

β -INTACT-VAE: IDENTIFYING AND ESTIMATING CAUSAL EFFECTS UNDER LIMITED OVERLAP

Pengzhou (Abel) Wu & Kenji Fukumizu

Department of Statistical Science, The Graduate University for Advanced Studies
& The Institute of Statistical Mathematics

Tachikawa, Tokyo

{wu.pengzhou, fukumizu}@ism.ac.jp

ABSTRACT

As an important problem in causal inference, we discuss the identification and estimation of treatment effects (TEs) under limited overlap; that is, when subjects with certain features belong to a single treatment group. We use a latent variable to model a prognostic score which is widely used in biostatistics and sufficient for TEs; i.e., we build a generative prognostic model. We prove that the latent variable recovers a prognostic score, and the model identifies individualized treatment effects. The model is then learned as β -Intact-VAE—a new type of variational autoencoder (VAE). We derive the TE error bounds that enable representations balanced for treatment groups conditioned on individualized features. The proposed method is compared with recent methods using (semi-)synthetic datasets.

1 INTRODUCTION

Causal inference (Imbens & Rubin, 2015; Pearl, 2009), i.e, inferring causal effects of interventions, is a fundamental field of research. In this work, we focus on treatment effects (TEs) based on a set of observations comprising binary labels T for treatment/control (non-treated), outcome Y , and other covariates X . Typical examples include estimating the effects of public policies or new drugs based on the personal records of the subjects. The fundamental difficulty of causal inference is that we never observe *counterfactual* outcomes that would have been if we had made the other decision (treatment or control). While randomized controlled trials (RCTs) control biases through randomization and are ideal protocols for causal inference, they often have ethical and practical issues, or suffer from expensive costs. Thus, causal inference from observational data is important.

Causal inference from observational data has other challenges as well. One is *confounding*: there may be variables, called confounders, that causally affect both the treatment and the outcome, and spurious correlation/bias follows. The other is the systematic *imbalance* (difference) of the distributions of the covariates between the treatment and control groups—that is, X depends on T , which introduces bias in estimation. A majority of studies on causal inference, including the current work, have relied on unconfoundedness; this means that the confounding can be controlled by conditioning on the covariates. The more covariates are collected the more likely unconfoundedness holds; however, more covariates tends to introduce a stronger imbalance between treatment and control.

The current work studies the issue of imbalance in estimating individualized TEs conditioned on X . Classical approaches aim for *covariate balance*, X independent of T , by matching and re-weighting (Stuart, 2010; Rosenbaum, 2020). Machine learning methods have also been exploited; there are semi-parametric methods—e.g., Van der Laan & Rose (2018, TMLE)—which improve finite sample performance, as well as non-parametric methods—e.g., Wager & Athey (2018, CF). Notably, from Johansson et al. (2016), there has been a recent increase in interest in *balanced representation learning* (BRL) to learn representations Z of the covariates, such that Z independent of T .

The most serious form of imbalance is the *limited (or weak) overlap of covariates*, which means that sample points with certain covariate values belong to a single treatment group. In this case, a straightforward estimation of TEs is not possible at non-overlapping covariate values due to lack of data. There are works that provide robustness to limited overlap (Armstrong & Kolesár, 2021), trim non-overlapping data points (Yang & Ding, 2018), weight data points by overlap (Li & Li, 2019), or study convergence rates depending on overlap (Hong et al., 2020). Limited overlap is particularly relevant to machine learning methods that exploit high-dimensional covariates. This is because, with higher-dimensional covariates, overlap is harder to satisfy and verify (D’Amour et al., 2020).

To address imbalance and limited overlap, we use a prognostic score (Hansen, 2008); it is a sufficient statistic of outcome predictors and is among the key concepts of sufficient scores for TE estimation. As a function of covariates, it can map some non-overlapping values to an overlapping value in a space of lower-dimensions. For individualized TEs, we consider *conditionally balanced representation* Z , such that Z is independent of T given X —which, as we will see, is a necessary condition for a balanced prognostic score. Moreover, prognostic score modeling can benefit from methods in predictive analytics and exploit rich literature, particularly in medicine and health (Hajage et al., 2017). Thus, it is promising to combine the predictive power of prognostic modeling and machine learning. With this idea, our method builds on a generative prognostic model that models the prognostic score as a latent variable and factorizes to the score distribution and outcome distribution.

As we consider latent variables and causal inference, *identification* is an issue that must be discussed before estimation is considered. “Identification” means that the parameters of interest (in our case, representation function and TEs) are uniquely determined and expressed using the true observational distribution. Without identification, a consistent estimator is impossible to obtain, and a model would fail silently; in other words, the model may fit perfectly but will return an estimator that converges to a wrong one, or does not converge at all (Lewbel, 2019, particularly Sec. 8). Identification is even more important for causal inference; because, unlike usual (non-causal) model misspecification, causal assumptions are often unverifiable through observables (White & Chalak, 2013). Thus, it is critical to specify the theoretical conditions for identification, and then the applicability of the methods can be judged by knowledge of an application domain.

A major strength of our generative model is that the latent variable is identifiable. This is because the factorization of our model is naturally realized as a combination of identifiable VAE (Khemakhem et al., 2020a, iVAE) and conditional VAE (Sohn et al., 2015, CVAE). Based on model identifiability, we develop two identification results for individualized TEs under limited overlap. A similar VAE architecture was proposed in Wu & Fukumizu (2020b); the current study is different in setting, theory, learning objective, and experiments. The previous work studies unobserved confounding but not limited overlap, with different set of assumptions and identification theories. The current study further provides bounds on individualized TE error, and the bounds justify a conditionally balancing term controlled by hyperparameter β , as an interpolation between the two identifications.

In summary, we study the identification (Sec. 3) and estimation (Sec. 4) of individualized TEs under limited overlap. Our approach is based on recovering prognostic scores from observed variables. To this end, our method exploits recent advances in identifiable representation—particularly iVAE. The code is in Supplementary Material, and the proofs are in Sec. A. Our main contributions are:

- 1) TE identification under limited overlap of X , via prognostic scores and an identifiable model;
- 2) bounds on individualized TE error, which justify our conditional BRL;
- 3) a new regularized VAE, β -Intact-VAE, realizing the identification and conditional balance;
- 4) experimental comparison to the state-of-the-art methods on (semi-)synthetic datasets.

1.1 RELATED WORK

Limited overlap. Under limited overlap, Luo et al. (2017) estimate the average TE (ATE) by reducing covariates to a linear prognostic score. Farrell (2015) estimates a constant TE under a partial linear outcome model. D’Amour & Franks (2021) study the identification of ATE by a general class of scores, given the (linear) propensity score and prognostic score. Machine learning studies on this topic have focused on finding overlapping regions (Oberst et al., 2020; Dai & Stultz, 2020), or indicating possible failure under limited overlap (Jesson et al., 2020), but not remedies. An exception is Johansson et al. (2020), which provides bounds under limited overlap. To the best of our knowledge, our method is the first machine learning method that provides identification under limited overlap.

Prognostic scores have been recently combined with machine learning approaches, mainly in the biostatistics community. For example, Huang & Chan (2017) estimate individualized TE by reducing covariates to a linear score which is a joint propensity-prognostic score. Tarr & Imai (2021) use SVM to minimize the worst-case bias due to prognostic score imbalance. However, in the machine learning community, few methods consider prognostic scores; Zhang et al. (2020a) and Hassanpour & Greiner (2019) learn outcome predictors, without mentioning prognostic score—while Johansson et al. (2020) conceptually, but not formally, connects BRL to prognostic score. Our work is the first to formally connect generative learning and prognostic scores for TE estimation.

Identifiable representation. Recently, independent component analysis (ICA) and representation learning—both ill-posed inverse problems—meet together to yield nonlinear ICA and identifiable representation; for example, using VAEs (Khemakhem et al., 2020a), and energy models (Khemakhem et al., 2020b). The results are exploited in causal discovery (Wu & Fukumizu, 2020a) and out-of-distribution (OOD) generalization (Sun et al., 2020). This study is the first to explore identifiable representations in TE identification.

BRL and related methods amount to a major direction. Early BRL methods include BLR/BNN (Johansson et al., 2016) and TARnet/CFR (Shalit et al., 2017). In addition, Yao et al. (2018) exploit the local similarity between data points. Shi et al. (2019) use similar architecture to TARnet, considering the importance of treatment probability. There are also methods that use GAN (Yoon et al., 2018, GANITE) and Gaussian processes (Alaa & van der Schaar, 2017). Our method shares the idea of BRL, and further extends to conditional balance—which is natural for individualized TE.

