{1-t} w.r.t. the distribution $p(\mathbf{x}, 1-t)$. The proof of Lemma 1 is deferred to Section A.1. Note that equation (2) corresponds to the result presented in Shalit et al. (2017), while equation (3) is our new result, which supplements the case when L denotes the absolute loss.

3.1 Wasserstein Distance Guided Error Bounds

Previous causal learning models commonly adopt the Wasserstein distance guided approach to seek representation balancing. In this subsection, we first give the definition of Wasserstein distance (Cuturi & Doucet, 2014) by introducing the Integral Probability Metric (IPM) (Sriperumbudur et al., 2012) defined in Definition 2. Then we state the theorem regarding the upper bounds for counterfactual error ϵ_{CF} and ITE error ϵ_{PEHE} using Wasserstein distance in Theorem 1.

Definition 2. *Let \mathcal{G} be a function family consisting of functions $g : \mathcal{S} \rightarrow \mathbb{R}$. For a pair of distributions p_1, p_2 over \mathcal{S} , the Integral Probability Metric is defined as*

$$IPM_{\mathcal{G}}(p_1, p_2) := \sup_{g \in \mathcal{G}} \left| \int_{\mathcal{S}} g(s)(p_1(s) - p_2(s)) ds \right|.$$

If \mathcal{G} is the family of 1-Lipschitz functions, we can obtain the so-called 1-Wasserstein distance, denoted by $Wass(p_1, p_2)$. Next, we present the bounds for counterfactual error ϵ_{CF} and ITE error ϵ_{PEHE} using Wasserstein distance in Theorem 1.

Theorem 1. *Let $\Phi : \mathcal{X} \rightarrow \mathcal{R}$ be an invertible representation with Ψ being its inverse. Define $\sigma_y^2 = \min\{\sigma_{y^t}^2(p(\mathbf{x}, t)), \sigma_{y^t}^2(p(\mathbf{x}, 1-t))\}$ and $A_y = \max\{A_{y^t}(p(\mathbf{x}, t)), A_{y^t}(p(\mathbf{x}, 1-t))\} \forall t \in \{0, 1\}$, where $\sigma_{y^t}^2(p(\mathbf{x}, t)) = \int_{\mathcal{X} \times \{0,1\} \times \mathcal{Y}} (y^t - \tau^t(\mathbf{x}))^2 p(y^t | \mathbf{x}) p(\mathbf{x}, t) dy^t d\mathbf{x} dt$ and $A_{y^t}(p(\mathbf{x}, t)) = \int_{\mathcal{X} \times \{0,1\} \times \mathcal{Y}} |y^t - \tau^t(\mathbf{x})| p(y^t | \mathbf{x}) p(\mathbf{x}, t) dy^t d\mathbf{x} dt \forall t \in \{0, 1\}$. Let $p_{\Phi}^{T=1}(\mathbf{r})$, $p_{\Phi}^{T=0}(\mathbf{r})$ be as defined before, $h : \mathcal{R} \times \{0, 1\} \rightarrow \mathcal{Y}$, $u := \Pr(T = 1)$ and \mathcal{G} be the family of 1-Lipschitz functions. Assume there exists a constant $B_{\Phi} \geq 0$, such that for $t \in \{0, 1\}$, the function $g_{\Phi, h}(\mathbf{r}, t) := \frac{1}{B_{\Phi}} \cdot \ell_{h, \Phi}(\Psi(\mathbf{r}), t) \in \mathcal{G}$. Given a loss function L , we have*

$$\epsilon_{CF}(h, \Phi) \leq (1-u) \cdot \epsilon_F^{T=1}(h, \Phi) + u \cdot \epsilon_F^{T=0}(h, \Phi) + B_{\Phi} \cdot Wass(p_{\Phi}^{T=1}, p_{\Phi}^{T=0}). \quad (4)$$

Let loss function L be the squared loss. Then we have:

$$\epsilon_{PEHE}(h, \Phi) \leq 2(\epsilon_F^{T=1}(h, \Phi) + \epsilon_F^{T=0}(h, \Phi) + B_{\Phi} \cdot Wass(p_{\Phi}^{T=1}, p_{\Phi}^{T=0}) - 2\sigma_y^2). \quad (5)$$

Let loss function L be the absolute loss. Then we have:

$$\epsilon_{PEHE}(h, \Phi) \leq \epsilon_F^{T=1}(h, \Phi) + \epsilon_F^{T=0}(h, \Phi) + B_{\Phi} \cdot Wass(p_{\Phi}^{T=1}, p_{\Phi}^{T=0}) + 2A_y. \quad (6)$$

Theorem 1 reveals that the ITE error is closely tied to the factual error ϵ_F and the Wasserstein distance between treated and controlled groups in the representation space. This theorem provides a theoretical foundation for representation balancing models based on group distance minimization (Section 4.1.1). The proof of Theorem 1 is deferred to Section A.2. Note that equation (5) corresponds to the result presented in Shalit et al. (2017), while equation (6) is our new result, which supplements the case when L denotes the absolute loss.

3.2 \mathcal{H} -divergence Guided Error Bounds

In most representation balancing literature, the models mainly rely on Wasserstein distance guided error bounds as discussed in Section 3.1. In this subsection, we will focus on establishing \mathcal{H} -divergence guided error bounds for counterfactual and ITE estimations in representation balancing approach. We first give the definition of \mathcal{H} -divergence (Ben-David et al., 2006) in Definition 3. Then we state the theorem regarding the upper bounds for counterfactual error ϵ_{CF} and ITE error ϵ_{PEHE} using \mathcal{H} -divergence in Theorem 2.

Definition 3. Given a pair of distributions p_1, p_2 over \mathcal{S} , and a hypothesis binary function class \mathcal{H} , the \mathcal{H} -divergence between p_1 and p_2 is defined as

$$d_{\mathcal{H}}(p_1, p_2) := 2 \sup_{\eta \in \mathcal{H}} |Pr_{p_1}[\eta(s) = 1] - Pr_{p_2}[\eta(s) = 1]|. \quad (7)$$

Theorem 2. Let $\Phi : \mathcal{X} \rightarrow \mathcal{R}$ be an invertible representation with Ψ being its inverse. Define $\sigma_y^2 = \min\{\sigma_{y^t}^2(p(\mathbf{x}, t)), \sigma_{y^t}^2(p(\mathbf{x}, 1 - t))\}$ and $A_y = \max\{A_{y^t}(p(\mathbf{x}, t)), A_{y^t}(p(\mathbf{x}, 1 - t))\} \forall t \in \{0, 1\}$, where $\sigma_{y^t}^2(p(\mathbf{x}, t)) = \int_{\mathcal{X} \times \{0, 1\} \times \mathcal{Y}} (y^t - \tau^t(\mathbf{x}))^2 p(y^t | \mathbf{x}) p(\mathbf{x}, t) dy^t d\mathbf{x} dt$ and $A_{y^t}(p(\mathbf{x}, t)) = \int_{\mathcal{X} \times \{0, 1\} \times \mathcal{Y}} |y^t - \tau^t(\mathbf{x})| p(y^t | \mathbf{x}) p(\mathbf{x}, t) dy^t d\mathbf{x} dt \forall t \in \{0, 1\}$. Let $p_{\Phi}^{T=1}(\mathbf{r}), p_{\Phi}^{T=0}(\mathbf{r})$ be as defined before, $h : \mathcal{R} \times \{0, 1\} \rightarrow \mathcal{Y}$, $u := Pr(T = 1)$ and \mathcal{H} be the family of binary functions. Assume that there exists a constant $K \geq 0$ such that $\int_{\mathcal{Y}} L(y, y') dy \leq K \forall y' \in \mathcal{Y}$. Given a loss function L , we have

$$\epsilon_{CF}(h, \Phi) \leq (1 - u) \cdot \epsilon_F^{T=1}(h, \Phi) + u \cdot \epsilon_F^{T=0}(h, \Phi) + \frac{K}{2} d_{\mathcal{H}}(p_{\Phi}^{T=1}, p_{\Phi}^{T=0}). \quad (8)$$

Let loss function L be the squared loss. Then we have:

$$\epsilon_{PEHE}(h, \Phi) \leq 2(\epsilon_F^{T=1}(h, \Phi) + \epsilon_F^{T=0}(h, \Phi) + \frac{K}{2} d_{\mathcal{H}}(p_{\Phi}^{T=1}, p_{\Phi}^{T=0}) - 2\sigma_y^2). \quad (9)$$

Let loss function L be the absolute loss. Then we have:

$$\epsilon_{PEHE}(h, \Phi) \leq \epsilon_F^{T=1}(h, \Phi) + \epsilon_F^{T=0}(h, \Phi) + \frac{K}{2} d_{\mathcal{H}}(p_{\Phi}^{T=1}, p_{\Phi}^{T=0}) + 2A_y. \quad (10)$$

Theorem 2 reveals that the ITE error is closely connected with the factual error ϵ_F and the \mathcal{H} -divergence between treated and controlled samples in the representation space. This new theoretical result provides a theoretical foundation for representation balancing models based on individual propensity confusion (Section 4.1.2). The proof of Theorem 2 is deferred to Section A.3.

4 Method

In the preceding section, we have stated the theoretical foundations for representation balancing methods, which are the Wasserstein distance guided error bounds (results in Shalit et al. (2017)) and \mathcal{H} -divergence guided error bounds (Our results). Moving on to Section 4.1, we will begin by introducing representation balancing methods without decomposed patterns. Specifically, Section 4.1.1 revisits a Wasserstein distance based representation balancing network GNet, and Section 4.1.2 demonstrates how Theorem 2 can be connected with individual propensity confusion, helping us to build a \mathcal{H} -divergence based representation balancing network INet. Subsequently, in Section 4.2, we will introduce how to design a representation balancing method within the scheme of decomposed patterns, based on the PDIG and PPBR methods (Section 4.2.1). The final proposed model DIGNet is presented in Section 4.2.2.

4.1 Representation Balancing without Decomposed Patterns

In representation balancing models, given the input data tuples $(\mathbf{x}, \mathbf{t}, \mathbf{y}) = \{(\mathbf{x}_i, t_i, y_i)\}_{i=1}^N$, the original covariates \mathbf{x} are extracted by some representation function $\Phi(\cdot)$, and representations $\Phi(\mathbf{x})$ are then fed into the outcome functions $h^1(\cdot) := h(\cdot, 1)$ and $h^0(\cdot) := h(\cdot, 0)$ that estimate the potential outcome y^1 and y^0 , respectively. Finally, the factual outcome can be predicted by $h^t(\cdot) = th^1(\cdot) + (1-t)h^0(\cdot)$, and the corresponding outcome loss is

$$\mathcal{L}_y(\mathbf{x}, \mathbf{t}, \mathbf{y}; \Phi, h^t) = \frac{1}{N} \sum_{i=1}^N L(h^t(\Phi(\mathbf{x}_i)), y_i). \quad (11)$$

The loss function \mathcal{L}_y approximates the factual error ϵ_F appeared in Theorems 1 and 2. Minimizing \mathcal{L}_y also corresponds to the Principle I as mentioned in the Introduction.

4.1.1 GNet: Group Distance Minimization Guided Network

The *group distance minimization* focuses on learning representations that minimize the distance between the treated and controlled groups, and the corresponding theoretical foundation is supported by Wasserstein distance guided counterfactual and ITE error bounds (Theorem 1). Previous causal inference methods (e.g., Shalit et al. (2017); Yao et al. (2018); Zhang et al. (2020); Huang et al. (2022)) commonly adopt Wasserstein distance to achieve group distance minimization. Specifically, these methods aim to minimize the empirical approximation of $\mathcal{L}_G(\mathbf{x}, \mathbf{t}; \Phi) = \text{Wass}(\{\Phi(\mathbf{x}_i)\}_{i:t_i=0}, \{\Phi(\mathbf{x}_i)\}_{i:t_i=1})$ to learn balancing patterns. If we denote $\Phi_E(\cdot)$ by the feature extractor that extracts the original covariates \mathbf{x} , then the objective function designed on Theorem 1 is

$$\min_{\Phi_E, h^t} \mathcal{L}_y(\mathbf{x}, \mathbf{t}, \mathbf{y}; \Phi_E, h^t) + \alpha_1 \mathcal{L}_G(\mathbf{x}, \mathbf{t}; \Phi_E). \quad (12)$$

Since the objective is to learn balancing patterns by minimizing the distributional distance between treated and controlled groups, i.e., group distance minimization, we refer to a model with the objective in equation (12) as **GNet**. For the reader's convenience, we illustrate the structure of GNet in Figure 3(a). Note that CFRNet (Shalit et al., 2017) is also the category of GNet.

4.1.2 INet: Individual Propensity Confusion Guided Network

In the field of causal inference, the propensity score plays a central role because it characterizes the probability that one receives treatment (Rosenbaum & Rubin, 1983). For example, the propensity score has been widely employed in prior literature for matching (Caliendo & Kopeinig, 2008) or weighting (Austin & Stuart, 2015) purposes. In this paper, we emphasize that the propensity score also plays an important role in representation balancing, where it serves as a natural indicator of the adequacy of learned balancing patterns. Specifically, we propose the concept of individual propensity confusion, which aims to learn representations that are difficult to utilize for characterizing the propensity of each individual being treated or controlled. **The underlying theoretical foundation is upon the \mathcal{H} -divergence guided ITE error bounds derived in Theorem 2. Specifically, equations 9 and 10 in Theorem 2 highlight the significance of minimizing the generalization bound associated with factual outcome error and the \mathcal{H} -divergence between treated and controlled representations in reducing ITE errors. Subsequently, we will present the details of achieving representation balancing by reducing the \mathcal{H} -divergence between treated and controlled samples in the representation space.**

Let $\mathbf{1}(a)$ be an indicator function that gives 1 if a is true, and \mathcal{H} be the family of binary functions as defined in Theorem 2. To achieve representation balancing, the objective function designed on Theorem 2 should aim to minimize the empirical \mathcal{H} -divergence $\hat{d}_{\mathcal{H}}(p_{\Phi}^{T=1}, p_{\Phi}^{T=0})$ such that

$$\hat{d}_{\mathcal{H}}(p_{\Phi}^{T=1}, p_{\Phi}^{T=0}) = 2 \left(1 - \min_{\eta \in \mathcal{H}} \left[\frac{1}{N} \sum_{i:\eta(\Phi(\mathbf{x}_i))=0} \mathbf{1}[t_i = 1] + \frac{1}{N} \sum_{i:\eta(\Phi(\mathbf{x}_i))=1} \mathbf{1}[t_i = 0] \right] \right). \quad (13)$$

The “min” part in equation (13) indicates that the optimal classifier $\eta^* \in \mathcal{H}$ minimizes the classification error between the estimated treatment $\eta^*(\Phi(\mathbf{x}_i))$ and the observed treatment t_i , i.e., discriminating whether

$\Phi(\mathbf{x}_i)$ is a control ($T = 0$) or treatment ($T = 1$). Equation (13) suggests that $\hat{d}_{\mathcal{H}}(p_{\Phi}^{T=1}, p_{\Phi}^{T=0})$ will be large if η^* can easily distinguish whether $\Phi(\mathbf{x}_i)$ is treated or controlled. In contrast, $\hat{d}_{\mathcal{H}}(p_{\Phi}^{T=1}, p_{\Phi}^{T=0})$ will be small if it is hard for η^* to determine whether $\Phi(\mathbf{x}_i)$ is treated or controlled. Therefore, the prerequisite of a small \mathcal{H} -divergence is to find a map Φ such that any classifier $\eta \in \mathcal{H}$ will get confused about the probability of $\Phi(\mathbf{x}_i)$ being treated or controlled. To achieve this goal, similar to the strategy of empirical approximation of \mathcal{H} -divergence (Ganin et al., 2016), we define a discriminator $\pi(\mathbf{r}) : \mathcal{R} \rightarrow [0, 1]$ that estimates the propensity score of \mathbf{r} , which can be regarded as a surrogate for $\eta(\mathbf{r})$. The classification error for the i^{th} individual can be empirically approximated by the cross-entropy loss between $\pi(\Phi(\mathbf{x}_i))$ and t_i :

$$\mathcal{L}_t(t_i, \pi(\Phi(\mathbf{x}_i))) = -[t_i \log \pi(\Phi(\mathbf{x}_i)) + (1 - t_i) \log(1 - \pi(\Phi(\mathbf{x}_i)))] . \quad (14)$$

As a consequence, we aim to find an optimal discriminator π^* for equation (13) such that π^* maximizes the probability that treatment is correctly classified:

$$\max_{\pi \in \mathcal{H}} \mathcal{L}_I(\mathbf{x}, \mathbf{t}; \Phi, \pi) = \max_{\pi \in \mathcal{H}} \left[-\frac{1}{N} \sum_{i=1}^N \mathcal{L}_t(t_i, \pi(\Phi(\mathbf{x}_i))) \right] . \quad (15)$$

Given the feature extractor $\Phi_E(\cdot)$, the objective of INet can be formulated as a min-max game:

$$\min_{\Phi_E, h^t} \max_{\pi} \mathcal{L}_y(\mathbf{x}, \mathbf{t}, \mathbf{y}; \Phi_E, h^t) + \alpha_2 \mathcal{L}_I(\mathbf{x}, \mathbf{t}; \Phi_E, \pi) . \quad (16)$$

In the maximization, the discriminator π is trained to maximize the probability that treatment is correctly classified. This forces $\pi(\Phi_E(\mathbf{x}_i))$ closer to the true propensity score $e(\mathbf{x}_i)$. In the minimization, the feature extractor Φ_E is trained to fool the discriminator π . This confuses π such that $\pi(\Phi_E(\mathbf{x}_i))$ cannot correctly specify the true propensity score $e(\mathbf{x}_i)$. Eventually, the representations are balanced as the adversarial process makes it difficult for π to determine the propensity of each individual being treated or controlled. We refer to this process as *individual propensity confusion*. Such an adversarial learning technique has been widely used in domain adaptation (e.g., Ganin et al. (2016); Tzeng et al. (2017)) and fair representation learning (e.g., Edwards & Storkey (2015); Madras et al. (2018)) to learn domain-invariant and fair representations. For the reader’s convenience, we illustrate the structure of INet in Figure 3(b).

4.2 Representation Balancing with Decomposed Patterns

4.2.1 The Proposed PDIG and PPBR Methods

PDIG. Although Theorems 1 and 2 provide solid theoretical foundation for GNet (model proposed by Shalit et al. (2017)) and INet (model proposed by us), both of these model types still encounter the inherent trade-off between representation balancing and outcome modeling. To this end, we expect to capture more effective balancing patterns by learning **P**atterns **D**ecomposed with **I**ndividual propensity confusion and **G**roup distance minimization (**PDIG**). More specifically, the covariates \mathbf{x} are extracted by the feature extractor $\Phi_E(\cdot)$, and then $\Phi_E(\mathbf{x})$ are fed into two distinct balancing networks $\Phi_G(\cdot)$ and $\Phi_I(\cdot)$ for group distance minimization and individual propensity confusion, respectively. In summary, PDIG decomposes the balancing patterns into two distinct parts, group distance minimization and individual propensity confusion, which are respectively achieved by the following loss functions:

$$\min_{\Phi_G} \mathcal{L}_G(\mathbf{x}, \mathbf{t}; \Phi_G \circ \Phi_E) \quad (17)$$

$$\min_{\Phi_I} \max_{\pi} \mathcal{L}_I(\mathbf{x}, \mathbf{t}; \Phi_I \circ \Phi_E, \pi) . \quad (18)$$

Here, \circ denotes the composition of two functions, indicating that $\Phi(\cdot)$ in $\mathcal{L}_G(\mathbf{x}, \mathbf{t}; \Phi)$ and $\mathcal{L}_I(\mathbf{x}, \mathbf{t}; \Phi, \pi)$ are replaced by $\Phi_G(\Phi_E(\cdot))$ and $\Phi_I(\Phi_E(\cdot))$, respectively.

PPBR. Motivated by the discussion in Section 1, we design a framework that is capable of capturing **P**atterns of **P**re-balancing and **B**alancing **R**epresentations (**PPBR**) to mitigate potential over-balancing issue mentioned in the Introduction, aiming to preserve information that is useful for outcome predictions.

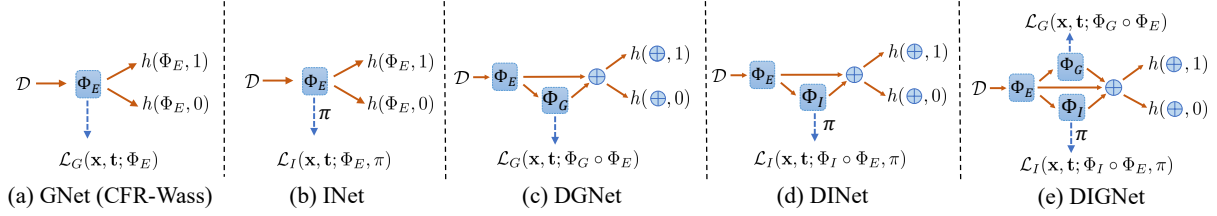


Figure 3: Illustrations of the network architecture of the five models studied in Section 5.

In the PPBR method, the balancing patterns $\Phi_G(\Phi_E(\mathbf{x}))$ and $\Phi_I(\Phi_E(\mathbf{x}))$ are first learned over Φ_G and Φ_I , while Φ_E is remained fixed as pre-balancing patterns. Furthermore, we concatenate the balancing patterns $\Phi_G(\Phi_E(\mathbf{x}))$ and $\Phi_I(\Phi_E(\mathbf{x}))$ with the pre-balancing representations $\Phi_E(\mathbf{x})$ as attributes for outcome prediction. As a result, the proxy features used for outcome predictions are $\Phi_E(\mathbf{x}) \oplus \Phi_G(\Phi_E(\mathbf{x})) \oplus \Phi_I(\Phi_E(\mathbf{x}))$, where \oplus indicates the concatenation by column. For example, if $\mathbf{a} = [1, 2]$ and $\mathbf{b} = [3, 4]$, then $\mathbf{a} \oplus \mathbf{b} = [1, 2, 3, 4]$. Consequently, representation balancing is accomplished over Φ_G and Φ_I , rather than Φ_E . Even if there may be a loss of information relevant to outcome prediction in Φ_G and Φ_I , the pre-balancing patterns Φ_E can still effectively preserve such information. Finally, the objective function with regard to outcome modeling under PPBR method becomes

$$\min_{\Phi_E, \Phi_I, \Phi_G, h^t} \mathcal{L}_y(\mathbf{x}, \mathbf{t}, \mathbf{y}; \Phi_E \oplus (\Phi_I \circ \Phi_E) \oplus (\Phi_G \circ \Phi_E), h^t). \quad (19)$$

4.2.2 The Proposed DIGNet

Combining with PDIG and PPBR, we propose a new neural **Network** model that incorporates **D**ecomposed patterns with **I**ndividual propensity confusion and **G**roup distance minimization, which we call **DIGNet**. The objective of DIGNet is separated into four stages:

$$\min_{\Phi_G} \alpha_1 \mathcal{L}_G(\mathbf{x}, \mathbf{t}; \Phi_G \circ \Phi_E), \quad (20)$$

$$\max_{\pi} \alpha_2 \mathcal{L}_I(\mathbf{x}, \mathbf{t}; \Phi_I \circ \Phi_E, \pi), \quad (21)$$

$$\min_{\Phi_I} \alpha_2 \mathcal{L}_I(\mathbf{x}, \mathbf{t}; \Phi_I \circ \Phi_E, \pi), \quad (22)$$

$$\min_{\Phi_E, \Phi_I, \Phi_G, h^t} \mathcal{L}_y(\mathbf{x}, \mathbf{t}, \mathbf{y}; \Phi_E \oplus (\Phi_I \circ \Phi_E) \oplus (\Phi_G \circ \Phi_E), h^t). \quad (23)$$

Within each iteration, DIGNet minimizes the group distance through equation 20, and plays an adversarial game to achieve propensity confusion through equation 21 and equation 22. In equation 23, DIGNet updates both the pre-balancing patterns Φ_E and balancing patterns Φ_I, Φ_G , along with the outcome function h^t to minimize the outcome prediction loss. For the reader's convenience, we illustrate the structure of DIGNet in Figure 3(e).

4.3 Insights of Representation Balancing with Decomposed Patterns

Our proposed DIGNet model builds upon the PDIG and PPBR methods. The PPBR method is relatively straightforward, as it forms more flexible predictor $(\Phi_E \oplus (\Phi_I \circ \Phi_E))$ (or $(\Phi_E \oplus (\Phi_G \circ \Phi_E))$) compared to the solely predictor Φ_E . Therefore, incorporating both pre-balancing and balancing patterns is helpful in enhancing the model's complexity and its ability to capture more useful information for outcome prediction. However, there still remains further exploration to better understand why the PDIG method is effective. The DIGNet model aims to learn balancing patterns based on both Wasserstein distance and \mathcal{H} -divergence. At first glance, one might assume that incorporating both distances could be redundant, as one distance seems naturally to imply the other. In this section, we gain some insights of these two divergence metrics. First, we provide a systematic discussion on the properties of Wasserstein distance and \mathcal{H} -divergence. In

addition, we utilize a toy example to illustrate their distinct abilities in capturing distributional disparity. Further, we use this example to aid readers in better understanding the trade-off problem encountered in representation balancing models (Figure 5). Finally, we establish a connection between our method and [the Elastic Net method and Multi-task learning approach](#), which offers valuable insights and explanations regarding the intuition behind involving both metrics as regularizations.

4.3.1 Properties of Wasserstein Distance and \mathcal{H} -Divergence

Wasserstein distance and \mathcal{H} -Divergence possess distinct theoretical properties. The effectiveness of the Wasserstein distance in measuring distributional differences for classification tasks in domain adaptation has been demonstrated in Shen et al. (2018). Furthermore, Shalit et al. (2017) highlights the potential of Wasserstein distance in representation balancing models for ITE estimation, which significantly outperforms traditional ITE estimation methods. Wasserstein distance is also widely adopted in other research domains, such as fair representation learning (Jiang et al., 2020), as discussed in Section 1.1. Its prevalence stems from its strong capability to capture better diversities compared to \mathcal{H} -Divergence (Shui et al., 2020). Studies have proven that under certain conditions, it is possible to bound \mathcal{H} -Divergence using Wasserstein distance (Villani et al., 2009; Shui et al., 2020), which provides a reasonable explanation for the overall superiority of the Wasserstein distance in learning domain-invariant features (Zhiri et al., 2022). However, it is important to note that this bound does not hold in general (Chae & Walker, 2020), suggesting that a smaller \mathcal{H} -divergence does not necessarily imply a smaller Wasserstein distance. To better illustrate the difference between these two measures, we provide a concrete example below.

Toy example. Consider the following three probability density functions $p_1(x)$, $p_2(x)$, and $p_3(x)$ defined over $x \in [0, 1]$:

$$p_1(x) = \begin{cases} 2.5, & \text{if } 0 \leq x < 0.25 \\ 0.5, & \text{if } 0.25 \leq x < 0.5 \\ 0.5, & \text{if } 0.5 \leq x < 0.75 \\ 0.5, & \text{if } 0.75 \leq x \leq 1 \end{cases} \quad p_2(x) = \begin{cases} 0.5, & \text{if } 0 \leq x < 0.25 \\ 2.5, & \text{if } 0.25 \leq x < 0.5 \\ 0.5, & \text{if } 0.5 \leq x < 0.75 \\ 0.5, & \text{if } 0.75 \leq x \leq 1 \end{cases} \quad p_3(x) = \begin{cases} 0.5, & \text{if } 0 \leq x < 0.25 \\ 0.5, & \text{if } 0.25 \leq x < 0.5 \\ 2.5, & \text{if } 0.5 \leq x < 0.75 \\ 0.5, & \text{if } 0.75 \leq x \leq 1 \end{cases}.$$

The above three distributions are depicted in Figure 4. Further, we set the classifier in \mathcal{H} -divergence as $\eta(x) = \mathbf{1}\{x \geq p\}$, and set the order in Wasserstein distance as $p = 1$. By utilizing the definitions of \mathcal{H} -divergence and 1-Wasserstein distance, one can make a direct comparison between the discrepancy in (p_1, p_2) and the discrepancy in (p_1, p_3) :

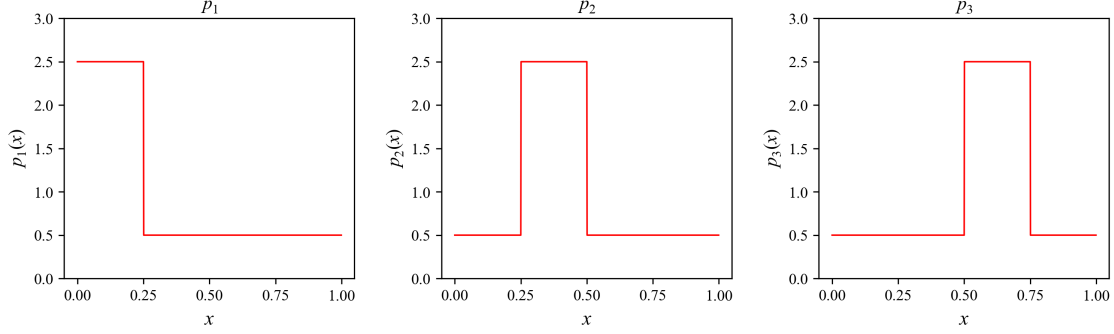
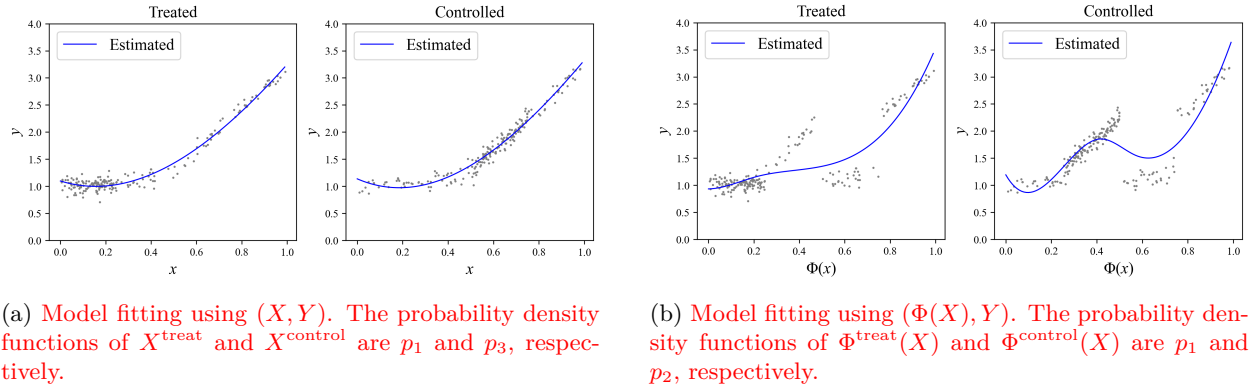
$$\begin{aligned} d_{\mathcal{H}}(p_1, p_2) &= d_{\mathcal{H}}(p_1, p_3); \\ \text{Wass}(p_1, p_2) &< \text{Wass}(p_1, p_3). \end{aligned} \tag{24}$$

Equation 24 confirms that Wasserstein distance is able to capture more diverse distributional disparities compared to \mathcal{H} -divergence. However, in the subsequent content, we will demonstrate that such an advantage might be a limitation in causal representation learning due to the trade-off problem.