More. Our work lays conceptual and theoretical foundations of VAE methods for TEs (e.g., CEVAE Louizos et al., 2017; Lu et al., 2020). See Sec. D for more related works, there we also make detailed comparisons to CFR and CEVAE, which are well-known machine learning methods.

2 SETUP AND PRELIMINARIES

2.1 COUNTERFACTUALS, TREATMENT EFFECTS, AND IDENTIFICATION

Following Imbens & Rubin (2015), we assume there exist *potential outcomes* $Y(t) \in \mathbb{R}^d$, $t \in \{0, 1\}$. $Y(t)$ is the outcome that would have been observed if the treatment value $T = t$ was applied. We see $Y(t)$ as the hidden variables that give the *factual outcome* Y under *factual assignment* $T = t$. Formally, $Y(t)$ is defined by the *consistency of counterfactuals*: $Y = Y(t)$ if $T = t$; or simply $Y = Y(T)$. The *fundamental problem of causal inference* is that, for a unit under research, we can observe only one of $Y(0)$ or $Y(1)$ —w.r.t. the treatment value applied. That is, “factual” refers to Y or T , which is *observable*; or estimators built on the observables. We also observe relevant covariate(s) $X \in \mathcal{X} \subseteq \mathbb{R}^m$, which is associated with individuals, with distribution $\mathcal{D} := (X, Y, T) \sim p(\mathbf{x}, \mathbf{y}, t)$. We use upper-case (e.g. T) to denote random variables, and lower-case (e.g. t) for realizations.

The expected potential outcome is denoted by $\mu_t(\mathbf{x}) = \mathbb{E}(Y(t)|X = \mathbf{x})$ conditioned on $X = \mathbf{x}$. The estimands in this work are the conditional ATE (CATE) and ATE, defined, respectively, by:

$$\tau(\mathbf{x}) = \mu_1(\mathbf{x}) - \mu_0(\mathbf{x}), \quad \nu = \mathbb{E}(\tau(X)). \quad (1)$$

CATE is seen as an *individual-level*, personalized, treatment effect, given highly discriminative X .

Standard results (Rubin, 2005)(Hernan & Robins, 2020, Ch. 3) show sufficient conditions for TE identification in general settings. They are *Exchangeability*: $Y(t) \perp\!\!\!\perp T|X$, and *Overlap*: $p(t|\mathbf{x}) > 0$ for any $\mathbf{x} \in \mathcal{X}$. Both are required for $t \in \{0, 1\}$. When t appears in statements without quantification, we always mean “for both t ”. Often, *Consistency* is also listed; however, as mentioned, it is better known as the well-definedness of counterfactuals. Exchangeability means, just as in RCTs, but additionally given X , that there is no correlation between factual T and potential $Y(t)$. Note that the popular assumption $Y(0), Y(1) \perp\!\!\!\perp T|X$ is stronger than $Y(t) \perp\!\!\!\perp T|X$ and is not necessary for identification (Hernan & Robins, 2020, pp. 15). Overlap means that the supports of $p(\mathbf{x}|t = 0)$ and $p(\mathbf{x}|t = 1)$ should be the same, and this ensures that there are data for $\mu_t(\mathbf{x})$ on any (\mathbf{x}, t) .

We rely on consistency and exchangeability, but in Sec. 3.2, will relax the condition of the overlapping covariate to allow some non-overlapping values \mathbf{x} —that is, covariate X is *limited-overlapping*. In this paper, we also discuss overlapping variables other than X (e.g., prognostic scores), and provide a definition for any random variable V with support \mathcal{V} as follows:

Definition 1. V is *Overlapping* if $p(t|V = v) > 0$ for any $t \in \{0, 1\}, v \in \mathcal{V}$. If the condition is violated at some value v , then v is *non-overlapping* and V is *limited-overlapping*.

2.2 PROGNOSTIC SCORES

Our method aims to recover a prognostic score (Hansen, 2008), adapted to account for both t as in Definition 2. On the other hand, balancing scores (Rosenbaum & Rubin, 1983) $\mathbf{b}(X)$ are defined by $T \perp\!\!\!\perp X|\mathbf{b}(X)$, of which the propensity score $p(t = 1|X)$ is a special case. See Sec. B.1 for detail.

Definition 2. A PGS is $\{p(X, t)\}_{t \in \{0,1\}}$ such that $Y(t) \perp\!\!\!\perp X | p(X, t)$, where $p(\mathbf{x}, t)$ ($p_t(\mathbf{x})$ hereafter) is a function defined on $\mathcal{X} \times \{0, 1\}$. A PGS is called *balanced* (and a *bPGS*) if $p_0 = p_1$.

We say a PGS is overlapping, if *both* $p_0(X)$ and $p_1(X)$ are overlapping. Obviously, a bPGS $p(X)$ is a conditionally balanced representation (defined as $Z \perp\!\!\!\perp T | X$ in Introduction) and is thus named. We often write t of the function argument in subscripts.

We use bPGS or PGS to construct representations for CATE estimation. **Why not balancing scores?** While balancing scores $b(X)$ have been widely used in causal inference, PGSs are more suitable for discussing overlap. Our purpose is to recover an overlapping score for limited-overlapping X . It is known that overlapping $b(X)$ implies overlapping X (D’Amour et al., 2020), which counters our purpose. In contrast, overlapping bPGS does not imply overlapping $b(X)$. **Example.** Let $T = \mathbb{I}(X + \epsilon > 0)$ and $Y = \mathbf{f}(|X|, T) + \mathbf{e}$, where \mathbb{I} is the indicator function, ϵ and \mathbf{e} are exogenous zero-mean noises, and the support of X is on the entire real line while ϵ is bounded. Now, X itself is a balancing score and $|X|$ is a bPGS; and $|X|$ is overlapping but X is not. Moreover, with theoretical and experimental evidence, it is recently conjectured that PGSs maximize overlap among a class of sufficient scores, including $b(X)$ (D’Amour & Franks, 2021). In general, Hajage et al. (2017) show that prognostic score methods perform better—or as well as—propensity score methods.

Below is a corollary of Proposition 5 in Hansen (2008); note that $p_t(X)$ satisfies exchangeability.

Proposition 1 (Identification via PGS). *If $p_t(X)$ is a PGS and $Y | p_{\hat{t}}(X), T \sim p_{Y | p_{\hat{t}}, T}(\mathbf{y} | P, t)$ where $\hat{t} \in \{0, 1\}$ is a counterfactual assignment, then CATE and ATE are identified, using (1) and*

$$\mu_{\hat{t}}(\mathbf{x}) = \mathbb{E}(Y(\hat{t}) | p_{\hat{t}}(X), X = \mathbf{x}) = \mathbb{E}(Y | p_{\hat{t}}(\mathbf{x}), T = \hat{t}) = \int p_{Y | p_{\hat{t}}, T}(\mathbf{y} | p_{\hat{t}}(\mathbf{x}), \hat{t}) \mathbf{y} d\mathbf{y} \quad (2)$$

With the knowledge of p_t and $p_{Y | p_{\hat{t}}, T}$, we choose one of p_0, p_1 and set $t = \hat{t}$ in the density function, w.r.t the $\mu_{\hat{t}}$ of interest. This counterfactual assignment resolves the problem of non-overlap at \mathbf{x} . Note that a sample point with $X = \mathbf{x}$ may not have $T = \hat{t}$.

We consider additive noise models for $Y(t)$, which ensures the existence of PGSs.

(G1)¹ (Additive noise model) the data generating process (DGP) for Y is $Y = \mathbf{f}^*(\mathbf{m}(X, T), T) + \mathbf{e}$ where \mathbf{f}^*, \mathbf{m} are functions and \mathbf{e} is a zero-mean exogenous (external) noise.

The DGP is causal and defines potential outcomes by $Y(t) := \mathbf{f}_t^*(\mathbf{m}_t(X)) + \mathbf{e}$, and specifies $\mathbf{m}(X, T), T$, and \mathbf{e} as the only direct causes of Y . Particularly, $\mathbf{m}_t(X)$ is a sufficient statistics of X for $Y(t)$. For example, 1) $\mathbf{m}_t(X)$ can be the component(s) of X that affect $Y(t)$ directly, or 2) if $Y(t) | X$ follows a generalized linear model, then $\mathbf{m}_t(X)$ can be the linear predictor of $Y(t)$.