Understanding the trade-off. The above example serves as simple evidence that supports the conclusion that Wasserstein distance can capture better diversities between distributions compared to \mathcal{H} -divergence (Shui et al., 2020). However, as discussed in Section 1, it is important to note that achieving a more balanced distribution does not necessarily ensure favorable generalization to counterfactuals. This is because the pursuit of balanced representations may inadvertently lead to a loss of information useful for factual outcome estimates. We will now use the above example to gain further understanding on this matter.

Consider a simple data-generating process where X^{treat} , the covariate in the treated group, follows the distribution $p_1(x)$, and X^{control} , the covariate in the controlled group, follows the distribution $p_3(x)$. Let the potential outcomes are generated by $Y^1 = \tau_1(X) + \epsilon_1$ and $Y^0 = \tau_0(X) + \epsilon_0$, where $\epsilon_1 \sim \mathcal{N}(0, 0.1)$ and $\epsilon_0 \sim \mathcal{N}(0, 0.1)$. Let the true potential outcome functions $\tau^1(x)$ and $\tau^0(x)$ be as follows:

$$\tau^1(x) = \tau^0(x) = (x^2 + 1)\mathbf{1}\{0 \leq x < 0.5\} + (4x - 0.75)\mathbf{1}\{0.5 \leq x \leq 1\}. \tag{25}$$

Figure 4: Distributions of $p_1(x)$, $p_2(x)$, and $p_3(x)$ in the example of Section 4.3.1.(a) Model fitting using (X, Y) . The probability density functions of X^{treat} and X^{control} are p_1 and p_3 , respectively.(b) Model fitting using $(\Phi(X), Y)$. The probability density functions of $\Phi^{\text{treat}}(X)$ and $\Phi^{\text{control}}(X)$ are p_1 and p_2 , respectively.Figure 5: Model fitting using (X, Y) and $(\Phi(X), Y)$ based on the example in Section 4.3.1.

In addition, consider a representation function $\Phi(x)$ such that

$$\Phi(x) = x\mathbf{1}\{0 \leq x < 0.25\} + (x + 0.25)\mathbf{1}\{0.25 \leq x < 0.5\} + (x - 0.25)\mathbf{1}\{0.5 \leq x < 0.75\} + x\mathbf{1}\{0.75 \leq x \leq 1\}. \quad (26)$$

We can find Φ achieves representation balancing under Wasserstein distance measure, but does not under \mathcal{H} -divergence measure. In original data, x^{treat} follows p_1 and x^{control} follows p_3 . After mapping x to $\Phi(x)$, $\Phi^{\text{treat}}(x)$ follows p_1 and $\Phi^{\text{control}}(x)$ follows p_2 . Consequently, based on the results in equation 24, we have

$$\begin{aligned} d_{\mathcal{H}}(p_{\Phi}^{\text{treat}}, p_{\Phi}^{\text{control}}) &= d_{\mathcal{H}}(p_X^{\text{treat}}, p_X^{\text{control}}); \\ \text{Wass}(p_{\Phi}^{\text{treat}}, p_{\Phi}^{\text{control}}) &< \text{Wass}(p_X^{\text{treat}}, p_X^{\text{control}}). \end{aligned} \quad (27)$$

We now investigate the fitting performance of models using (x, y) and $(\Phi(x), y)$ to check whether there is a loss of outcome-related information during representation balancing. In Figure 5a and Figure 5b, we present scatter plots of samples from (x, y) and $(\Phi(x), y)$ respectively, depicted as gray points. Following the approach of Kennedy (2023), we employ smoothing spline functions to fit these samples, and the estimated functions are illustrated in blue.

In Figure 5a, we observe that both τ_1 and τ_0 are well fitted using (x, y) , with their estimates being very close to each other. This is consistent with the setup of $\tau_1 = \tau_0$. In contrast, Figure 5b reveals that the fittings of τ_1 and τ_0 are inadequate using $(\Phi(x), y)$, resulting in substantially different estimates. The result of different estimates violates the setup of $\tau_1 = \tau_0$. In this case, a model based on Wasserstein distance would retain Φ due to its achievement of representation balancing. Unfortunately, Φ suffers from a loss of valuable information that is crucial for outcome prediction. In contrast, a model based on \mathcal{H} -divergence would not keep Φ since it does not contribute to reducing the domain distance compared to the original data. Fortunately, the original data preserve the information necessary for outcome modeling. Therefore,

this example not only emphasizes the significance of incorporating both metrics but also highlights the importance of considering both pre-balancing patterns and balancing patterns.

4.3.2 Connection with other machine learning methods

In the previous sections, we have discussed the trade-off between factual outcome prediction and representation balancing in classic representation learning models. As part of our proposed improvements, DIGNet involves learning two distinct representations using Wasserstein distance and \mathcal{H} -divergence separately and concatenates the learned representations for outcome modeling. In this section, we will explore more detailed connections between our design and other machine learning methods.

Connection with Elastic Net: balancing on two discrepancies. Our DIGNet model involves two discrepancy metrics: Wasserstein distance and \mathcal{H} -divergence. We will now provide additional explanations on its connection with the Elastic Net method. In supervised learning, a regularization term is often incorporated during model training to mitigate the bias-variance trade-off. In the case of linear regression, Lasso (Tibshirani, 1996) and Ridge (Hoerl & Kennard, 1970) are proposed to improve the Ordinary Least Squares (OLS) method, with Lasso involving l_1 regularization while Ridge involving l_2 regularization:

$$\begin{aligned} \text{Lasso: } \min_{\beta \in \mathcal{R}^d} \frac{1}{N} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \alpha \|\beta\|_1 &= \min_{\beta \in \mathcal{R}^d} \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \beta - y_i)^2 + \alpha \sum_{j=1}^d |\beta_j|. \\ \text{Ridge: } \min_{\beta \in \mathcal{R}^d} \frac{1}{N} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \alpha \|\beta\|_2^2 &= \min_{\beta \in \mathcal{R}^d} \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \beta - y_i)^2 + \alpha \sum_{j=1}^d \beta_j^2. \end{aligned}$$

The different properties between l_1 regularization and l_2 regularization lead to distinct advantages and disadvantages between Lasso method and Ridge method. Given their differences, a method of Elastic Net (Zou & Hastie, 2005) is proposed by combining both l_1 regularization and l_2 regularization:

$$\text{Elastic Net: } \min_{\beta \in \mathcal{R}^d} \frac{1}{N} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \alpha_1 \|\beta\|_1 + \alpha_2 \|\beta\|_2^2 = \min_{\beta \in \mathcal{R}^d} \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \beta - y_i)^2 + \alpha_1 \sum_{j=1}^d |\beta_j| + \alpha_2 \sum_{j=1}^d \beta_j^2.$$

The Elastic Net method integrates the strengths of two distinct approaches: the l_1 regularization term enforces sparsity, while the l_2 regularization maintains the grouping effect (Zhou, 2013; Narisetty, 2020). The Elastic Net has also motivated some research studies to adopt the idea of combining l_1 and l_2 regularizations in of deep neural networks (DNNs) (Kang et al., 2017; Chen et al., 2018; Hu et al., 2023; Xu et al., 2023a). Notably, a recent study (Xu et al., 2023a) presents an excess risk bound for Elastic Net Regularized DNNs. This finding provides supporting evidence that incorporating both l_1 and l_2 regularizations in a DNN model is reasonable. The insights gained from (Xu et al., 2023a) shed light on the theoretical explanation of our method, and even pave the way for exploring the integration of different divergence metrics in other research areas, such as domain adaptation, transfer learning, and fair representation learning.

Connection with multi-task learning: balancing on two representations. Our DIGNet model performs representation balancing on two distinct representations using Wasserstein distance and \mathcal{H} -divergence separately, and the learned representations are then concatenated for outcome modeling. We will now provide additional explanations regarding its connection with the multi-task learning method. In multi-task learning, distinct representations are learned for different tasks, with each task involving its own objective function. An important step in multi-task learning is integrating the information from these separately learned representations into a unified representation. One common approach is to concatenate the task-specific representations to form a joint representation, which effectively preserves the information from each task for outcome modeling. Li et al. (2018); Baltrušaitis et al. (2018); Crawshaw (2020); Yan et al. (2021); Xu et al. (2023b). For example, in an E-commerce application Liu et al. (2023), diverse types of user footprints are encoded using different representations with diverse objectives. The learned representations are then concatenated to make the final target prediction. Similarly, in another application Wu et al. (2018), user and product attentions are separately learned on two distinct representations, which are later concatenated

for the final outcome prediction. In tasks such as image and text classification Hao et al. (2023), concatenating multiple representations has been shown useful in enhancing the classification performance by leveraging the complementary information provided by each representation. Furthermore, a recent study on multi-view learning (Li et al., 2024) has also demonstrated that concatenating both the non-attention and attention representations of each view can prevent information loss in the final classification task.

Summary of strengths and limitations. Our method combines Wasserstein distance and \mathcal{H} -divergence for representation balancing to capture different types of balancing patterns compared to classic representation balancing models. This shares a similar intuition with the Elastic Net, which combines l_1 and l_2 regularizations to learn features with different properties. Notably, the two regularizations in Elastic Net are learned on a single parameter space with one objective, this provides more interpretability but might introduce a new trade-off. Different from Elastic Net, our DIGNet model concatenates the two distinct representations that are learned from two different tasks: Wasserstein distance guided and \mathcal{H} -divergence guided representation balancing. This aligns with the principle of multi-task learning. The concatenation fusion technique is extensively employed in numerous multi-task learning studies (Baltrušaitis et al., 2018), as it effectively preserves and integrates information from different tasks, leading to improved performance in the final prediction task (Hao et al., 2023; Li et al., 2024). However, it is crucial to acknowledge that this straightforward concatenation approach can present challenges when interpreting the specific role of each representation and can also increase model complexity (Jia et al., 2020).

5 Experiments

In non-randomized observational data, the ground truth regarding treatment effects remains inaccessible due to the absence of counterfactual information. Therefore, we use simulated data and semi-synthetic benchmark data to test the performance of our methods and other baseline models. In this section, we primarily investigate the three following questions:

Q1. Is PDIG helpful in ITE estimation through Path I in the Introduction, i.e., learning more effective balancing patterns without affecting factual outcome prediction?

Q2. Is PPBR helpful in ITE estimation through Path II in the Introduction, i.e., improving factual outcome prediction without affecting learning balancing patterns?

Q3. Can the proposed DIGNet model outperform other baseline models on benchmark dataset?

Ablation models. To investigate Q1 and Q2, we conducted ablation studies and designed two ablation models, **DGNet** and **DINet**, where DGNet (or DINet) can be considered as DIGNet without PDIG, and GNet (or INet) can be considered as DGNet (or DINet) without PPBR. The structures of DGNet and DINet are shown in Figure 3(c) and Figure 3(d), and the objectives of DGNet and DINet are deferred to Section A.5.

5.1 Experimental Settings

Simulation data. Previous causal inference works assess the model effectiveness by varying the distribution imbalance of covariates in treated and controlled groups at different levels (Yao et al., 2018; Yoon et al., 2018; Du et al., 2021). As suggested by Assaad et al. (2021), we draw 1000 observational data points from the following data generating strategy:

$$\begin{aligned} \mathbf{X}_i &\sim \mathcal{N}(\mathbf{0}, \sigma^2 \cdot [\rho \mathbf{1}_p \mathbf{1}_p' + (1 - \rho) \mathbf{I}_p]), \\ T_i | \mathbf{X}_i &\sim \text{Bernoulli}(1/(1 + \exp(-\gamma \mathbf{X}_i))), \\ Y_i^0 &= \beta_0' \mathbf{X}_i + \xi_i, \quad Y_i^1 = \beta_1' \mathbf{X}_i + \xi_i, \quad \xi_i \sim \mathcal{N}(0, 1). \end{aligned}$$

Here, $\mathbf{1}_p$ denotes the p -dimensional all-ones vector and \mathbf{I}_p denotes the identity matrix of size p . We fix $p = 10, \rho = 0.3, \sigma^2 = 2, \beta_0' = [0.3, \dots, 0.3], \beta_1' = [1.3, \dots, 1.3]$ and vary $\gamma \in \{0.25, 0.5, 0.75, 1, 1.5, 2, 3\}$ to yield different levels of selection bias. As seen in Figure 6, selection bias becomes more severe with γ increasing.

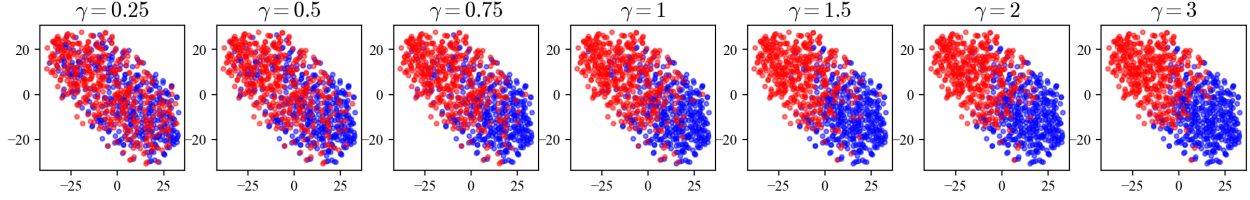


Figure 6: T-SNE visualizations of the covariates as γ varies. Red represents the treatment group and blue represents the control group. A larger γ indicates a greater imbalance between the two groups.

For each γ , we repeat the above data generating process to generate 30 different datasets, with each dataset split by the ratio of 56%/24%/20% as training/validation/test sets.

Semi-synthetic data. The IHDP dataset, introduced by Hill (2011), originates from the Infant Health and Development Program (IHDP). This program conducted a randomized controlled experiment in 1985 to investigate whether there is a positive causal effect of frequent high-quality child care and home visits (treatment) on cognitive scores (outcome). The collected data comprise 25-dimensional pre-treatment covariates, including measurements on the infants (e.g., birth weight, gender, head circumference), as well as measurements on the mothers during pregnancy (e.g., age, marital status, education, smoking and drinking habits). In order to create an observational dataset that involves selection bias, Hill excluded a subpopulation (children with nonwhite mothers) from the treated group. Consequently, the IHDP dataset exhibits a covariate shift, resulting in imbalanced treated and controlled groups. The final IHDP dataset consists of 747 samples, comprising 139 treated samples and 608 controlled samples. The potential outcomes were generated using setting A in the NPCI package Dorie (2021). We use the same 1000 datasets as used in Shalit et al. (2017), with each dataset split by the ratio of 63%/27%/10% as training/validation/test sets.

Models and metrics. In simulation experiments, we perform comprehensive comparisons between INet, GNet, DINet, DGNet, and DIGNet in terms of the mean and standard error for the following metrics: $\sqrt{\epsilon_{PEHE}}$, $\sqrt{\epsilon_{CF}}$, and $\sqrt{\epsilon_F}$ with L defined in Definition 1 being the squared loss, as well as the empirical approximations of $Wass(p_{\Phi}^{T=1}, p_{\Phi}^{T=0})$ and $d_{\mathcal{H}}(p_{\Phi}^{T=1}, p_{\Phi}^{T=0})$ (denoted by $Wass$ and $\hat{d}_{\mathcal{H}}$, respectively). Note that as shown in Figure 3, $Wass$ is over Φ_E for GNet while over Φ_G for DGNet and DIGNet; $\hat{d}_{\mathcal{H}}$ is over Φ_E for INet while over Φ_I for DINet and DIGNet. To analyze the source of gain and ensure fair comparison in simulation studies, we fix hyperparameters across all models. This way is consistent with Curth & van der Schaar (2021). We apply an early stopping rule to all models as Shalit et al. (2017) do. In IHDP experiment, we use $\sqrt{\epsilon_{PEHE}}$, as well as an additional metric $\epsilon_{ATE} = |\hat{\tau}_{ATE} - \tau_{ATE}|$ to evaluate performances of various causal models (see them in Table 6). More descriptions of the implementation details, as well as the analysis of training time, training stability, and hyperparameter sensitivity, are deferred to Section A.4.

Device. All the experiments are run on Dell 7920 with one 16-core Intel Xeon Gold 6250 3.90GHz CPU and three NVIDIA Quadro RTX 6000 GPUs.

5.2 Results and Analysis

5.2.1 Preliminary Experimental Results

In this part, we first make a general comparison between different models with the degree of covariate imbalance increasing, and the relevant results are shown in Figure 7. There are four main observations:

1. DIGNet attains the lowest $\sqrt{\epsilon_{PEHE}}$ across all datasets, while GNet have inferior performances than other models;
2. DINet and DGNet outperform INet and GNet regarding $\sqrt{\epsilon_{CF}}$ and $\sqrt{\epsilon_{PEHE}}$;

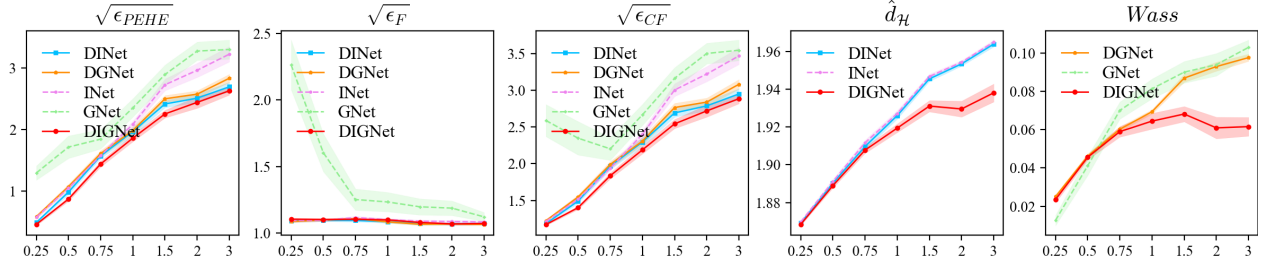


Figure 7: Plots of model performances on test set for different metrics as γ varies in $\{0.25, 0.5, 0.75, 1, 1.5, 2, 3\}$. Each graph shows the average of 30 runs with standard errors shaded. Lower lines indicate lower values of the metric.

3. INet, DNet, and DGNet have comparable performance to DIGNet in terms of factual outcome estimations ($\sqrt{\epsilon_F}$), but cannot compete with DIGNet in terms of counterfactual estimations ($\sqrt{\epsilon_{CF}}$) or ITE estimations ($\sqrt{\epsilon_{PEHE}}$);
4. DIGNet achieves smaller \hat{d}_H (or $Wass$) than DNet and INet (or DGNet and GNet), especially when the covariate shift problem is severe (e.g., when $\gamma > 1$).

In conclusion, the above study has produced several noteworthy findings. Firstly, finding (1) reveals that our proposed DIGNet model consistently performs well in ITE estimation. Secondly, as indicated by finding (2), implementing the PPBR approach can enhance the predictive accuracy of factual and counterfactual outcomes. Lastly, findings (3) and (4) highlight the role of PDIG structure in enhancing the simultaneous reinforcement and complementarity of group distance minimization and individual propensity confusion, resulting in more balanced representations. Our subsequent analysis will step beyond these preliminary conclusions to gain a deeper understanding of the effectiveness of the proposed methods.

5.2.2 Further Ablation Studies

So far our preliminary observations have show that the relationship between the ITE errors of each model is: $DIGNet < DNet < INet$ and $DIGNet < DGNet < GNet$. To further explore how PDIG and PPBR contribute to the improvement of ITE estimations, we choose the case with high selection bias ($\gamma = 3$) to analyze the source of gain for PDIG and PPBR. We report model performances (mean \pm std) of each specific metric averaged across 30 runs on test set in Table 1 and Table 2. We also report model performances (mean \pm std) averaged over 30 training and test sets in Table 3. Below we discuss the source of gain in detail.

Ablation study for PDIG. *The PDIG structure is manifest to be effective in capturing more effective balancing patterns, without affecting factual outcome predictions.* As depicted in Figure 7, DIGNet exhibits more balanced representations, irrespective of whether the discrepancy is measured by \hat{d}_H or $Wass$, while DIGNet, DNet, and DGNet demonstrate comparable estimates of factual outcomes ($\sqrt{\epsilon_F}$). Two additional pieces of specific evidence can be observed from Table 1: (1) Despite the absence of PDIG in DNet and DGNet when compared to DIGNet, these three models exhibit very similar performance regarding $\sqrt{\epsilon_F}$, with the performance being 1.07 ± 0.01 . This indicates that PDIG does not impact the factual estimation. (2) DIGNet achieves smaller \hat{d}_H with a $|1.94/1.96 - 1| = 1.0\%$ reduction (or $Wass$ with a $|0.06/0.10 - 1| = 40\%$ reduction) compared with DNet (or DGNet). This indicates that PDIG enables the model to learn more effective balancing patterns. The above two points indicate that PDIG can capture more effective balancing patterns, without affecting factual outcome predictions. This advantage translates into superior counterfactual estimation, with DIGNet reducing $\sqrt{\epsilon_{CF}}$ by $|2.89/2.95 - 1| = 2.0\%$ and $|2.89/3.08 - 1| = 6.2\%$ compared to DNet and DGNet, respectively. Correspondingly, DIGNet also shows superiority in treatment effect estimation ($\sqrt{\epsilon_{PEHE}}$ and ϵ_{ATE}) compared to DNet (or DGNet), as demonstrated in Table 3.

Table 1: Ablation study for PDIG: Mean \pm std of each metric averaged across 30 runs on test set when $\gamma = 3$. Lower value is better.

	$\sqrt{\epsilon_F}$	$\sqrt{\epsilon_{CF}}$	$\hat{d}_{\mathcal{H}}$	W_{ass}
DIGNet	1.07 ± 0.01	2.89 ± 0.07	1.94 ± 0.00	0.06 ± 0.00
DINet	1.07 ± 0.01	2.95 ± 0.07	1.94 ± 0.00	-
DGNet	1.07 ± 0.01	3.08 ± 0.07	-	0.10 ± 0.00

Table 2: Ablation study for PPBR: Mean \pm std of each metric averaged across 30 runs on test set when $\gamma = 3$. Lower value is better.

	$\sqrt{\epsilon_F}$	$\sqrt{\epsilon_{CF}}$	$\hat{d}_{\mathcal{H}}$	W_{ass}
DGNet	1.07 ± 0.01	3.08 ± 0.07	-	0.10 ± 0.00
GNet	1.12 ± 0.03	3.55 ± 0.14	-	0.10 ± 0.00
DINet	1.07 ± 0.01	2.95 ± 0.07	1.96 ± 0.00	-
INet	1.08 ± 0.01	3.47 ± 0.12	1.96 ± 0.00	-

Table 3: Training- & test- set $\sqrt{\epsilon_{PEHE}}$ & ϵ_{ATE} when $\gamma = 3$. Mean \pm standard error of 30 runs.

	Training set		Test set	
	$\sqrt{\epsilon_{PEHE}}$	ϵ_{ATE}	$\sqrt{\epsilon_{PEHE}}$	ϵ_{ATE}
GNet	3.30 ± 0.15	2.58 ± 0.14	3.30 ± 0.16	2.59 ± 0.14
INet	3.24 ± 0.11	2.46 ± 0.09	3.22 ± 0.12	2.47 ± 0.10
DGNet	2.86 ± 0.06	2.15 ± 0.03	2.83 ± 0.07	2.15 ± 0.04
DINet	2.70 ± 0.06	2.12 ± 0.04	2.69 ± 0.08	2.13 ± 0.05
DIGNet	2.66 ± 0.07	2.04 ± 0.05	2.63 ± 0.07	2.03 ± 0.04

Table 4: Training- & test- set $\sqrt{\epsilon_{PEHE}}$ & ϵ_{ATE} on IHDP. Mean \pm standard error of 100 runs.

	Training set		Test set	
	$\sqrt{\epsilon_{PEHE}}$	ϵ_{ATE}	$\sqrt{\epsilon_{PEHE}}$	ϵ_{ATE}
GNet	0.71 ± 0.15	0.12 ± 0.01	0.77 ± 0.18	0.15 ± 0.02
INet	0.66 ± 0.09	0.13 ± 0.01	0.72 ± 0.11	0.15 ± 0.02
DGNet	0.53 ± 0.07	0.11 ± 0.01	0.60 ± 0.09	0.13 ± 0.01
DINet	0.57 ± 0.12	0.13 ± 0.01	0.60 ± 0.11	0.14 ± 0.01
DIGNet	0.42 ± 0.02	0.11 ± 0.01	0.45 ± 0.04	0.12 ± 0.01

Ablation study for PPBR. *The PPBR approach contributes to enhancing factual outcome predictions, without affecting learning balancing patterns.* From Table 2, we gain two important insights: (1) The difference in learned balancing patterns, measured by $\hat{d}_{\mathcal{H}}$ (or W_{ass}), between DINet and INet (or DGNet and GNet), is negligible. This implies that PPBR does not affect learning balancing patterns. (2) Compared with INet, DINet achieves smaller $\sqrt{\epsilon_F}$, with $|1.07/1.08 - 1| = 0.9\%$ error reduction. Similarly, compared with GNet, DGNet achieves smaller $\sqrt{\epsilon_F}$, with $|1.07/1.12 - 1| = 4.5\%$ error reduction. These two observations reveal that PPBR can improve factual outcome predictions, without affecting learning balancing patterns. Benefiting from the advantage of PPBR, the improvement is particularly pronounced in counterfactual estimation. Comparing DINet with INet, the reduction in $\sqrt{\epsilon_{CF}}$ amounts to $|2.95/3.47 - 1| = 15.0\%$. Similarly, comparing DGNet with GNet, the reduction is $|3.08/3.55 - 1| = 13.2\%$. Correspondingly, DINet (or DGNet) attains smaller treatment effect errors ($\sqrt{\epsilon_{PEHE}}$ and ϵ_{ATE}) compared with INet (or GNet), as shown in Table 3.

Significance analysis for the improvements. To assess the significance of the improvements observed in the above ablation studies, we conducted an additional significance analysis by recording the values of $\sqrt{\epsilon_{PEHE}}$ and ϵ_{ATE} for 30 runs of each of the 5 models (GNet, INet, DGNet, DINet, and DIGNet). Subsequently, we performed a t-test for GNet vs. DGNet, INet vs. DINet, DGNet vs. DIGNet, and DINet vs. DIGNet, to investigate the statistical significance of their differences. The relevant results are reported in Table 5. The results reveal a statistically significant difference between GNet and DGNet, INet and DINet, as well as DGNet and DIGNet. Note that the difference between DINet and DIGNet is not statistically significant, despite DIGNet exhibiting smaller treatment effect estimation errors on average compared to DINet.

5.2.3 Comparisons on IHDP benchmark.

In this part, we perform experiments on the IHDP benchmark dataset to compare the performances of different models. The corresponding results are reported in Table 4 and 6.

First, we report the ablation results on 1-100 IHDP datasets in Table 4, aiming to examine the consistent effectiveness of PDIG and PPBR. Specifically, Table 4 shows that DINet and DGNet are superior to INet and GNet but inferior to DIGNet concerning treatment effect estimation, suggesting that both PDIG and PPBR are advantageous for treatment effect estimation. For example, on the test set, DINet reduces $\sqrt{\epsilon_{PEHE}}$ by $|0.60/0.72 - 1| = 16.7\%$ for INet, and DIGNet reduces $\sqrt{\epsilon_{PEHE}}$ by $|0.45/0.60 - 1| = 25\%$ for DINet. This is consistent with the findings before: PDIG and PPBR are beneficial to treatment effect estimation.

Table 5: Significance analysis regarding the achieved improvements by comparing GNet and DGNet, INet and DINet, DGNet and DIGNet, DINet and DIGNet. The p-value ≤ 0.05 indicates difference is statistically significant.

	Training set				Test set			
	$\sqrt{\epsilon_{PEHE}}$		ϵ_{ATE}		$\sqrt{\epsilon_{PEHE}}$		ϵ_{ATE}	
	t-value	p-value	t-value	p-value	t-value	p-value	t-value	p-value
GNet vs. DGNet	2.7435	0.0081	2.9844	0.0042	2.7073	0.0089	2.9269	0.0049
INet vs. DINet	4.0812	0.0001	3.5222	0.0008	3.5665	0.0007	3.0824	0.0031
DGNet vs. DIGNet	2.0240	0.0476	1.8888	0.0639	2.0650	0.0434	2.0935	0.0407
DINet vs. DIGNet	0.4513	0.6535	1.3525	0.1815	0.6079	0.5456	1.5473	0.1272

Furthermore, we undergo comparisons between DIGNet and other causal models on 1-1000 IHDP datasets and report the results in Table 6. The results highlight the superior performance of the proposed DIGNet across all the models. Specifically, in comparison to the second-best method in test-sample performance, DIGNet achieves a substantial improvement, with error reduced by $|0.45/0.57 - 1| = 21\%$ in terms of $\sqrt{\epsilon_{PEHE}}$ and $|0.12/0.13 - 1| = 7.7\%$ in terms of ϵ_{ATE} . Moreover, it is worth noting that DIGNet consistently achieves the lowest errors across various datasets and metrics, revealing its robust performance. We also conduct an additional experiments on another benchmark dataset Twins. The details and results are deferred to Section A.4

Table 6: Training- & test- set $\sqrt{\epsilon_{PEHE}}$ & ϵ_{ATE} on IHDP. Mean \pm standard error of 1000 runs.