Under **(G1)**, 1) $\mathbf{m}_t(X)$ is a PGS; 2) $\mu_t(X) = \mathbf{f}_t^*(\mathbf{m}_t(X))$ is a PGS; 3) X is a (trivial) bPGS; and 4) $u(X) := (\mu_0(X), \mu_1(X))$ is a bPGS. The **essence of our method** is to recover the PGS $\mathbf{m}_t(X)$ as a representation, assuming $\mathbf{m}_t(X)$ is not higher-dimensional than Y and approximately balanced. Note that $\mu_t(X)$, our final target, is a low-dimensional PGS but not balanced, and we estimate it conditioning on the approximate bPGS $\mathbf{m}_t(X)$.

3 IDENTIFICATION UNDER GENERATIVE PROGNOSTIC MODEL

In Sec. 3.1, we specify the generative prognostic model $p(\mathbf{y}, \mathbf{z} | \mathbf{x}, t)$, and show its identifiability. In Sec. 3.2, we prove the identification of CATEs, which is one of our main contributions. The theoretical analysis involves only our generative model (i.e., prior and decoder), but not the encoder. The encoder is not part of the generative model and is involved as an approximate posterior in the estimation, which is studied in Sec. 4.

3.1 MODEL, ARCHITECTURE, AND IDENTIFIABILITY

Our goal is to build a model that can be learned by VAE from observational data to obtain a PGS, or better, a bPGS, via the latent variable Z . The generative prognostic model of the proposed method is in (3),

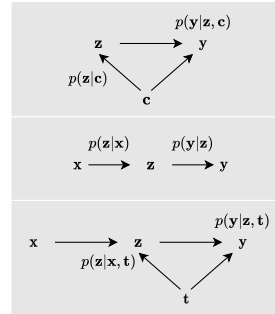


Figure 1: CVAE, iVAE, and Intact-VAE: Graphical models of the decoders.

¹The labels **G**, **M**, or **D** mean Generating process (of Y), probabilistic Model, or Distribution (of X). We introduce assumptions when appropriate but compile them in one place in Sec. C.1.

where $\theta := (\mathbf{f}, \mathbf{h}, \mathbf{k})$ contains the functional parameters. The first factor $p_{\mathbf{f}}(\mathbf{y}|\mathbf{z}, t)$, our decoder, models $p_{Y|p_t, T}(\mathbf{y}|P, t)$ in (2) and is an additive noise model, with $\epsilon \sim p_{\epsilon}$ as the exogenous noise. The second factor $p_{\lambda}(\mathbf{z}|\mathbf{x}, t)$, our conditional prior, models $p_T(X)$ and is a factorized Gaussian, with $\lambda_T(X) := \text{diag}^{-1}(\mathbf{k}_T(X))(\mathbf{h}_T(X), -\frac{1}{2})^T$ as its natural parameter in the exponential family, where $\text{diag}(\cdot)$ gives a diagonal matrix from a vector.

$$\begin{aligned} p_{\theta}(\mathbf{y}, \mathbf{z}|\mathbf{x}, t) &= p_{\mathbf{f}}(\mathbf{y}|\mathbf{z}, t)p_{\lambda}(\mathbf{z}|\mathbf{x}, t), \\ p_{\mathbf{f}}(\mathbf{y}|\mathbf{z}, t) &= p_{\epsilon}(\mathbf{y} - \mathbf{f}_t(\mathbf{z})), \quad p_{\lambda}(\mathbf{z}|\mathbf{x}, t) \sim \mathcal{N}(\mathbf{z}; \mathbf{h}_t(\mathbf{x}), \text{diag}(\mathbf{k}_t(\mathbf{x}))). \end{aligned} \quad (3)$$

We denote $n := \dim(Z)$. For inference, the ELBO is given by the standard variational lower bound

$$\log p(\mathbf{y}|\mathbf{x}, t) \geq \mathbb{E}_{\mathbf{z} \sim q} \log p_{\mathbf{f}}(\mathbf{y}|\mathbf{z}, t) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}, \mathbf{y}, t) \| p_{\lambda}(\mathbf{z}|\mathbf{x}, t)). \quad (4)$$

Note that the encoder q conditions on all the observables (X, Y, T) ; this fact plays an important role in Sec. 4.1. Full parameterization of the encoder and decoder is also given in Sec. 4.1. This architecture is called *Intact-VAE* (*I*dentifiable *t*reatment-*c*onditional VAE). See Figure 1 for comparison in terms of graphical models (which have *not* causal implications here). See Sec. C.2 for more expositions and Sec. B.2 for basics of VAEs.

Our model identifiability extends the theory of iVAE, and the following conditions are inherited.

(M1) i) \mathbf{f}_t is injective, and ii) \mathbf{f}_t is differentiable.

(D1) $\lambda_t(X)$ is non-degenerate, i.e., the linear hull of its support is $2n$ -dimensional.

Under **(M1)** and **(D1)**, we obtain the following identifiability of the parameters in the model: if $p_{\theta}(\mathbf{y}|\mathbf{x}, t) = p_{\theta'}(\mathbf{y}|\mathbf{x}, t)$, we have, for any \mathbf{y}_t in the image of \mathbf{f}_t :

$$\mathbf{f}_t^{-1}(\mathbf{y}_t) = \text{diag}(\mathbf{a})\mathbf{f}_t^{\prime -1}(\mathbf{y}_t) + \mathbf{b} =: \mathcal{A}_t(\mathbf{f}_t^{\prime -1}(\mathbf{y}_t)) \quad (5)$$

where $\text{diag}(\mathbf{a})$ is an invertible n -diagonal matrix and \mathbf{b} is an n -vector, both of which depend on $\lambda_t(\mathbf{x})$ and $\lambda_t'(\mathbf{x})$. The essence of the result is that $\mathbf{f}_t' = \mathbf{f}_t \circ \mathcal{A}_t$; that is, \mathbf{f}_t can be identified (learned) up to an affine transformation \mathcal{A}_t . See Sec. A for the proof and a relaxation of **(D1)**. In this paper, symbol $'$ (prime) always indicates another parameter (variable, etc.): $\theta' = (\mathbf{f}', \lambda')$.

3.2 IDENTIFICATIONS UNDER LIMITED-OVERLAPPING COVARIATE

In this subsection, we present two results of CATE identification based on the recovery of equivalent bPGS and PGS, respectively. Since PGSs are functions of X , the theory assumes a noiseless prior for simplicity, i.e., $\mathbf{k}(X) = \mathbf{0}$; the prior $Z_{\lambda, t} \sim p_{\lambda}(\mathbf{z}|\mathbf{x}, t)$ degenerates to function $\mathbf{h}_t(X)$.

PGSs with dimensionality lower than or equal to $d = \dim(Y)$ are essential to address limited overlapping, as shown below. We set $n = d$ because μ_t is a PGS of the same dimension as Y under **(G1)**. In practice, $n = d$ means that we seek a low-dimensional representation of X . We introduce

(G1') (Low-dimensional PGS) **(G1)** is true, and $\mu_t = \mathbf{j}_t \circ \mathbf{p}_t$ for some \mathbf{p}_t and injective \mathbf{j}_t ,

which is equivalent to **(G1)** because $\mu_t = \mathbf{j}_t \circ \mathbf{p}_t$ is trivially satisfied with \mathbf{j}_t is identity and $\mathbf{p}_t = \mu_t$. **(G1')** is used instead in this subsection. First, it explicitly restricts $\dim(\mathbf{p}_t)$ via injectivity, which ensures that $n = \dim(Y) \geq \dim(\mathbf{p}_t)$. Second, it reminds us that, possibly, the decomposition is not unique; and, clearly, all \mathbf{p}_t that satisfy **(G1')** are PGSs. For example, if \mathbf{f}_t^* is injective, then $\mathbf{j}_t = \mathbf{f}_t^*$ and $\mathbf{p}_t = \mathbf{m}_t$ satisfies $\mu_t = \mathbf{j}_t \circ \mathbf{p}_t$. Finally, it is then natural to introduce

(G2) (Low-dimensional bPGS) **(G1)** is true, and $\mu_t = \mathbf{j}_t \circ \mathbf{p}$ for some \mathbf{p} and injective \mathbf{j}_t ,

which is stronger than **(G1)**, gives bPGS $\mathbf{p}(X)$, and ensures that $n \geq \dim(\mathbf{p})$. **(G2)** is satisfied if \mathbf{f}_t^* is injective and $\mathbf{m}_0 = \mathbf{m}_1$. **(G2)** implies $\mu_1 = \mathbf{i} \circ \mu_0$ where $\mathbf{i} := \mathbf{j}_1 \circ \mathbf{j}_0^{-1}$; in words, CATEs are given by μ_0 and an invertible function. See Sec. C.3 for real-world examples and more discussions.