	Training set		Test set	
	$\sqrt{\epsilon_{PEHE}}$	ϵ_{ATE}	$\sqrt{\epsilon_{PEHE}}$	ϵ_{ATE}
OLS/LR ₁ (Johansson et al., 2016)	5.8 \pm .3	.73 \pm .04	5.8 \pm .3	.94 \pm .06
OLS/LR ₂ (Johansson et al., 2016)	2.4 \pm .1	.14 \pm .01	2.5 \pm .1	.31 \pm .02
k-NN (Crump et al., 2008)	2.1 \pm .1	.14 \pm .01	4.1 \pm .2	.79 \pm .05
BART (Chipman et al., 2010)	2.1 \pm .1	.23 \pm .01	2.3 \pm .1	.34 \pm .02
CF (Wager & Athey, 2018)	3.8 \pm .2	.18 \pm .01	3.8 \pm .2	.40 \pm .03
CEVAE (Louizos et al., 2017)	2.7 \pm .1	.34 \pm .01	2.6 \pm .1	.46 \pm .02
SITE (Yao et al., 2018)	.69 \pm .0	.22 \pm .01	.75 \pm .0	.24 \pm .01
GANITE (Yoon et al., 2018)	1.9 \pm .4	.43 \pm .05	2.4 \pm .4	.49 \pm .05
BLR (Johansson et al., 2016)	5.8 \pm .3	.72 \pm .04	5.8 \pm .3	.93 \pm .05
BNN (Johansson et al., 2016)	2.2 \pm .1	.37 \pm .03	2.1 \pm .1	.42 \pm .03
TARNet (Shalit et al., 2017)	.88 \pm .0	.26 \pm .01	.95 \pm .0	.28 \pm .01
CFR-Wass (GNet) (Shalit et al., 2017)	.73 \pm .0	.12 \pm .01	.81 \pm .0	.15 \pm .01
Dragonnet (Shi et al., 2019)	1.3 \pm .4	.14 \pm .01	1.3 \pm .5	.20 \pm .05
MBRL (Huang et al., 2022)	.52 \pm .0	.12 \pm .01	.57 \pm .0	.13 \pm .01
DIGNet (Ours)	.41 \pm .0	.11 \pm .01	.46 \pm .0	.12 \pm .01

6 Conclusion

This paper establishes a theoretical foundation by deriving counterfactual and ITE error bounds based on \mathcal{H} -divergence. This theoretical foundation builds a connection between representation balancing and individual propensity confusion. Furthermore, based on individual propensity confusion and group distance minimization, we suggest learning decomposed patterns for representation balancing models using the PDIG and PPBR methods. Further, building upon PDIG and PPBR, we propose a novel model DIGNet, for treatment effect estimation. Comprehensive experiments verify that PDIG and PPBR follow different pathways to improve counterfactual and ITE estimation. In particular, PDIG enables the model to capture more effective balancing patterns without affecting factual outcome prediction, while PPBR contributes to improving factual outcome predictions without influencing learning balancing patterns. We hope these findings can constitute an important step to inspire more research concerning the generalization of representation balancing models for counterfactual and ITE estimation.

Limitations and future works. Our paper verifies the effectiveness of PDIG and PPBR in improving ITE estimation, it is also important to step beyond our empirical insights into future theoretical studies aimed at addressing the trade-off challenge mentioned in the introduction, e.g., exploring the possibility of deriving tighter theoretical error bounds based on learning decomposed patterns, and involving the orthogonal machine learning (Chernozhukov et al., 2018; Oprescu et al., 2019; Nie & Wager, 2021) into the representation

learning model to improve model’s robustness to the misspecification. Furthermore, it remains challenging to analytically determine the best divergence metric for representation balancing methods. A promising avenue for future theoretical investigations would involve developing new distributional divergences or exploring a unified theory that enables models to select appropriate divergence metrics based on the distinct data. **Empirical studies can focus on discouraging the redundancy of the concatenation fusion of each decomposed pattern and improving the efficacy of the multi-task learning objectives.** While we have followed the same approach as previous studies by evaluating model performance using simulated and semi-synthetic data, it is crucial for future research to explore the creation of appropriate benchmark datasets (Athey & Wager, 2019; Curth et al., 2021) for assessing the performance of ITE estimation methods in real-world scenarios.

References

- Tameem Adel, Isabel Valera, Zoubin Ghahramani, and Adrian Weller. One-network adversarial fairness. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 2412–2420, 2019.
- Joshua D Angrist, Guido W Imbens, and Donald B Rubin. Identification of causal effects using instrumental variables. *Journal of the American statistical Association*, 91(434):444–455, 1996.
- Serge Assaad, Shuxi Zeng, Chenyang Tao, Shounak Datta, Nikhil Mehta, Ricardo Henao, Fan Li, and Lawrence Carin. Counterfactual representation learning with balancing weights. In *International Conference on Artificial Intelligence and Statistics*, pp. 1972–1980. PMLR, 2021.
- Susan Athey and Stefan Wager. Estimating treatment effects with causal forests: An application. *Observational Studies*, 5(2):37–51, 2019.
- Peter C Austin and Elizabeth A Stuart. Moving towards best practice when using inverse probability of treatment weighting (iptw) using the propensity score to estimate causal treatment effects in observational studies. *Statistics in medicine*, 34(28):3661–3679, 2015.
- Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2):423–443, 2018.
- Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations for domain adaptation. *Advances in neural information processing systems*, 19, 2006.
- Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1):151–175, 2010.
- Alex Beutel, Jilin Chen, Zhe Zhao, and Ed H Chi. Data decisions and theoretical implications when adversarially learning fair representations. *arXiv preprint arXiv:1707.00075*, 2017.
- Ioana Bica, Ahmed M Alaa, Craig Lambert, and Mihaela Van Der Schaar. From real-world patient data to individualized treatment effects using machine learning: current and future methods to address underlying challenges. *Clinical Pharmacology & Therapeutics*, 109(1):87–100, 2021a.
- Ioana Bica, Daniel Jarrett, Alihan Hüyük, and Mihaela van der Schaar. Learning “what-if” explanations for sequential decision-making. In *International Conference on Learning Representations*, 2021b. URL <https://openreview.net/forum?id=h0de3QWtGG>.
- Stephen Burgess, Dylan S Small, and Simon G Thompson. A review of instrumental variable estimators for mendelian randomization. *Statistical methods in medical research*, 26(5):2333–2355, 2017.
- Marco Caliendo and Sabine Kopeinig. Some practical guidance for the implementation of propensity score matching. *Journal of economic surveys*, 22(1):31–72, 2008.
- Defu Cao, James Enouen, Yujing Wang, Xiangchen Song, Chuizheng Meng, Hao Niu, and Yan Liu. Estimating treatment effects from irregular time series observations with hidden confounders. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 6897–6905, 2023.

- Minwoo Chae and Stephen G Walker. Wasserstein upper bounds of the total variation for smooth densities. Statistics & Probability Letters, 163:108771, 2020.
- Pin-Yu Chen, Yash Sharma, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. Ead: elastic-net attacks to deep neural networks via adversarial examples. In Proceedings of the AAAI conference on artificial intelligence, volume 32, 2018.
- Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. Double/debiased machine learning for treatment and structural parameters. The Econometrics Journal, 21(1):C1–C68, 2018.
- Hugh A Chipman, Edward I George, and Robert E McCulloch. Bart: Bayesian additive regression trees. The Annals of Applied Statistics, 4(1):266–298, 2010.
- Zhixuan Chu, Stephen L. Rathbun, and Sheng Li. Graph infomax adversarial learning for treatment effect estimation with networked observational data. In KDD, pp. 176–184, 2021. URL <https://doi.org/10.1145/3447548.3467302>.
- Michael Crawshaw. Multi-task learning with deep neural networks: A survey. arXiv preprint arXiv:2009.09796, 2020.
- Richard K Crump, V Joseph Hotz, Guido W Imbens, and Oscar A Mitnik. Nonparametric tests for treatment effect heterogeneity. The Review of Economics and Statistics, 90(3):389–405, 2008.
- Alicia Curth and Mihaela van der Schaar. On inductive biases for heterogeneous treatment effect estimation. Advances in Neural Information Processing Systems, 34:15883–15894, 2021.
- Alicia Curth and Mihaela Van Der Schaar. In search of insights, not magic bullets: Towards demystification of the model selection dilemma in heterogeneous treatment effect estimation. In Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pp. 6623–6642. PMLR, 23–29 Jul 2023.
- Alicia Curth, David Svensson, Jim Weatherall, and Mihaela van der Schaar. Really doing great at estimating cate? a critical look at ml benchmarking practices in treatment effect estimation. In Thirty-fifth conference on neural information processing systems datasets and benchmarks track (round 2), 2021.
- Marco Cuturi and Arnaud Doucet. Fast computation of wasserstein barycenters. In International conference on machine learning, pp. 685–693. PMLR, 2014.
- Pedro Domingos. A unified bias-variance decomposition. In Proceedings of 17th international conference on machine learning, pp. 231–238. Morgan Kaufmann Stanford, 2000.
- Vincent Dorie. Nonparametric methods for causal inference. <https://github.com/vdorie/npci>, 2021.
- Xin Du, Lei Sun, Wouter Duivesteijn, Alexander Nikolaev, and Mykola Pechenizkiy. Adversarial balancing-based representation learning for causal effect inference with observational data. Data Mining and Knowledge Discovery, 35(4):1713–1738, 2021.
- Harrison Edwards and Amos Storkey. Censoring representations with an adversary. arXiv preprint arXiv:1511.05897, 2015.
- Max H Farrell. Robust inference on average treatment effects with possibly more covariates than observations. Journal of Econometrics, 189(1):1–23, 2015.
- Rui Feng, Yang Yang, Yuehan Lyu, Chenhao Tan, Yizhou Sun, and Chunping Wang. Learning fair representations via an adversarial framework. arXiv preprint arXiv:1904.13341, 2019.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. The journal of machine learning research, 17(1):2096–2030, 2016.

- Stuart Geman, Elie Bienenstock, and René Doursat. Neural networks and the bias/variance dilemma. Neural computation, 4(1):1–58, 1992.
- Ruocheng Guo, Jundong Li, Yichuan Li, K Selçuk Candan, Adrienne Raglin, and Huan Liu. Ignite: A minimax game toward learning individual treatment effects from networked observational data. In IJCAI, pp. 4534–4540, 2020a.
- Ruocheng Guo, Jundong Li, and Huan Liu. Learning individual causal effects from networked observational data. In Proceedings of the 13th International Conference on Web Search and Data Mining, pp. 232–240, 2020b.
- Yaru Hao, Xiao-Yuan Jing, Runhang Chen, and Wei Liu. Learning enhanced specific representations for multi-view feature learning. Knowledge-Based Systems, 272:110590, 2023.
- Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. Advances in neural information processing systems, 29, 2016.
- Negar Hassanpour and Russell Greiner. Counterfactual regression with importance sampling weights. In IJCAI, pp. 5880–5887, 2019a.
- Negar Hassanpour and Russell Greiner. Learning disentangled representations for counterfactual regression. In International Conference on Learning Representations, 2019b.
- Jennifer L Hill. Bayesian nonparametric modeling for causal inference. Journal of Computational and Graphical Statistics, 20(1):217–240, 2011.
- Arthur E Hoerl and Robert W Kennard. Ridge regression: Biased estimation for nonorthogonal problems. Technometrics, 12(1):55–67, 1970.
- Cong Hu, Yuanbo Li, Zhenhua Feng, and Xiaojun Wu. Attention-guided evolutionary attack with elastic-net regularization on face recognition. Pattern recognition, 143:109760, 2023.
- Yiyan Huang, Cheuk Hang Leung, Xing Yan, Qi Wu, Nanbo Peng, Dongdong Wang, and Zhixiang Huang. The causal learning of retail delinquency. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pp. 204–212, 2021.
- Yiyan Huang, Cheuk Hang Leung, Shumin Ma, Qi Wu, Dongdong Wang, and Zhixiang Huang. Moderately-balanced representation learning for treatment effects with orthogonality information. In Pacific Rim International Conference on Artificial Intelligence, pp. 3–16. Springer, 2022.
- Andrew Jesson, Sören Mindermann, Yarin Gal, and Uri Shalit. Quantifying ignorance in individual-level causal-effect estimates under hidden confounding. In International Conference on Machine Learning, pp. 4829–4838. PMLR, 2021.
- Xiaodong Jia, Xiao-Yuan Jing, Xiaoke Zhu, Songcan Chen, Bo Du, Ziyun Cai, Zhenyu He, and Dong Yue. Semi-supervised multi-view deep discriminant representation learning. IEEE transactions on pattern analysis and machine intelligence, 43(7):2496–2509, 2020.
- Ray Jiang, Aldo Pacchiano, Tom Stepleton, Heinrich Jiang, and Silvia Chiappa. Wasserstein fair classification. In Uncertainty in artificial intelligence, pp. 862–872. PMLR, 2020.
- Fredrik Johansson, Uri Shalit, and David Sontag. Learning representations for counterfactual inference. In International conference on machine learning, pp. 3020–3029. PMLR, 2016.
- Fredrik D. Johansson, Uri Shalit, Nathan Kallus, and David Sontag. Generalization bounds and representation learning for estimation of potential outcomes and causal effects. Journal of Machine Learning Research, 23(166):1–50, 2022a. URL <http://jmlr.org/papers/v23/19-511.html>.
- Fredrik D Johansson, Uri Shalit, Nathan Kallus, and David Sontag. Generalization bounds and representation learning for estimation of potential outcomes and causal effects. Journal of Machine Learning Research, 23(166):1–50, 2022b.

- Nathan Kallus, Xiaojie Mao, and Angela Zhou. Interval estimation of individual-level causal effects under unobserved confounding. In The 22nd international conference on artificial intelligence and statistics, pp. 2281–2290. PMLR, 2019.
- Guoliang Kang, Jun Li, and Dacheng Tao. Shakeout: A new approach to regularized deep neural network training. IEEE transactions on pattern analysis and machine intelligence, 40(5):1245–1258, 2017.
- Edward H Kennedy. Towards optimal doubly robust estimation of heterogeneous causal effects. Electronic Journal of Statistics, 17(2):3008–3049, 2023.
- Kun Kuang, Peng Cui, Bo Li, Meng Jiang, Shiqiang Yang, and Fei Wang. Treatment effect estimation with data-driven variable decomposition. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 31, 2017.
- Ananya Kumar, Aditi Raghunathan, Robbie Matthew Jones, Tengyu Ma, and Percy Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution. In International Conference on Learning Representations, 2021.
- Ananya Kumar, Tengyu Ma, Percy Liang, and Aditi Raghunathan. Calibrated ensembles can mitigate accuracy tradeoffs under distribution shift. In Uncertainty in Artificial Intelligence, pp. 1041–1051. PMLR, 2022.
- Sören R Künnel, Jasjeet S Sekhon, Peter J Bickel, and Bin Yu. Metalearners for estimating heterogeneous treatment effects using machine learning. Proceedings of the national academy of sciences, 116(10):4156–4165, 2019.
- Jinxing Li, Chuha Zhou, Xiaoqiang Ji, Mu Li, Guangming Lu, Yong Xu, and David Zhang. Multi-view instance attention fusion network for classification. Information Fusion, 101:101974, 2024.
- Yingming Li, Ming Yang, and Zhongfei Zhang. A survey of multi-view representation learning. IEEE transactions on knowledge and data engineering, 31(10):1863–1883, 2018.
- Qi Liu, Zhilong Zhou, Gangwei Jiang, Tiezheng Ge, and Defu Lian. Deep task-specific bottom representation network for multi-task recommendation. In Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, pp. 1637–1646, 2023.
- Mingsheng Long, Jianmin Wang, Guiguang Ding, Jiaguang Sun, and Philip S Yu. Transfer joint matching for unsupervised domain adaptation. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1410–1417, 2014.
- Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In International conference on machine learning, pp. 97–105. PMLR, 2015.
- Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Deep transfer learning with joint adaptation networks. In International conference on machine learning, pp. 2208–2217. PMLR, 2017.
- Christos Louizos, Uri Shalit, Joris M Mooij, David Sontag, Richard Zemel, and Max Welling. Causal effect inference with deep latent-variable models. Advances in neural information processing systems, 30, 2017.
- David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. Learning adversarially fair and transferable representations. In International Conference on Machine Learning, pp. 3384–3393. PMLR, 2018.
- Divyat Mahajan, Ioannis Mitliagkas, Brady Neal, and Vasilis Syrgkanis. Empirical analysis of model selection for heterogeneous causal effect estimation. International Conference on Learning Representations, 2024.
- Maggie Makar, Fredrik Johansson, John Guttag, and David Sontag. Estimation of bounds on potential outcomes for decision making. In International Conference on Machine Learning, pp. 6661–6671. PMLR, 2020.

- Wang Miao, Wenjie Hu, Elizabeth L Ogburn, and Xiao-Hua Zhou. Identifying effects of multiple treatments in the presence of unmeasured confounding. Journal of the American Statistical Association, 118(543):1953–1967, 2023.
- Naveen Naidu Narisetty. Bayesian model selection for high-dimensional data. In Handbook of statistics, volume 43, pp. 207–248. Elsevier, 2020.
- Xinkun Nie and Stefan Wager. Quasi-oracle estimation of heterogeneous treatment effects. Biometrika, 108(2):299–319, 2021.
- Ana Rita Nogueira, João Gama, and Carlos Abreu Ferreira. Causal discovery in machine learning: Theories and applications. Journal of Dynamics & Games, 8(3), 2021.
- Miruna Oprescu, Vasilis Syrgkanis, and Zhiwei Steven Wu. Orthogonal random forest for causal inference. In International Conference on Machine Learning, pp. 4932–4941. PMLR, 2019.
- Judea Pearl. Causality. Cambridge university press, 2009.
- Zhaozhi Qian, Yao Zhang, Ioana Bica, Angela Wood, and Mihaela van der Schaar. Synctwin: Treatment effect estimation with longitudinal outcomes. Advances in Neural Information Processing Systems, 34:3178–3190, 2021.
- Paul R Rosenbaum and Donald B Rubin. The central role of the propensity score in observational studies for causal effects. Biometrika, 70(1):41–55, 1983.
- Donald B Rubin. Causal inference using potential outcomes: Design, modeling, decisions. Journal of the American Statistical Association, 100(469):322–331, 2005.
- Uri Shalit, Fredrik D Johansson, and David Sontag. Estimating individual treatment effect: generalization bounds and algorithms. In International Conference on Machine Learning, pp. 3076–3085. PMLR, 2017.
- Jian Shen, Yanru Qu, Weinan Zhang, and Yong Yu. Wasserstein distance guided representation learning for domain adaptation. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 32, 2018.
- Claudia Shi, David Blei, and Victor Veitch. Adapting neural networks for the estimation of treatment effects. Advances in neural information processing systems, 32, 2019.
- Changjian Shui, Fan Zhou, Christian Gagné, and Boyu Wang. Deep active learning: Unified and principled method for query and training. In International Conference on Artificial Intelligence and Statistics, pp. 1308–1318. PMLR, 2020.
- Bharath K Sriperumbudur, Kenji Fukumizu, Arthur Gretton, Bernhard Schölkopf, and Gert RG Lanckriet. On the empirical estimation of integral probability metrics. Electronic Journal of Statistics, 6:1550–1599, 2012.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. Journal of the Royal Statistical Society Series B: Statistical Methodology, 58(1):267–288, 1996.
- Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 7167–7176, 2017.
- Giorgio Valentini and Thomas G Dietterich. Bias-variance analysis of support vector machines for the development of svm-based ensemble methods. Journal of Machine Learning Research, 5(Jul):725–775, 2004.
- Cédric Villani et al. Optimal transport: old and new, volume 338. Springer, 2009.
- Matthew J Vowels, Necati Cihan Camgoz, and Richard Bowden. D’ya like dags? a survey on structure learning and causal discovery. ACM Computing Surveys, 55(4):1–36, 2022.

- Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using random forests. Journal of the American Statistical Association, 113(523):1228–1242, 2018.
- Anpeng Wu, Kun Kuang, Bo Li, and Fei Wu. Instrumental variable regression with confounder balancing. In International Conference on Machine Learning, pp. 24056–24075. PMLR, 2022.
- Zhen Wu, Xin-Yu Dai, Cunyan Yin, Shujian Huang, and Jiajun Chen. Improving review representations with user attention and product attention for sentiment classification. In Proceedings of the AAAI conference on artificial intelligence, volume 32, 2018.
- Lihu Xu, Fang Yao, Qiuran Yao, and Huiming Zhang. Non-asymptotic guarantees for robust statistical learning under infinite variance assumption. Journal of Machine Learning Research, 24(92):1–46, 2023a.
- Peng Xu, Xiatian Zhu, and David A Clifton. Multimodal learning with transformers: A survey. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2023b.
- Xiaoqiang Yan, Shizhe Hu, Yiqiao Mao, Yangdong Ye, and Hui Yu. Deep multi-view learning methods: A review. Neurocomputing, 448:106–129, 2021.
- Zitong Yang, Yaodong Yu, Chong You, Jacob Steinhardt, and Yi Ma. Rethinking bias-variance trade-off for generalization of neural networks. In International Conference on Machine Learning, pp. 10767–10777. PMLR, 2020.
- Liuyi Yao, Sheng Li, Yaliang Li, Mengdi Huai, Jing Gao, and Aidong Zhang. Representation learning for treatment effect estimation from observational data. Advances in Neural Information Processing Systems, 31, 2018.
- Jinsung Yoon, James Jordon, and Mihaela Van Der Schaar. Ganite: Estimation of individualized treatment effects using generative adversarial nets. In International Conference on Learning Representations, 2018.
- Junkun Yuan, Xu Ma, Ruoxuan Xiong, Mingming Gong, Xiangyu Liu, Fei Wu, Lanfen Lin, and Kun Kuang. Instrumental variable-driven domain generalization with unobserved confounders. ACM Transactions on Knowledge Discovery from Data, 17(8):1–21, 2023.
- Alessio Zanga, Elif Ozkirimli, and Fabio Stella. A survey on causal discovery: Theory and practice. International Journal of Approximate Reasoning, 151:101–129, 2022.
- Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations. In International conference on machine learning, pp. 325–333. PMLR, 2013.
- Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating unwanted biases with adversarial learning. In Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society, pp. 335–340, 2018.
- Yao Zhang, Alexis Bellot, and Mihaela Schaar. Learning overlapping representations for the estimation of individualized treatment effects. In International Conference on Artificial Intelligence and Statistics, pp. 1005–1014. PMLR, 2020.
- Han Zhao and Geoffrey J Gordon. Inherent tradeoffs in learning fair representations. Journal of Machine Learning Research, 23(57):1–26, 2022.
- Han Zhao, Amanda Coston, Tameem Adel, and Geoffrey J Gordon. Conditional learning of fair representations. In International Conference on Learning Representations, 2019a.
- Han Zhao, Remi Tachet Des Combes, Kun Zhang, and Geoffrey Gordon. On learning invariant representations for domain adaptation. In International conference on machine learning, pp. 7523–7532. PMLR, 2019b.
- Han Zhao, Chen Dan, Bryon Aragam, Tommi S Jaakkola, Geoffrey J Gordon, and Pradeep Ravikumar. Fundamental limits and tradeoffs in invariant representation learning. Journal of Machine Learning Research, 23(340):1–49, 2022.

YUAN Zhiri, HU Xixu, WU Qi, MA Shumin, Cheuk Hang Leung, Xin Shen, and Yiyan Huang. A unified domain adaptation framework with distinctive divergence analysis. 2022.

Ding-Xuan Zhou. On grouping effect of elastic net. *Statistics & Probability Letters*, 83(9):2108–2112, 2013.

Indre Zliobaite. On the relation between accuracy and fairness in binary classification. *arXiv preprint arXiv:1505.05723*, 2015.

Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 67(2):301–320, 2005.

A Appendix

A.1 Proof of Lemma 1

Proof of L taking the squared loss, i.e., $L(y_1, y_2) = (y_1 - y_2)^2$:

Proof. We denote $\epsilon_{PEHE}(f) = \epsilon_{PEHE}(h, \Phi)$, $\epsilon_F(f) = \epsilon_F(h, \Phi)$, $\epsilon_{CF}(f) = \epsilon_{CF}(h, \Phi)$ for $f(\mathbf{x}, t) = h(\Phi(\mathbf{x}), t)$.

$$\begin{aligned} & \epsilon_F(f) \\ &= \int_{\mathcal{X} \times \{0,1\} \times \mathcal{Y}} (f(\mathbf{x}, t) - y^t)^2 p(y^t | \mathbf{x}) p(\mathbf{x}, t) dy^t d\mathbf{x} dt \\ &= \int_{\mathcal{X} \times \{0,1\} \times \mathcal{Y}} (f(\mathbf{x}, t) - \tau^t(\mathbf{x}))^2 p(y^t | \mathbf{x}) p(\mathbf{x}, t) dy^t d\mathbf{x} dt \\ &\quad + \int_{\mathcal{X} \times \{0,1\} \times \mathcal{Y}} (\tau^t(\mathbf{x}) - y^t)^2 p(y^t | \mathbf{x}) p(\mathbf{x}, t) dy^t d\mathbf{x} dt \\ &\quad + 2 \int_{\mathcal{X} \times \{0,1\} \times \mathcal{Y}} (f(\mathbf{x}, t) - \tau^t(\mathbf{x})) (\tau^t(\mathbf{x}) - y^t) p(y^t | \mathbf{x}) p(\mathbf{x}, t) dy^t d\mathbf{x} dt \end{aligned} \quad (28)$$

$$= \int_{\mathcal{X} \times \{0,1\}} (f(\mathbf{x}, t) - \tau^t(\mathbf{x}))^2 p(\mathbf{x}, t) d\mathbf{x} dt + \sigma_{y^t}^2(p(\mathbf{x}, t)) \quad (29)$$

Equation (29) is by the definition of $\sigma_{y^t}^2(p(\mathbf{x}, t))$ in Lemma 1 and equation (28) equaling zero since $\tau^t(\mathbf{x}) = \int_{\mathcal{Y}} y^t p(y^t | \mathbf{x}) dy^t$. A similar result can be obtained for ϵ_{CF} :

$$\epsilon_{CF}(f) = \int_{\mathcal{X} \times \{0,1\}} (f(\mathbf{x}, t) - \tau^t(\mathbf{x}))^2 p(\mathbf{x}, 1 - t) d\mathbf{x} dt + \sigma_{y^t}^2(p(\mathbf{x}, 1 - t)).$$

$$\begin{aligned} & \epsilon_{PEHE}(f) \\ &= \int_{\mathcal{X}} ((f(\mathbf{x}, 1) - f(\mathbf{x}, 0)) - (\tau^1(\mathbf{x}) - \tau^0(\mathbf{x})))^2 p(\mathbf{x}) d\mathbf{x} \\ &\leq 2 \int_{\mathcal{X}} ((f(\mathbf{x}, 1) - \tau^1(\mathbf{x}))^2 + (f(\mathbf{x}, 0) - \tau^0(\mathbf{x}))^2) p(\mathbf{x}) d\mathbf{x} \end{aligned} \quad (30)$$

$$\begin{aligned} &= 2 \int_{\mathcal{X}} (f(\mathbf{x}, 1) - \tau^1(\mathbf{x}))^2 p(\mathbf{x}, T = 1) d\mathbf{x} + 2 \int_{\mathcal{X}} (f(\mathbf{x}, 0) - \tau^0(\mathbf{x}))^2 p(\mathbf{x}, T = 0) d\mathbf{x} \\ &\quad + 2 \int_{\mathcal{X}} (f(\mathbf{x}, 1) - \tau^1(\mathbf{x}))^2 p(\mathbf{x}, T = 0) d\mathbf{x} + 2 \int_{\mathcal{X}} (f(\mathbf{x}, 0) - \tau^0(\mathbf{x}))^2 p(\mathbf{x}, T = 1) d\mathbf{x} \end{aligned} \quad (31)$$

$$\begin{aligned} &= 2 \int_{\mathcal{X} \times \{0,1\}} (f(\mathbf{x}, t) - \tau^t(\mathbf{x}))^2 p(\mathbf{x}, t) d\mathbf{x} dt + 2 \int_{\mathcal{X} \times \{0,1\}} (f(\mathbf{x}, t) - \tau^t(\mathbf{x}))^2 p(\mathbf{x}, 1 - t) d\mathbf{x} dt \\ &= 2(\epsilon_F(f) - \sigma_{y^t}^2(p(\mathbf{x}, t))) + 2(\epsilon_{CF}(f) - \sigma_{y^t}^2(p(\mathbf{x}, 1 - t))). \end{aligned} \quad (32)$$

Inequality (30) is by $(x + y)^2 \leq 2(x^2 + y^2)$; equation (31) is by $p(\mathbf{x}) = p(\mathbf{x}, T = 0) + p(\mathbf{x}, T = 1)$. By (equation 32) and the definition of σ_y^2 in Lemma 1, we have

$$\epsilon_{PEHE}(h, \Phi) \leq 2(\epsilon_{CF}(h, \Phi) + \epsilon_F(h, \Phi) - 2\sigma_y^2).$$

□

Proof of L taking the absolute loss, i.e., $L(y_1, y_2) = |y_1 - y_2|$:

Proof. We denote $\epsilon_{PEHE}(f) = \epsilon_{PEHE}(h, \Phi)$, $\epsilon_F(f) = \epsilon_F(h, \Phi)$, $\epsilon_{CF}(f) = \epsilon_{CF}(h, \Phi)$ for $f(\mathbf{x}, t) = h(\Phi(\mathbf{x}), t)$.