With **(G1')** or **(G2)**, overlapping X can be relaxed to overlapping bPGS or PGS plus the following:

(M2) (Score partition preserving) For any $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$, if $\mathbf{p}_t(\mathbf{x}) = \mathbf{p}_t(\mathbf{x}')$, then $\mathbf{h}_t(\mathbf{x}) = \mathbf{h}_t(\mathbf{x}')$.

Note that **(M2)** is only required for the optimal \mathbf{h} specified in Proposition 2 or Theorem 1. The intuition is that \mathbf{p}_t maps each non-overlapping \mathbf{x} to an overlapping value, and \mathbf{h}_t preserves this property through learning. This is non-trivial because, for a given t , some values of X are unobserved due to limited overlap. Thus, **(M2)** can be seen as a weak form of OOD generalization: the NNs for \mathbf{h} can

learn the OOD score partition. While unnecessary for us, linear \mathbf{p}_t and \mathbf{h}_t trivially imply **(M2)** and are often assumed, e.g., in Huang & Chan (2017); Luo et al. (2017); D’Amour & Franks (2021).

Our first identification, Proposition 2, relies on **(G2)** and our generative model, *without* model identifiability (so differentiable \mathbf{f}_t is not needed).

Proposition 2 (Identification via recovery of bPGS). *Suppose we have DGP **(G2)** and model (3) with $n = d$. Assume **(M1)**-i) and **(M3)** (PS matching) let $\mathbf{h}_0(X) = \mathbf{h}_1(X)$ and $\mathbf{k}(X) = \mathbf{0}$. Then, if $\mathbb{E}_{p_\theta}(Y|X, T) = \mathbb{E}(Y|X, T)$, we have*

- 1) (Recovery of bPGS) $\mathbf{z}_{\lambda, t} = \mathbf{h}_t(\mathbf{x}) = \mathbf{v}(\mathbf{p}(\mathbf{x}))$ on overlapping \mathbf{x} , where $\mathbf{v} : \mathcal{P} \rightarrow \mathbb{R}^n$ is an injective function, and $\mathcal{P} := \{\mathbf{p}(\mathbf{x}) | \text{overlapping } \mathbf{x}\}$;
- 2) (CATE identification) if $\mathbf{p}(X)$ in **(G2)** is overlapping, and **(M2)** is satisfied, then $\mu_t(\mathbf{x}) = \hat{\mu}_t(\mathbf{x}) := \mathbb{E}_{p_{\lambda}(Z|\mathbf{x}, t)} \mathbb{E}_{p_f}(Y|Z, t) = \mathbf{f}_t(\mathbf{h}_t(\mathbf{x}))$, for any $t \in \{0, 1\}$ and $\mathbf{x} \in \mathcal{X}$.

In essence, i) the true DGP is identified up to an invertible mapping \mathbf{v} , such that $\mathbf{f}_t = \hat{\mathbf{f}}_t \circ \mathbf{v}^{-1}$ and $\mathbf{h} = \mathbf{v} \circ \mathbf{p}$; and ii) \mathbf{p}_t is recovered up to \mathbf{v} , and $Y(t) \perp\!\!\!\perp X | \mathbf{p}_t(X)$ is preserved—with *same* \mathbf{v} for both t . Theorem 1 below also achieves the essence i) and ii), under $\mathbf{p}_0 \neq \mathbf{p}_1$.

The existence of bPGS is preferred, because it satisfies overlap and **(M2)** more easily than PGS which requires the conditions for each of the two functions of PGS. However, the existence of low-dimensional bPGS is uncertain in practice when our knowledge of the DGP is limited. Thus, we depend on Theorem 1 based on the model identifiability to work under PGS which generally exists.

Theorem 1 (Identification via recovery of PGS). *Suppose we have DGP **(G1’)** and model (3) with $n = d$. For the model, assume **(M1)** and **(M3’)** (Noise matching) let $p_e = p_\epsilon$ and $\mathbf{k}(X) = k\mathbf{k}'(X)$, $k \rightarrow 0$. Assume further that **(D1)** and **(D2)** (Balance from data) $\mathcal{A}_0 = \mathcal{A}_1$ in (5). Then, if $p_\theta(\mathbf{y}|\mathbf{x}, t) = p(\mathbf{y}|\mathbf{x}, t)$; conclusions 1) and 2) in Proposition 2 hold with \mathbf{p} replaced with \mathbf{p}_t in **(G1’)**; and the domain of \mathbf{v} becomes $\mathcal{P} := \{\mathbf{p}_t(\mathbf{x}) | p(t, \mathbf{x}) > 0\}$.*

Theorem 1 implies that, without bPGS, we need to know or learn the distribution of hidden noise ϵ to have $p_e = p_\epsilon$. Proposition 2 and Theorem 1 achieve recovery and identification in a complementary manner; the former starts from the prior by $\mathbf{p}_0 = \mathbf{p}_1$ and $\mathbf{h}_0 = \mathbf{h}_1$, while the latter starts from the decoder by $\mathcal{A}_0 = \mathcal{A}_1$ and $p_e = p_\epsilon$. We see that $\mathcal{A}_0 = \mathcal{A}_1$ acts as a kind of balance because it replaces $\mathbf{p}_0 = \mathbf{p}_1$ in Proposition 2. We show in Sec. A a sufficient and necessary condition **(D2’)** on data that ensures $\mathcal{A}_0 = \mathcal{A}_1$. Note that the singularities due to $k \rightarrow 0$ (e.g., $\lambda \rightarrow \mathbf{0}$) cancel out in (5). See Sec. C.4 for more on the complementarity between the two identifications.

4 ESTIMATION BY β -INTACT-VAE

4.1 PRIOR AS BPGS, POSTERIOR AS PGS, AND β AS REGULARIZATION STRENGTH

In Sec. 3.2, we see that the existence of bPGS (Proposition 2) is preferable in identifying the true DGP up to an equivalent expression—while Theorem 1 allows us to deal with PGS by adding other conditions. In learning our model with data, we formally require **(G1)** and further expect that **(G2)** holds approximately; the latter is true when \mathbf{f}_t^* is injective and $\mathbf{m}_0 \approx \mathbf{m}_1$ ($\mathbf{m}_t(X)$ is an approximate bPGS). Instead of the trivial regression $\mu_t(X) = \mathbb{E}(Y|X, T = t)$, we want to recover the approximate bPGS $\mathbf{m}_t(X)$. This idea is common in practice. For example, in a real-world nutrition study (Huang & Chan, 2017), a reduction of 11 covariates recovers a 1-dimensional linear bPGS.

We consider two ways to recover an approximate bPGS by a VAE. One is to use a prior which does not depend on t , indicating a preference for bPGS. Namely, we set $\lambda_0 = \lambda_1$, denote $\Lambda(X) := \lambda(X)$ and have $p_\Lambda(\mathbf{z}|\mathbf{x})$ as the prior in (3). The decoder and encoder are factorized Gaussians:

$$p_{f, g}(\mathbf{y}|\mathbf{z}, t) = \mathcal{N}(\mathbf{y}; \mathbf{f}_t(\mathbf{z}), \text{diag}(\mathbf{g}_t(\mathbf{z}))), \quad q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}, t) = \mathcal{N}(\mathbf{z}; \mathbf{r}_t(\mathbf{x}, \mathbf{y}), \text{diag}(\mathbf{s}_t(\mathbf{x}, \mathbf{y}))), \quad (6)$$

where $\phi = (\mathbf{r}, \mathbf{s})$. The other is to introduce a hyperparameter β in the ELBO as in β -VAE (Higgins et al., 2017). The modified ELBO with β , up to the additive constant, is derived as:

$$\mathbb{E}_{\mathcal{D}}\{-\beta D_{\text{KL}}(q_\phi \| p_\Lambda) - \mathbb{E}_{\mathbf{z} \sim q_\phi}[(\mathbf{y} - \mathbf{f}_t(\mathbf{z}))^2 / 2\mathbf{g}_t(\mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_\phi} \log |\mathbf{g}_t(\mathbf{z})|\}. \quad (7)$$

For convenience, here and in \mathcal{L}_f in Sec. 4.2, we omit the summation as if Y is univariate. The encoder q_ϕ depends on t and can realize a PGS. With β , we control the trade-off between the first and second terms: the former is the divergence of the posterior from the balanced prior, and the latter is the reconstruction of the outcome. Note that a larger β encourages the conditional

balance $Z \perp\!\!\!\perp T | X$ on the posterior. By choosing β appropriately, e.g., by validation, the ELBO can recover an approximate bPGS while fitting the outcome well. In summary, we base the estimation on Proposition 2 and bPGS as much as possible, but step into Theorem 1 and noise modeling required by $p_e = p_\epsilon$ when necessary.