$$\begin{aligned} \epsilon_F(f) &= \int_{\mathcal{X} \times \{0,1\} \times \mathcal{Y}} |f(\mathbf{x}, t) - y^t| p(y^t | \mathbf{x}) p(\mathbf{x}, t) dy^t d\mathbf{x} dt \\ &\geq \int_{\mathcal{X} \times \{0,1\} \times \mathcal{Y}} |f(\mathbf{x}, t) - \tau^t(\mathbf{x})| p(y^t | \mathbf{x}) p(\mathbf{x}, t) dy^t d\mathbf{x} dt \\ &\quad - \int_{\mathcal{X} \times \{0,1\} \times \mathcal{Y}} |\tau^t(\mathbf{x}) - y^t| p(y^t | \mathbf{x}) p(\mathbf{x}, t) dy^t d\mathbf{x} dt \end{aligned} \quad (33)$$

$$= \int_{\mathcal{X} \times \{0,1\}} |f(\mathbf{x}, t) - \tau^t(\mathbf{x})| p(\mathbf{x}, t) d\mathbf{x} dt - A_{y^t}(p(\mathbf{x}, t)). \quad (34)$$

Inequality (33) is by $|x + y| \geq |x| - |y|$, equation (34) is by the definition of $A_{y^t}(p(\mathbf{x}, t))$ in Lemma 1. A similar result can be obtained for ϵ_{CF} :

$$\epsilon_{CF}(f) \geq \int_{\mathcal{X} \times \{0,1\}} |f(\mathbf{x}, t) - \tau^t(\mathbf{x})| p(\mathbf{x}, 1 - t) d\mathbf{x} dt - A_{y^t}(p(\mathbf{x}, 1 - t)).$$

$$\begin{aligned} \epsilon_{PEHE}(f) &= \int_{\mathcal{X}} |(f(\mathbf{x}, 1) - f(\mathbf{x}, 0)) - (\tau^1(\mathbf{x}) - \tau^0(\mathbf{x}))| p(\mathbf{x}) d\mathbf{x} \\ &\leq \int_{\mathcal{X}} (|f(\mathbf{x}, 1) - \tau^1(\mathbf{x})| + |f(\mathbf{x}, 0) - \tau^0(\mathbf{x})|) p(\mathbf{x}) d\mathbf{x} \end{aligned} \quad (35)$$

$$= \int_{\mathcal{X}} |f(\mathbf{x}, 1) - \tau^1(\mathbf{x})| p(\mathbf{x}, T = 1) d\mathbf{x} + \int_{\mathcal{X}} |f(\mathbf{x}, 1) - \tau^1(\mathbf{x})| p(\mathbf{x}, T = 0) d\mathbf{x} \quad (36)$$

$$+ \int_{\mathcal{X}} |f(\mathbf{x}, 0) - \tau^0(\mathbf{x})| p(\mathbf{x}, T = 0) d\mathbf{x} + \int_{\mathcal{X}} |f(\mathbf{x}, 0) - \tau^0(\mathbf{x})| p(\mathbf{x}, T = 1) d\mathbf{x} \quad (37)$$

$$\begin{aligned} &= \int_{\mathcal{X} \times \{0,1\}} |f(\mathbf{x}, t) - \tau^t(\mathbf{x})| p(\mathbf{x}, t) d\mathbf{x} dt + \int_{\mathcal{X} \times \{0,1\}} |f(\mathbf{x}, t) - \tau^t(\mathbf{x})| p(\mathbf{x}, 1 - t) d\mathbf{x} dt \\ &\leq \epsilon_F(f) + A_{y^t}(p(\mathbf{x}, t)) + \epsilon_{CF}(f) + A_{y^t}(p(\mathbf{x}, 1 - t)). \end{aligned} \quad (38)$$

Inequality (35) is by $|x + y| \leq |x| + |y|$. Equation (36) and equation (37) are by $p(\mathbf{x}) = p(\mathbf{x}, T = 0) + p(\mathbf{x}, T = 1)$. By equation (38) and the definition of A_y in Lemma 1, we have

$$\begin{aligned} \epsilon_{PEHE}(h, \Phi) &\leq \epsilon_F(f) + A_{y^t}(p(\mathbf{x}, t)) + \epsilon_{CF}(f) + A_{y^t}(p(\mathbf{x}, 1 - t)) \\ &\leq \epsilon_{CF}(h, \Phi) + \epsilon_F(h, \Phi) + 2A_y. \end{aligned}$$

□

A.2 Proof of Theorem 1

Proof of equation (4):

Proof.

$$\begin{aligned}
& \epsilon_{CF}(h, \Phi) - [(1-u) \cdot \epsilon_F^{T=1}(h, \Phi) + u \cdot \epsilon_F^{T=0}(h, \Phi)] \\
&= [(1-u) \cdot \epsilon_{CF}^{T=1}(h, \Phi) + u \cdot \epsilon_{CF}^{T=0}(h, \Phi)] - [(1-u) \cdot \epsilon_F^{T=1}(h, \Phi) + u \cdot \epsilon_F^{T=0}(h, \Phi)] \\
&= (1-u) \cdot [\epsilon_{CF}^{T=1}(h, \Phi) - \epsilon_F^{T=1}(h, \Phi)] + u \cdot [\epsilon_{CF}^{T=0}(h, \Phi) - \epsilon_F^{T=0}(h, \Phi)] \\
&= (1-u) \int_{\mathcal{X}} \ell_{h,\Phi}(\mathbf{x}, 1)(p^{T=0}(\mathbf{x}) - p^{T=1}(\mathbf{x}))d\mathbf{x} + u \int_{\mathcal{X}} \ell_{h,\Phi}(\mathbf{x}, 0)(p^{T=1}(\mathbf{x}) - p^{T=0}(\mathbf{x}))d\mathbf{x} \\
&= (1-u) \int_{\mathcal{R}} \ell_{h,\Phi}(\Psi(\mathbf{r}), 1)(p_{\Phi}^{T=0}(\mathbf{r}) - p_{\Phi}^{T=1}(\mathbf{r}))d\mathbf{r} + u \int_{\mathcal{R}} \ell_{h,\Phi}(\Psi(\mathbf{r}), 0)(p_{\Phi}^{T=1}(\mathbf{r}) - p_{\Phi}^{T=0}(\mathbf{r}))d\mathbf{r} \quad (39)
\end{aligned}$$

$$\begin{aligned}
&= B_{\Phi} \cdot (1-u) \int_{\mathcal{R}} \frac{1}{B_{\Phi}} \ell_{h,\Phi}(\Psi(\mathbf{r}), 1)(p_{\Phi}^{T=0}(\mathbf{r}) - p_{\Phi}^{T=1}(\mathbf{r}))d\mathbf{r} \\
&\quad + B_{\Phi} \cdot u \int_{\mathcal{R}} \frac{1}{B_{\Phi}} \ell_{h,\Phi}(\Psi(\mathbf{r}), 0)(p_{\Phi}^{T=1}(\mathbf{r}) - p_{\Phi}^{T=0}(\mathbf{r}))d\mathbf{r} \\
&\leq B_{\Phi} \cdot (1-u) \sup_{g \in \mathcal{G}} \left| \int_{\mathcal{R}} g(\mathbf{r})(p_{\Phi}^{T=0}(\mathbf{r}) - p_{\Phi}^{T=1}(\mathbf{r}))d\mathbf{r} \right| \\
&\quad + B_{\Phi} \cdot u \cdot \sup_{g \in \mathcal{G}} \left| \int_{\mathcal{R}} g(\mathbf{r})(p_{\Phi}^{T=1}(\mathbf{r}) - p_{\Phi}^{T=0}(\mathbf{r}))d\mathbf{r} \right| \quad (40)
\end{aligned}$$

$$= B_{\Phi} \cdot \text{Wass}(p_{\Phi}^{T=1}, p_{\Phi}^{T=0}) \quad (41)$$

Equation (39) is by the change of formula, $p_{\Phi}^{T=0}(\mathbf{r}) = p^{T=0}(\Psi(\mathbf{r}))J_{\Psi}(\mathbf{r})$, $p_{\Phi}^{T=1}(\mathbf{r}) = p^{T=1}(\Psi(\mathbf{r}))J_{\Psi}(\mathbf{r})$, where $J_{\Psi}(\mathbf{r})$ is the absolute of the determinant of the Jacobian of $\Psi(\mathbf{r})$. Equation (41) is by Definition 2. \square

Proof of equation (5):

Proof.

$$\begin{aligned}
& \epsilon_{PEHE}(h, \Phi) \\
& \leq 2(\epsilon_{CF}(h, \Phi) + \epsilon_F(h, \Phi) - 2\sigma_y^2). \quad (42)
\end{aligned}$$

$$\leq 2(\epsilon_F^{T=1}(h, \Phi) + \epsilon_F^{T=0}(h, \Phi) + B_{\Phi} \cdot \text{Wass}(p_{\Phi}^{T=1}, p_{\Phi}^{T=0}) - 2\sigma_y^2). \quad (43)$$

Inequality (42) is by equation (2) in Lemma 1. Inequality (43) is by equation 4 in Theorem 1. \square

Proof of equation (6):

Proof.

$$\begin{aligned}
& \epsilon_{PEHE}(h, \Phi) \\
& \leq \epsilon_{CF}(h, \Phi) + \epsilon_F(h, \Phi) + 2A_y \quad (44)
\end{aligned}$$

$$\leq \epsilon_F^{T=1}(h, \Phi) + \epsilon_F^{T=0}(h, \Phi) + B_{\Phi} \cdot \text{Wass}(p_{\Phi}^{T=1}, p_{\Phi}^{T=0}) + 2A_y \quad (45)$$

Inequality (44) is by equation (3) in Lemma 1. Inequality (45) is by equation 4 in Theorem 1. \square

A.3 Proof of Theorem 2

We first introduce Lemma 2 that is useful for proving Theorem 2.

Lemma 2. *Let \mathcal{G} that is defined in Definition 2 be the family of binary functions. Then we obtain $\sup_{\eta \in \mathcal{H}} \left| \int_{\mathcal{S}} \eta(s)(p_1(s) - p_2(s))ds \right| = \frac{1}{2}d_{\mathcal{H}}(p_1, p_2)$.*

Proof. Let $\mathbb{I}(\cdot)$ denotes an indicator function.

$$\begin{aligned}
& d_{\mathcal{H}}(p_1, p_2) \\
&= 2 \sup_{\eta \in \mathcal{H}} \left| \int_{\eta(s)=1} (p_1(s) - p_2(s)) ds \right| \\
&= 2 \sup_{\eta \in \mathcal{H}} \left| \int_{\mathcal{S}} \mathbb{I}(\eta(s) = 1) (p_1(s) - p_2(s)) ds \right| \\
&= 2 \sup_{\eta \in \mathcal{H}} \left| \int_{\mathcal{S}} \eta(s) (p_1(s) - p_2(s)) ds \right|
\end{aligned} \tag{46}$$

The last equation is because an indicator function is also a binary function. \square

Proof of equation (8):

Proof.

$$\begin{aligned}
& \epsilon_{CF}(h, \Phi) - [(1-u) \cdot \epsilon_F^{T=1}(h, \Phi) + u \cdot \epsilon_F^{T=0}(h, \Phi)] \\
&= (1-u) \int_{\mathcal{R}} \ell_{h,\Phi}(\Psi(\mathbf{r}), 1) (p_{\Phi}^{T=0}(\mathbf{r}) - p_{\Phi}^{T=1}(\mathbf{r})) d\mathbf{r} + u \int_{\mathcal{R}} \ell_{h,\Phi}(\Psi(\mathbf{r}), 0) (p_{\Phi}^{T=1}(\mathbf{r}) - p_{\Phi}^{T=0}(\mathbf{r})) d\mathbf{r}
\end{aligned} \tag{47}$$

$$\begin{aligned}
& \leq (1-u) \int_{p_{\Phi}^{T=0} > p_{\Phi}^{T=1}} \ell_{h,\Phi}(\Psi(\mathbf{r}), 1) (p_{\Phi}^{T=0}(\mathbf{r}) - p_{\Phi}^{T=1}(\mathbf{r})) d\mathbf{r} \\
& \quad + u \int_{p_{\Phi}^{T=1} > p_{\Phi}^{T=0}} \ell_{h,\Phi}(\Psi(\mathbf{r}), 0) (p_{\Phi}^{T=1}(\mathbf{r}) - p_{\Phi}^{T=0}(\mathbf{r})) d\mathbf{r}
\end{aligned} \tag{48}$$

$$\leq (1-u) K \int_{p_{\Phi}^{T=0} > p_{\Phi}^{T=1}} (p_{\Phi}^{T=0}(\mathbf{r}) - p_{\Phi}^{T=1}(\mathbf{r})) d\mathbf{r} + u \cdot K \int_{p_{\Phi}^{T=1} > p_{\Phi}^{T=0}} (p_{\Phi}^{T=1}(\mathbf{r}) - p_{\Phi}^{T=0}(\mathbf{r})) d\mathbf{r} \tag{49}$$

$$\begin{aligned}
&= (1-u) K \int_{\mathcal{R}} \mathbb{I}(p_{\Phi}^{T=0} > p_{\Phi}^{T=1}) (p_{\Phi}^{T=0}(\mathbf{r}) - p_{\Phi}^{T=1}(\mathbf{r})) d\mathbf{r} \\
& \quad + u \cdot K \int_{\mathcal{R}} \mathbb{I}(p_{\Phi}^{T=1} > p_{\Phi}^{T=0}) (p_{\Phi}^{T=1}(\mathbf{r}) - p_{\Phi}^{T=0}(\mathbf{r})) d\mathbf{r} \\
&\leq (1-u) K \sup_{\eta \in \mathcal{H}} \left| \int_{\mathcal{R}} \eta(\mathbf{r}) (p_{\Phi}^{T=1}(\mathbf{r}) - p_{\Phi}^{T=0}(\mathbf{r})) d\mathbf{r} \right| \\
& \quad + u \cdot K \cdot \sup_{\eta \in \mathcal{H}} \left| \int_{\mathcal{R}} \eta(\mathbf{r}) (p_{\Phi}^{T=1}(\mathbf{r}) - p_{\Phi}^{T=0}(\mathbf{r})) d\mathbf{r} \right|
\end{aligned} \tag{50}$$

$$\begin{aligned}
&\leq K \cdot \sup_{\eta \in \mathcal{H}} \left| \int_{\mathcal{R}} \eta(\mathbf{r}) ((p_{\Phi}^{T=1}(\mathbf{r}) - p_{\Phi}^{T=0}(\mathbf{r}))) d\mathbf{r} \right| \\
&= \frac{K}{2} d_{\mathcal{H}}(p_{\Phi}^{T=1}, p_{\Phi}^{T=0})
\end{aligned} \tag{51}$$

Equation (47) is derived in the same way as equation (39). Equation (48) is by $\ell_{h,\Phi} \geq 0$ for all \mathbf{r} and t . Inequality (49) is by the definition of K in Theorem 2. Inequality (50) is because an indicator function is also a binary function. Equation (51) is by Lemma 2. \square

Proof of equation (9):

Proof.

$$\begin{aligned}
& \epsilon_{PEHE}(h, \Phi) \\
&\leq 2(\epsilon_{CF}(h, \Phi) + \epsilon_F(h, \Phi) - 2\sigma_y^2)
\end{aligned} \tag{52}$$

$$\leq 2(\epsilon_F^{T=1}(h, \Phi) + \epsilon_F^{T=0}(h, \Phi) + \frac{K}{2} d_{\mathcal{H}}(p_{\Phi}^{T=1}, p_{\Phi}^{T=0}) - 2\sigma_y^2) \tag{53}$$

Inequality (52) is by equation 2 in Lemma 1. Inequality (53) is by equation 8 in Theorem 2. \square

Proof of equation (10):

Proof.

$$\begin{aligned} & \epsilon_{PEHE}(h, \Phi) \\ & \leq \epsilon_{CF}(h, \Phi) + \epsilon_F(h, \Phi) + 2A_y \end{aligned} \quad (54)$$

$$\leq \epsilon_F^{T=1}(h, \Phi) + \epsilon_F^{T=0}(h, \Phi) + \frac{K}{2} d_{\mathcal{H}}(p_{\Phi}^{T=1}, p_{\Phi}^{T=0}) + 2A_y \quad (55)$$

Inequality (54) is by equation 3 in Lemma 1. Inequality (55) is by equation 8 in Theorem 2. \square

A.4 Additional Experimental details

Additional results on Twins Benchmark. To investigate the applicability of our model DIGNet to benchmark datasets beyond the commonly used IHDP benchmark, we conducted additional comparisons with several baseline models, including linear, tree, matching, and representation learning methods, on the Twins benchmark, as presented in Table 7.