Note also that the parameters \mathbf{g} and \mathbf{k} , which model the outcome noise and express the uncertainty of the prior, respectively, are both learned by the ELBO. This deviates from the theoretical conditions described in Sec. 3.2, but it is more practical and yields better results in our experiments. See Sec. C.5 for more ideas and connections behind the ELBO.

Once the VAE is learned² by the ELBO, the estimate of the expected potential outcomes is given by:

$$\hat{\mu}_{\hat{t}}(\mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \mathbf{f}_{\hat{t}}(\mathbf{z}) = \mathbb{E}_{\mathcal{D}|\mathbf{x} \sim p(\mathbf{y}, t|\mathbf{x})} \mathbb{E}_{\mathbf{z} \sim q_\phi} \mathbf{f}_{\hat{t}}(\mathbf{z}), \quad \hat{t} \in \{0, 1\}, \quad (8)$$

where $q(\mathbf{z}|\mathbf{x}) := \mathbb{E}_{p(\mathbf{y}, t|\mathbf{x})} q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}, t)$ is the aggregated posterior. We mainly consider the case where \mathbf{x} is observed in the data, and the sample of (Y, T) is taken from the data given $X = \mathbf{x}$. When \mathbf{x} is not in the data, we replace q_ϕ with p_Λ in (8) (see Sec. C.7 for details and E for results). Note that \hat{t} in (8) indicates a counterfactual assignment that may not be the same as the factual $T = t$ in the data. That is, we set $T = \hat{t}$ in the decoder. The assignment is not applied to the encoder which is learned from factual X, Y, T (see also the explanation of $\epsilon_{CF,t}$ in Sec. 4.2). The overall **algorithm** steps are i) train the VAE using (7), and ii) infer CATE $\hat{\tau}(\mathbf{x}) = \hat{\mu}_1(\mathbf{x}) - \hat{\mu}_0(\mathbf{x})$ by (8).

4.2 CONDITIONALLY BALANCED REPRESENTATION LEARNING

We formally justify our ELBO (7) from the BRL viewpoint. We show that the conditional BRL via the KL (first) term of the ELBO results from bounding a CATE error; particularly, the error due to the imprecise recovery of \mathbf{j}_t in (**G1'**) is controlled by the ELBO. Previous works (Shalit et al., 2017; Lu et al., 2020) instead focus on unconditional balance and bound PEHE which is marginalized on X . Sec. 5.2 experimentally shows the advantage of our bounds and ELBO. Further, we connect the bounds to identification and consider noise modeling through $\mathbf{g}_t(\mathbf{z})$. Sec Sec. D.3 for detailed comparisons to previous works. In Sec. E.4, we empirically validate our bounds, and, particularly, the bounds are more useful under weaker overlap.

We introduce the objective that we bound. Using (8) to estimate CATE, $\hat{\tau}_f(\mathbf{z}) := \mathbf{f}_1(\mathbf{z}) - \mathbf{f}_0(\mathbf{z})$ is marginalized on $q(\mathbf{z}|\mathbf{x})$. On the other hand, the *true CATE*, given the covariate \mathbf{x} or score \mathbf{z} , is:

$$\tau(\mathbf{x}) = \mathbf{j}_1(\mathbf{p}_1(\mathbf{x})) - \mathbf{j}_0(\mathbf{p}_0(\mathbf{x})), \quad \tau_j(\mathbf{z}) = \mathbf{j}_1(\mathbf{z}) - \mathbf{j}_0(\mathbf{z}), \quad (9)$$

where \mathbf{j}_t is associated with an approximate bPGS \mathbf{p}_t (say, \mathbf{m}_t) as the target of recovery by our VAE. Accordingly, given \mathbf{x} , the *error of posterior CATE*, with or without knowing \mathbf{p}_t , is defined as

$$\epsilon_f^*(\mathbf{x}) := \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} (\hat{\tau}_f(\mathbf{z}) - \tau(\mathbf{x}))^2; \quad \epsilon_f(\mathbf{x}) := \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} (\hat{\tau}_f(\mathbf{z}) - \tau_j(\mathbf{z}))^2. \quad (10)$$

We bound ϵ_f instead of ϵ_f^* because the error between $\tau(X)$ and $\tau_j(Z)$ is small—if the score recovery works well, then $\mathbf{z} \approx \mathbf{p}_0(\mathbf{x}) \approx \mathbf{p}_1(\mathbf{x})$ in (9). We consider the error between $\hat{\tau}_f$ and τ_j below. We define the risks of outcome regression, into which ϵ_f is decomposed.

Definition 3 (CATE risks). Let $Y(\hat{t})|\mathbf{p}_{\hat{t}}(X) \sim p_{Y(\hat{t})|\mathbf{p}_{\hat{t}}}(\mathbf{y}|P)$ and $q_t(\mathbf{z}|\mathbf{x}) := q(\mathbf{z}|\mathbf{x}, t) = \mathbb{E}_{p(\mathbf{y}|\mathbf{x}, t)} q_\phi$. The *potential outcome loss* at (\mathbf{z}, t) , *factual risk*, and *counterfactual risk* are:

$$\begin{aligned} \mathcal{L}_f(\mathbf{z}, \hat{t}) &:= \mathbb{E}_{p_{Y(\hat{t})|\mathbf{p}_{\hat{t}}}(\mathbf{y}|P=\mathbf{z})} (\mathbf{y} - \mathbf{f}_{\hat{t}}(\mathbf{z}))^2 / \mathbf{g}_{\hat{t}}(\mathbf{z}) = \mathbf{g}_{\hat{t}}(\mathbf{z})^{-1} \int (\mathbf{y} - \mathbf{f}_{\hat{t}}(\mathbf{z}))^2 p_{Y(\hat{t})|\mathbf{p}_{\hat{t}}}(\mathbf{y}|\mathbf{z}) d\mathbf{y}; \\ \epsilon_{F,t}(\mathbf{x}) &:= \mathbb{E}_{q_t(\mathbf{z}|\mathbf{x})} \mathcal{L}_f(\mathbf{z}, t); \quad \epsilon_{CF,t}(\mathbf{x}) := \mathbb{E}_{q_{1-t}(\mathbf{z}|\mathbf{x})} \mathcal{L}_f(\mathbf{z}, t). \end{aligned}$$

With $Y(t)$ involved, \mathcal{L}_f is a potential outcome loss on \mathbf{f} , weighted by \mathbf{g} . The factual and counterfactual counterparts, $\epsilon_{F,t}$ and $\epsilon_{CF,t}$, are defined accordingly. In $\epsilon_{F,t}$, unit $\mathbf{u} = (\mathbf{x}, \mathbf{y}, t)$ is involved in the learning of $q_t(\mathbf{z}|\mathbf{x})$, as well as in $\mathcal{L}_f(\mathbf{z}, t)$ since $Y(t) = \mathbf{y}$ for the unit. In $\epsilon_{CF,t}$, however, unit $\mathbf{u}' = (\mathbf{x}, \mathbf{y}', 1-t)$ is involved in $q_{1-t}(\mathbf{z}|\mathbf{x})$, but not in $\mathcal{L}_f(\mathbf{z}, t)$ since $Y(t) \neq \mathbf{y}' = Y(1-t)$.

Thus, *the regression error (second) term in ELBO (7) controls $\epsilon_{F,t}$ via factual data*. On the other hand, $\epsilon_{CF,t}$ is not estimable due to the unobservable $Y(1-T)$, but is bounded by $\epsilon_{F,t}$ plus $MD(\mathbf{x})$ in Theorem 2 below—which, in turn, bounds ϵ_f by decomposing it to $\epsilon_{F,t}$, $\epsilon_{CF,t}$, and \mathbf{V}_Y .

²As usual, we expect the variational inference and optimization procedure to be (near) optimal; that is, consistency of VAE. *Consistent estimation* using the prior is a direct corollary of the consistent VAE. See Sec. C.6 for formal statements and proofs. Under Gaussian models, it is possible to prove the consistency of the posterior estimation, as shown in Bonhomme & Weidner (2021).

Theorem 2 (CATE error bound). Assume $|\mathcal{L}_f(\mathbf{z}, t)| \leq M$ and $|\mathbf{g}_t(\mathbf{z})| \leq G$, then:

$$\epsilon_f(\mathbf{x}) \leq 2[G(\epsilon_{F,0}(\mathbf{x}) + \epsilon_{F,1}(\mathbf{x}) + M\mathbf{D}(\mathbf{x})) - \mathbf{V}_Y(\mathbf{x})] \quad (11)$$

where $\mathbf{D}(\mathbf{x}) := \sum_t \sqrt{D_{\text{KL}}(q_t \| q_{1-t})}/2$, and $\mathbf{V}_Y(\mathbf{x}) := \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \sum_t \mathbb{E}_{p_{Y(t)|\mathbf{z}}(\mathbf{y}|\mathbf{z})} (\mathbf{y} - \mathbf{j}_t(\mathbf{z}))^2$.

$\mathbf{D}(\mathbf{x})$ measures the imbalance between $q_t(\mathbf{z}|\mathbf{x})$ and is symmetric for t . Correspondingly, the KL term in ELBO (7) is symmetric for t and balances $q_t(\mathbf{z}|\mathbf{x})$ by encouraging $Z \perp\!\!\!\perp T | X$ for the posterior. $\mathbf{V}_Y(\mathbf{x})$ reflects the intrinsic variance in the DGP and can not be controlled. Estimating G, M is nontrivial. Instead, we rely on β in the ELBO (7) to weight the terms. We do not need two hyperparameters since G is implicitly controlled by the third term, a norm constraint, in ELBO.

5 EXPERIMENTS

We compare our method with existing methods on three types of datasets. Here, we present two experiments; the remaining one on the Pokec dataset is deferred to Sec. E.3. As in previous works (Shalit et al., 2017; Louizos et al., 2017), we report the absolute error of ATE $\epsilon_{ate} := |\mathbb{E}_{\mathcal{D}}(y(1) - y(0)) - \mathbb{E}_{\mathcal{D}}\hat{\tau}(\mathbf{x})|$ and, as a surrogate of square CATE error $\epsilon_{cate}(\mathbf{x}) = \mathbb{E}_{\mathcal{D}|\mathbf{x}}[(y(1) - y(0)) - \hat{\tau}(\mathbf{x})]^2$, the empirical PEHE $\epsilon_{pehe} := \mathbb{E}_{\mathcal{D}}\epsilon_{cate}(\mathbf{x})$ (Hill, 2011), which is the average square CATE error.

Unless otherwise indicated, for each function f, g, h, k, r, s in ELBO (7), we use a multilayer perceptron, with $200 * 3$ hidden units (width 200, 3 layers), and ELU activations (Clevert et al., 2015). $\Lambda = (h, k)$ depends only on X . The Adam optimizer with initial learning rate 10^{-4} and batch size 100 is employed. All experiments use early-stopping of training by evaluating the ELBO on a validation set. More details on hyper-parameters and settings are given in each experiment.

5.1 SYNTHETIC DATASET

$$W|X \sim \mathcal{N}(\mathbf{h}(X), \mathbf{k}(X)); T|X \sim \text{Bern}(\text{Logi}(\omega l(X))); Y|W, T \sim \mathcal{N}(f_T(W), g_T(W)). \quad (12)$$

We generate synthetic datasets following (12). Both $X \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma})$ and W are factorized Gaussians. $\boldsymbol{\mu}, \boldsymbol{\sigma}$ are randomly sampled. The functions h, k, l are linear. Outcome models f_0, f_1 are built by NNs with invertible activations. Y is univariate, $\dim(X) = 30$, and $\dim(W)$ ranges from 1 to 5. W is a bPGS, but the dimensionality is not low enough to satisfy the injectivity in (G2), when $\dim(W) > 1$. We have 5 different overlap levels controlled by ω that multiplies the logit value. See Sec. E.1 for details and more results on synthetic datasets.

With the same $(\dim(W), \omega)$, we evaluate our method and CFR on 10 random DGPs, with different sets of functions f, g, h, k, l in (12). For each DGP, we sample 1500 data points, and split them into 3 equal sets for training, validation, and testing. We show our results for different hyperparameter β . For CFR, we try different balancing parameters and present the best results (see the Appendix for detail).

In each panel of Figure 2, we adjust one of $\omega, \dim(W)$, with the other fixed to the lowest. As implied by our theory, our method, with only 1-dimensional Z , performs much better in the left panel (where $\dim(W) = 1$ satisfies (G2)) than in the right panel (when $\dim(W) > 1$). Although CFR uses 200-dimensional representation, in the left panel our method performs much better than CFR; moreover, in the right panel CFR is not much better than ours. Further, our method is much more robust against different DGPs than CFR (see the error bars). Thus, the results indicate the power of identification and recovery of scores. (see Figure 3 also).

Under the lowest overlap level ($\omega = 22$), large $\beta (= 2.5, 3)$ shows the best results, which accords with the intuition and bounds in Sec. 4. When $\dim(W) > 1$, f_t in (12) is non-injective and learning

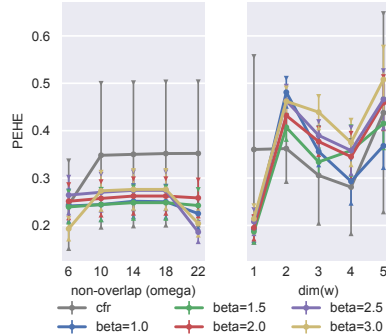


Figure 2: $\sqrt{\epsilon_{pehe}}$ on synthetic datasets. Each error bar is on 10 random DGPs.

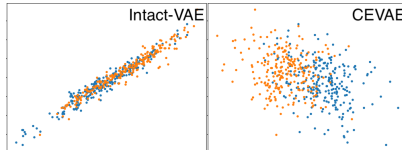


Figure 3: Plots of recovered - true latent. Blue: $T = 0$, Orange: $T = 1$.

of PGS is necessary, and thus, larger β has a negative effect. In fact, $\beta = 1$ is significantly better than $\beta = 3$ when $\dim(W) > 2$. We note that our method, with a higher-dimensional Z , outperforms or matches CFR also under $\dim(W) > 1$ (see Appendix Figure 5). Thus, the performance gap under $\dim(W) > 1$ in Figure 2 should be due to the capacity of NNs in β -Intact-VAE. In Appendix Figure 7 for ATE error, CFR drops performance w.r.t overlap levels. This is evidence that CFR and its unconditional balance overly focus on PEHE (see Sec. 5.2 for more explicit comparison).

When $\dim(W) = 1$, there are no better PSs than W , because f_t is invertible and no information can be dropped from W . Thus, our method stably learns Z as an approximate affine transformation of the true W , showing identification. An example is shown in Figure 3, and more plots are in Appendix Figure 9. For comparison, we run CEVAE, which is also based on VAE but without identification; CEVAE shows much lower quality of recovery. As expected, both recovery and estimation are better with the balanced prior $p_\Lambda(z|x)$, and we can see examples of bad recovery using $p_\Lambda(z|x, t)$ in Appendix Figure 10.

5.2 IHDP BENCHMARK DATASET

This experiment shows our conditional BRL matches state-of-the-art BRL methods and does not overly focus on PEHE. The IHDP (Hill, 2011) is a widely used benchmark dataset; while it is less known, its covariates are limited-overlapping, and thus it is used in Johansson et al. (2020) which considers limited overlap. The dataset is based on an RCT, but `Race` is artificially introduced as a confounder by removing all treated babies with nonwhite mothers in the data. Thus, `Race` is highly limited-overlapping, and other covariates that have high correlation to `Race`, e.g. `Birth weight` (Kelly et al., 2009), are also limited-overlapping. See Sec. E.2 for detail and more results.