The Twins dataset comprises records of twin births in the USA between 1989 and 1991. After preprocessing, each unit contains 30 covariates relevant to parents, pregnancy, and birth. The treatment $D = 1$ indicates the heavier twin, while $D = 0$ indicates the lighter twin. The binary outcome variable Y represents 1-year mortality. For more comprehensive details on this dataset and the limitation of IHDP, refer to Curth et al. (2021).

Notably, for ϵ_{ATE} , the simple linear or matching estimator performs best across different methods. On the other hand, when assessing ITE performance using the AUC of potential outcomes, representation learning models all demonstrate strong performance, with AUC values exceeding 0.800 on both training and test sets. The observation might stem from the fact that representation balancing models are based on ITE error bounds, rather than ATE error bounds, thereby optimizing for AUC instead of ϵ_{ATE} . Moreover, among all the models, our DIGNet achieves the second-best AUC results. The best results are achieved by MBRL, which involves the orthogonality information (similar to doubly robust estimators) in representation balancing. This, in turn, inspires us to explore ATE error bounds, or consider involving doubly robust methods in future research.

Table 7: Training- & test- set AUC & ϵ_{ATE} on Twins. Mean \pm standard error of 100 runs.

	Training set		Test set	
	AUC	ϵ_{ATE}	AUC	ϵ_{ATE}
OLS/LR ₁ Johansson et al. (2016)	.660 \pm .005	.004 \pm .003	.500 \pm .028	.007 \pm .006
OLS/LR ₂ Johansson et al. (2016)	.660 \pm .004	.004 \pm .003	.500 \pm .016	.007 \pm .006
k-NN Crump et al. (2008)	.609 \pm .010	.003 \pm .002	.492 \pm .012	.005 \pm .004
BART Chipman et al. (2010)	.506 \pm .014	.121 \pm .024	.500 \pm .011	.127 \pm .024
CEVAE Louizos et al. (2017)	.845 \pm .003	.022 \pm .002	.841 \pm .004	.032 \pm .003
SITE Yao et al. (2018)	.862 \pm .002	.016 \pm .001	.853 \pm .006	.020 \pm .002
BLR Johansson et al. (2016)	.611 \pm .009	.006 \pm .004	.510 \pm .018	.033 \pm .009
BNN Johansson et al. (2016)	.690 \pm .008	.006 \pm .003	.676 \pm .008	.020 \pm .007
TARNet Shalit et al. (2017)	.849 \pm .002	.011 \pm .002	.840 \pm .006	.015 \pm .002
CFR-Wass (GNet) Shalit et al. (2017)	.850 \pm .002	.011 \pm .002	.842 \pm .005	.028 \pm .003
MBRL (Huang et al., 2022)	.879 \pm .000	.003 \pm .000	.874 \pm .001	.007 \pm .00q
DIGNet (Ours)	.874 \pm .001	.004 \pm .001	.871 \pm .001	.008 \pm .001

Implementation details. In simulation studies, we ensure a fair comparison by fixing all the hyperparameters in all datasets across different models. The relevant details are stated in Table 8. In IHDP studies,

Table 8: Hyperparameters of different models in simulation studies.

	Φ_E	Φ_G	Φ_I	π	h^1	h^0	α_1	α_2	batchsize	iteration	learning rate	learning rate for π
Gnet	(100, 100, 100, 100)	—	—	—	(100, 100)	(100, 100)	0.1	—	100	300	$1e^{-3}$	—
Inet	(100, 100, 100, 100)	—	—	(100, 100, 100)	(100, 100)	(100, 100)	—	0.1	100	300	$1e^{-3}$	$1e^{-4}$
DGNet	(100, 100, 100, 100)	(100, 100)	—	—	(100, 100)	(100, 100)	0.1	—	100	300	$1e^{-3}$	—
DINet	(100, 100, 100, 100)	—	(100, 100)	(100, 100, 100)	(100, 100)	(100, 100)	—	0.1	100	300	$1e^{-3}$	$1e^{-4}$
DIGNet	(100, 100, 100, 100)	(100, 100)	(100, 100)	(100, 100, 100)	(100, 100)	(100, 100)	0.1	0.1	100	300	$1e^{-3}$	$1e^{-4}$

to compare with the baseline model CFR-Wass (GNet), we remain the hyperparameters of INet, DGNet, DINet and the early stopping rule the same as those used in CFR-Wass Shalit et al. (2017). Since DIGNet is more complex than other four models, we adjust the hyperparameters of Φ_E , Φ_G , Φ_I , α_1 , and α_2 for DIGNet as Shalit et al. (2017) do. The relevant details are stated in Table 9.

Table 9: Hyperparameters of different models in IHDP experiments.

	Φ_E	Φ_G	Φ_I	π	h^1	h^0	α_1	α_2	batchsize	iteration	learning rate	learning rate for π
Gnet	(100, 100, 100, 100)	—	—	—	(100, 100, 100)	(100, 100, 100)	1	—	100	600	$1e^{-3}$	—
Inet	(100, 100, 100, 100)	—	—	(200, 200, 200)	(100, 100, 100)	(100, 100, 100)	—	1	100	600	$1e^{-3}$	$1e^{-3}$
DGNet	(100, 100, 100, 100)	(100, 100)	—	—	(100, 100, 100)	(100, 100, 100)	1	—	100	600	$1e^{-3}$	—
DINet	(100, 100, 100, 100)	—	(100, 100)	(200, 200, 200)	(100, 100, 100)	(100, 100, 100)	—	1	100	600	$1e^{-3}$	$1e^{-3}$
DIGNet	(100, 100, 100, 100, 100, 100)	(100, 100, 100)	(100, 100, 100)	(200, 200, 200)	(100, 100, 100)	(100, 100, 100)	1	1	100	600	$1e^{-3}$	$1e^{-3}$

Analysis of training time and training stability. We record the time it took for different models to run through 100 IHDP datasets in Table 10, and each model is trained within 600 epochs. Following Shalit et al. (2017), all models adopt the early stopping rule. We also record the average early stopping epoch on 100 runs and the actual time on 100 runs, where (actual time) = (total time) \times (average early stopping epoch)/600. Not surprisingly, GNet took the least amount of time with 3096 seconds since the objective of GNet is the simplest. However, it is very interesting that the proposed methods, DGNet and DINet, are the first two to early stop. As a result, though DGNet and DINet have multi-objectives, they spent less actual training time but achieved better ITE estimation compared to GNet and INet. Since GNet and INet are actually DGNet and DINet with PPBR ablated, we find that PPBR component can help a model achieve better ITE estimates with less time. In addition, we find that DIGNet spent the longest time to optimize since it has the most complex objective. To further study the stability of the model training, we also plot the metrics $\sqrt{\epsilon_F}$, Wass, \hat{d}_H , and $\sqrt{\epsilon_{PEHE}}$ for the first 100 epochs of each model on the first IHDP dataset in Figure 8. We find that the training process of DIGNet is stable, even steadier than GNet and INet. From this perspective, we haven’t seen a difficulty of optimizing DIGNet.

Table 10: Training time records on 100 IHDP datasets.

Model	Time for 600 epochs	Avg early stopping	Actual time	$\sqrt{\epsilon_{PEHE}}$ on test set
GNet	3096s	240.61	1241s	0.77 \pm 0.18
INet	4042s	254.19	1712s	0.72 \pm 0.11
DGNet	3775s	169.17	1064s	0.60 \pm 0.09
DINet	3212s	157.98	846s	0.60 \pm 0.11
DIGNet	4984s	226.76	1884s	0.45 \pm 0.04

We also provide the ITE and ATE estimation results on 100 IHDP datasets when the combination of (α_1, α_2) in DIGNet objective varies in $\{0.1, 0.5, 1\}$. The relevant results are reported in Table 11, indicating our DIGNet model is robust to the hyperparameters varying.

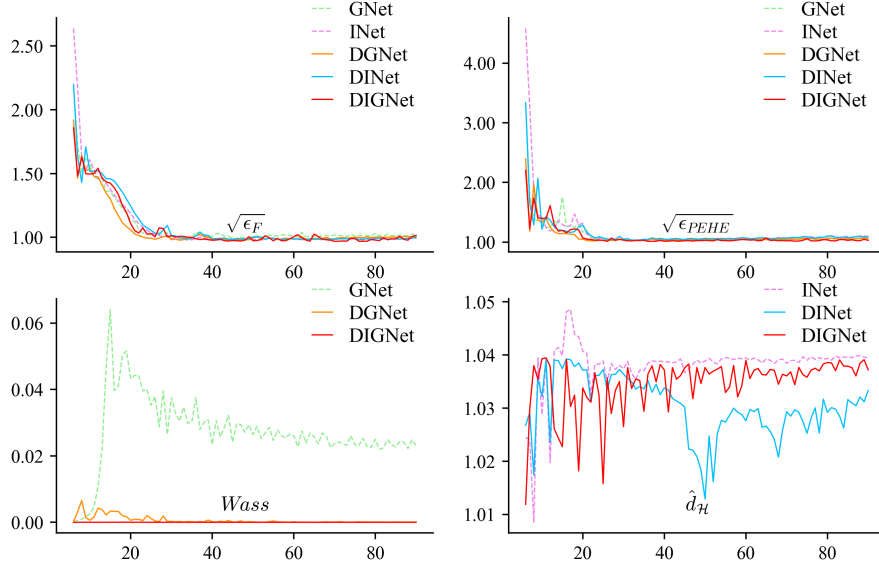


Figure 8: Training loss plots for the first 100 epochs on the first IHDP dataset.

Table 11: The results on 100 IHDP datasets with different combinations of (α_1, α_2) in DIGNet objective.

(α_1, α_2)	Training set		Test set	
	$\sqrt{\epsilon_{PEHE}}$	ϵ_{ATE}	$\sqrt{\epsilon_{PEHE}}$	ϵ_{ATE}
(0.1, 0.1)	0.407 ± 0.018	0.125 ± 0.015	0.434 ± 0.022	0.138 ± 0.016
(0.1, 0.5)	0.414 ± 0.026	0.120 ± 0.015	0.434 ± 0.028	0.123 ± 0.015
(0.1, 1)	0.416 ± 0.019	0.116 ± 0.014	0.452 ± 0.026	0.121 ± 0.015
(0.5, 0.1)	0.417 ± 0.023	0.130 ± 0.016	0.440 ± 0.026	0.137 ± 0.017
(0.5, 0.5)	0.407 ± 0.021	0.125 ± 0.015	0.416 ± 0.022	0.124 ± 0.015
(0.5, 1)	0.413 ± 0.020	0.126 ± 0.014	0.455 ± 0.028	0.133 ± 0.016
(1, 0.1)	0.411 ± 0.021	0.119 ± 0.015	0.439 ± 0.027	0.118 ± 0.015
(1, 0.5)	0.403 ± 0.020	0.118 ± 0.015	0.430 ± 0.026	0.128 ± 0.016
(1, 1)	0.402 ± 0.019	0.112 ± 0.014	0.437 ± 0.027	0.121 ± 0.015

A.5 Objectives of Different Models

Objective of GNet.

$$\min_{\Phi_E, h^t} \mathcal{L}_y(\mathbf{x}, \mathbf{t}, \mathbf{y}; \Phi_E, h^t) + \alpha_1 \mathcal{L}_G(\mathbf{x}, \mathbf{t}; \Phi_E).$$

Objective of INet.

$$\begin{aligned} & \max_{\pi} \alpha_2 \mathcal{L}_I(\mathbf{x}, \mathbf{t}; \Phi_E, \pi), \\ & \min_{\Phi_E, h^t} \mathcal{L}_y(\mathbf{x}, \mathbf{t}, \mathbf{y}; \Phi_E, h^t) + \alpha_2 \mathcal{L}_I(\mathbf{x}, \mathbf{t}; \Phi_E, \pi). \end{aligned}$$

Objective of DINet. Note that similar to DIGNet, the pre-balancing patterns are preserved by only updating Φ_I but fixing Φ_E in the second step.

$$\begin{aligned} \max_{\pi} \quad & \alpha_2 \mathcal{L}_I(\mathbf{x}, \mathbf{t}; \Phi_I \circ \Phi_E, \pi), \\ \min_{\Phi_I} \quad & \alpha_2 \mathcal{L}_I(\mathbf{x}, \mathbf{t}; \Phi_I \circ \Phi_E, \pi), \\ \min_{\Phi_E, \Phi_I, h^t} \quad & \mathcal{L}_y(\mathbf{x}, \mathbf{t}, \mathbf{y}; \Phi_E \oplus (\Phi_I \circ \Phi_E), h^t). \end{aligned}$$

Objective of DGNet. Note that similar to DIGNet, the pre-balancing patterns are preserved by only updating Φ_G but fixing Φ_E in the first step.

$$\begin{aligned} \min_{\Phi_G} \quad & \alpha_1 \mathcal{L}_G(\mathbf{x}, \mathbf{t}; \Phi_G \circ \Phi_E), \\ \min_{\Phi_E, \Phi_G, h^t} \quad & \mathcal{L}_y(\mathbf{x}, \mathbf{t}, \mathbf{y}; \Phi_E \oplus (\Phi_G \circ \Phi_E), h^t). \end{aligned}$$

Objective of DIGNet.

$$\begin{aligned} \min_{\Phi_G} \quad & \alpha_1 \mathcal{L}_G(\mathbf{x}, \mathbf{t}; \Phi_G \circ \Phi_E), \\ \max_{\pi} \quad & \alpha_2 \mathcal{L}_I(\mathbf{x}, \mathbf{t}; \Phi_I \circ \Phi_E, \pi), \\ \min_{\Phi_I} \quad & \alpha_2 \mathcal{L}_I(\mathbf{x}, \mathbf{t}; \Phi_I \circ \Phi_E, \pi), \\ \min_{\Phi_E, \Phi_I, \Phi_G, h^t} \quad & \mathcal{L}_y(\mathbf{x}, \mathbf{t}, \mathbf{y}; \Phi_E \oplus (\Phi_I \circ \Phi_E) \oplus (\Phi_G \circ \Phi_E), h^t). \end{aligned}$$