There is a linear bPGS (linear combination of the covariates). However, most of the covariates are binary, so the support of the bPGS is often on small and separated intervals. Thus, the Gaussian latent Z in our model is misspecified. We use higher-dimensional Z to address this, similar to Louizos et al. (2017). Specifically, we set $\dim(Z) = 50$, together with NNs of $50 * 2$ hidden units in the prior and encoder. We set $\beta = 1$ since it works well on synthetic datasets with limited overlap.

As shown in Table 1, β -Intact-VAE outperforms or matches the state-of-the-art methods; it has the best performance measured by both ϵ_{ate} and ϵ_{pehe} and matches CF and CFR respectively. Also notably, our method outperforms other generative models (CEVAE and GANITE) by large margins.

To show our conditional balance is preferable, we also modify our method and add two components for *unconditional* balance from CFR (see the Appendix), which is based on bounding PEHE and is controlled by another hyperparameter γ . In the modified version, the over-focus on PEHE of the unconditional balance is seen clearly—with different γ , it significantly affects PEHE, but barely affects ATE error. In fact, the unconditional balance, with larger γ , only worsens the performance. See also Appendix Figure 7 where CFR gives larger ATE errors with less overlap.

Table 1: Errors on IHDP over 1000 random DGPs. “Mod. *” indicates the modified version with unconditional balance of strength $\gamma = *$. *Italic* indicates where the modified version is significantly worse than the original. **Bold** indicates method(s) which is significantly better than others. The results of other methods are taken from Shalit et al. (2017), except for GANITE and CEVAE, the results of which are taken from original works.

Method	TMLE	BNN	CFR	CF	CEVAE	GANITE	Ours	Mod. 1	Mod. 0.2	Mod. 0.1	Mod. 0.05	Mod. 0.01
ϵ_{ate}	.30 \pm .01	.37 \pm .03	.25 \pm .01	.18 \pm .01	.34 \pm .01	.43 \pm .05	.180 \pm .007	.185 \pm .008	.185 \pm .008	.186 \pm .009	.183 \pm .008	.181 \pm .008
$\sqrt{\epsilon_{pehe}}$	5.0 \pm .2	2.2 \pm .1	.71 \pm .02	3.8 \pm .2	2.7 \pm .1	1.9 \pm .4	.709 \pm .024	1.175 \pm .046	.797 \pm .030	.748 \pm .028	.732 \pm .028	.719 \pm .027

6 CONCLUSION

We proposed a method for CATE estimation under limited overlap. Our method exploits identifiable VAE, a recent advance in generative models, and is fully motivated and theoretically justified by causal considerations: identification, prognostic score, and balance. Experiments show evidence that the injectivity of f_t in our model is possibly unnecessary because $\dim(Z) > \dim(Y)$ yields better results. A theoretical study of this is an interesting future direction. We have evidence that Intact-VAE works under unobserved confounding and believe that VAEs are suitable for *principled* causal inference owing to their probabilistic nature, if not compromised by ad hoc heuristics (Wu & Fukumizu, 2021).

REFERENCES

- Jason Abrevaya, Yu-Chin Hsu, and Robert P Lieli. Estimating conditional average treatment effects. *Journal of Business & Economic Statistics*, 33(4):485–505, 2015.
- Ahmed M Alaa and Mihaela van der Schaar. Bayesian inference of individualized treatment effects using multi-task gaussian processes. In *Advances in Neural Information Processing Systems*, pp. 3424–3432, 2017.
- Timothy B Armstrong and Michal Kolesár. Finite-sample optimal estimation and inference on average treatment effects under unconfoundedness. *Econometrica*, 89(3):1141–1177, 2021.
- Stéphane Bonhomme and Martin Weidner. Posterior average effects. *Journal of Business & Economic Statistics*, (just-accepted):1–38, 2021.
- Victor Chernozhukov and Christian Hansen. Quantile models with endogeneity. *Annu. Rev. Econ.*, 5(1):57–81, 2013.
- Denis Chetverikov and Daniel Wilhelm. Nonparametric instrumental variable estimation under monotonicity. *Econometrica*, 85(4):1303–1320, 2017.
- Denis Chetverikov, Andres Santos, and Azeem M Shaikh. The econometrics of shape restrictions. *Annual Review of Economics*, 10:31–63, 2018.
- Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015.
- Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in neural information processing systems*, pp. 2292–2300, 2013.
- Wangzhi Dai and Collin M Stultz. Quantifying common support between multiple treatment groups using a contrastive-vae. In *Machine Learning for Health*, pp. 41–52. PMLR, 2020.
- Alexander D’Amour and Alexander Franks. Deconfounding scores: Feature representations for causal effect estimation with weak overlap. *arXiv preprint arXiv:2104.05762*, 2021.
- Carl Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016.
- Alexander D’Amour, Peng Ding, Avi Feller, Lihua Lei, and Jasjeet Sekhon. Overlap in observational studies with high-dimensional covariates. *Journal of Econometrics*, 2020.
- Max H Farrell. Robust inference on average treatment effects with possibly more covariates than observations. *Journal of Econometrics*, 189(1):1–23, 2015.
- Joachim Freyberger and Joel L Horowitz. Identification and shape restrictions in nonparametric instrumental variables estimation. *Journal of Econometrics*, 189(1):41–53, 2015.
- Li Gan and Qi Li. Efficiency of thin and thick markets. *Journal of Econometrics*, 192(1):40–54, 2016.
- Prem K Gopalan and David M Blei. Efficient discovery of overlapping communities in massive networks. *Proceedings of the National Academy of Sciences*, 110(36):14534–14539, 2013.
- Sander Greenland. The effect of misclassification in the presence of covariates. *American journal of epidemiology*, 112(4):564–569, 1980.
- David Hajage, Yann De Rycke, Guillaume Chauvet, and Florence Tubach. Estimation of conditional and marginal odds ratios using the prognostic score. *Statistics in medicine*, 36(4):687–716, 2017.
- Ben B Hansen. The prognostic analogue of the propensity score. *Biometrika*, 95(2):481–488, 2008.
- Negar Hassanpour and Russell Greiner. Learning disentangled representations for counterfactual regression. In *International Conference on Learning Representations*, 2019.
- Miguel A. Hernan and James M. Robins. *Causal Inference: What If*. CRC Press, 1st edition, 2020. ISBN 978-1-4200-7616-5.

- Irina Higgins, Loïc Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *5th International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=Sy2fzU9gl>.
- Jennifer L Hill. Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, 20(1):217–240, 2011.
- Han Hong, Michael P Leung, and Jessie Li. Inference on finite-population treatment effects under limited overlap. *The Econometrics Journal*, 23(1):32–47, 2020.
- Ming-Yueh Huang and Kwun Chuen Gary Chan. Joint sufficient dimension reduction and estimation of conditional and average treatment effects. *Biometrika*, 104(3):583–596, 2017.
- Martin Huber and Kaspar Wüthrich. Local average and quantile treatment effects under endogeneity: a review. *Journal of Econometric Methods*, 8(1), 2018.
- Guido W Imbens and Donald B Rubin. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press, 2015.
- Dominik Janzing and Bernhard Scholkopf. Causal inference using the algorithmic markov condition. *IEEE Transactions on Information Theory*, 56(10):5168–5194, 2010.
- Andrew Jesson, Sören Mindermann, Uri Shalit, and Yarin Gal. Identifying causal-effect inference failure with uncertainty-aware models. *Advances in Neural Information Processing Systems*, 33, 2020.
- Fredrik Johansson, Uri Shalit, and David Sontag. Learning representations for counterfactual inference. In *International conference on machine learning*, pp. 3020–3029, 2016.
- Fredrik D Johansson, David Sontag, and Rajesh Ranganath. Support and invertibility in domain-invariant representations. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 527–536. PMLR, 2019.
- Fredrik D Johansson, Uri Shalit, Nathan Kallus, and David Sontag. Generalization bounds and representation learning for estimation of potential outcomes and causal effects. *arXiv preprint arXiv:2001.07426*, 2020.
- Nathan Kallus, Brenton Pennicooke, and Michele Santacatterina. More robust estimation of sample average treatment effects using kernel optimal matching in an observational study of spine surgical interventions. *arXiv preprint arXiv:1811.04274*, 2018.
- Yvonne Kelly, Lidia Panico, Mel Bartley, Michael Marmot, James Nazroo, and Amanda Sacker. Why does birthweight vary among ethnic groups in the uk? findings from the millennium cohort study. *Journal of public health*, 31(1):131–137, 2009.
- Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational autoencoders and nonlinear ica: A unifying framework. In *International Conference on Artificial Intelligence and Statistics*, pp. 2207–2217, 2020a.
- Ilyes Khemakhem, Ricardo Monti, Diederik Kingma, and Aapo Hyvarinen. Ice-beem: Identifiable conditional energy-based deep models based on nonlinear ica. *Advances in Neural Information Processing Systems*, 33, 2020b.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. URL <http://arxiv.org/abs/1312.6114>.
- Diederik P Kingma, Max Welling, et al. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019.
- Durk P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised learning with deep generative models. In *Advances in neural information processing systems*, pp. 3581–3589, 2014.

- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *5th International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=SJU4ayYgl>.
- Manabu Kuroki and Judea Pearl. Measurement bias and effect restoration in causal inference. *Biometrika*, 101(2):423–437, 2014.
- Jure Leskovec and Andrej Krevl. Snap datasets: Stanford large network dataset collection, 2014.
- Arthur Lewbel. The identification zoo: Meanings of identification in econometrics. *Journal of Economic Literature*, 57(4):835–903, 2019.
- Fan Li and Fan Li. Propensity score weighting for causal inference with multiple treatments. *The Annals of Applied Statistics*, 13(4):2389–2415, 2019.
- Zheng Li, Guannan Liu, and Qi Li. Nonparametric knn estimation with monotone constraints. *Econometric Reviews*, 36(6-9):988–1006, 2017.
- Christos Louizos, Uri Shalit, Joris M Mooij, David Sontag, Richard Zemel, and Max Welling. Causal effect inference with deep latent-variable models. In *Advances in Neural Information Processing Systems*, pp. 6446–6456, 2017.
- Danni Lu, Chenyang Tao, Junya Chen, Fan Li, Feng Guo, and Lawrence Carin. Reconsidering generative objectives for counterfactual reasoning. *Advances in Neural Information Processing Systems*, 33, 2020.
- Wei Luo, Yeying Zhu, and Debashis Ghosh. On estimating regression-based causal effects using sufficient dimension reduction. *Biometrika*, 104(1):51–65, 2017.
- Emile Mathieu, Tom Rainforth, Nana Siddharth, and Yee Whye Teh. Disentangling disentanglement in variational autoencoders. In *International Conference on Machine Learning*, pp. 4402–4412. PMLR, 2019.
- Wang Miao, Zhi Geng, and Eric J Tchetgen Tchetgen. Identifying causal effects with proxy variables of an unmeasured confounder. *Biometrika*, 105(4):987–993, 2018.
- Michael Oberst, Fredrik Johansson, Dennis Wei, Tian Gao, Gabriel Brat, David Sontag, and Kush Varshney. Characterization of overlap in observational studies. In *International Conference on Artificial Intelligence and Statistics*, pp. 788–798. PMLR, 2020.
- Judea Pearl. *Causality: models, reasoning and inference*. Cambridge University Press, 2009.
- Severi Rissanen and Pekka Marttinen. A critical look at the identifiability of causal effects with deep latent variable models. *NeurIPS 2021, to appear*, 2021.
- Paul R Rosenbaum. Modern algorithms for matching in observational studies. *Annual Review of Statistics and Its Application*, 7:143–176, 2020.
- 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.
- Claudia Shi, David Blei, and Victor Veitch. Adapting neural networks for the estimation of treatment effects. In *Advances in Neural Information Processing Systems*, pp. 2507–2517, 2019.
- Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. In *Advances in neural information processing systems*, pp. 3483–3491, 2015.

- Peter Sorrenson, Carsten Rother, and Ullrich Köthe. Disentanglement by nonlinear ica with general incompressible-flow networks (gin). In *International Conference on Learning Representations*, 2019.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- Jennifer E Starling, Catherine E Aiken, Jared S Murray, Annetee Nakimuli, and James G Scott. Monotone function estimation in the presence of extreme data coarsening: Analysis of preeclampsia and birth weight in urban uganda. *arXiv preprint arXiv:1912.06946*, 2019.
- Elizabeth A. Stuart. Matching Methods for Causal Inference: A Review and a Look Forward. *Statistical Science*, 25(1):1 – 21, 2010. doi: 10.1214/09-STS313. URL <https://doi.org/10.1214/09-STS313>.
- Xinwei Sun, Botong Wu, Chang Liu, Xiangyu Zheng, Wei Chen, Tao Qin, and Tie-yan Liu. Latent causal invariant model. *arXiv preprint arXiv:2011.02203*, 2020.
- Alexander Tarr and Kosuke Imai. Estimating average treatment effects with support vector machines. *arXiv preprint arXiv:2102.11926*, 2021.
- Mark J Van der Laan and Sherri Rose. *Targeted learning in data science: causal inference for complex longitudinal studies*. Springer, 2018.
- Victor Veitch, Yixin Wang, and David Blei. Using embeddings to correct for unobserved confounding in networks. In *Advances in Neural Information Processing Systems*, pp. 13792–13802, 2019.
- 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.
- Shanshan Wang, Liren Yang, Li Shang, Wenfang Yang, Cuifang Qi, Liyan Huang, Guilan Xie, Ruiqi Wang, and Mei Chun Chung. Changing trends of birth weight with maternal age: a cross-sectional study in xi’an city of northwestern china. *BMC Pregnancy and Childbirth*, 20(1):1–8, 2020.
- Halbert White and Karim Chalak. Identification and identification failure for treatment effects using structural systems. *Econometric Reviews*, 32(3):273–317, 2013.
- Pengzhou Wu and Kenji Fukumizu. Causal mosaic: Cause-effect inference via nonlinear ica and ensemble method. In *International Conference on Artificial Intelligence and Statistics*, pp. 1157–1167. PMLR, 2020a. URL <http://proceedings.mlr.press/v108/wu20b.html>.
- Pengzhou Wu and Kenji Fukumizu. Towards principled causal effect estimation by deep identifiable models. *arXiv preprint arXiv:2109.15062*, 2021.
- Pengzhou Abel Wu and Kenji Fukumizu. Identifying treatment effects under unobserved confounding by causal representation learning. *submitted to ICLR 2021*, 2020b. URL <https://openreview.net/forum?id=D3TNqCspFpM>.
- S Yang and P Ding. Asymptotic inference of causal effects with observational studies trimmed by the estimated propensity scores. *Biometrika*, 105(2):487–493, 03 2018. ISSN 0006-3444. doi: 10.1093/biomet/asy008. URL <https://doi.org/10.1093/biomet/asy008>.
- Liuyi Yao, Sheng Li, Yaliang Li, Mengdi Huai, Jing Gao, and Aidong Zhang. Representation learning for treatment effect estimation from observational data. In *Advances in Neural Information Processing Systems*, pp. 2633–2643, 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. URL <https://openreview.net/forum?id=ByKWUeWA->.
- Weijia Zhang, Lin Liu, and Jiuyong Li. Treatment effect estimation with disentangled latent factors. *arXiv preprint arXiv:2001.10652*, 2020a.
- 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, 2020b.

A PROOFS

We restate our model identifiability formally.

Lemma 1 (Model identifiability). *Given model (3) under (M1), for $T = t$, assume*

(D1') (Non-degenerated data for λ) there exist $2n + 1$ points $\mathbf{x}_0, \dots, \mathbf{x}_{2n} \in \mathcal{X}$ such that the $2n$ -square matrix $\mathbf{L}_t := [\gamma_{t,1}, \dots, \gamma_{t,2n}]$ is invertible, where $\gamma_{t,k} := \lambda_t(\mathbf{x}_k) - \lambda_t(\mathbf{x}_0)$.

Then, given $T = t$, the family is identifiable up to an equivalence class. That is, if $p_{\theta}(\mathbf{y}|\mathbf{x}, t) = p_{\theta'}(\mathbf{y}|\mathbf{x}, t)$, we have the relation between parameters: for any \mathbf{y}_t in the image of \mathbf{f}_t ,

$$\mathbf{f}_t^{-1}(\mathbf{y}_t) = \text{diag}(\mathbf{a})\mathbf{f}_t'^{-1}(\mathbf{y}_t) + \mathbf{b} =: \mathcal{A}_t(\mathbf{f}_t'^{-1}(\mathbf{y}_t)) \quad (13)$$

where $\text{diag}(\mathbf{a})$ is an invertible n -diagonal matrix and \mathbf{b} is a n -vector, both depend on λ_t and λ'_t .

Note, **(D1)** in the main text implies **(D1')**, see Sec. B.2.3 in Khemakhem et al. (2020a). The main part of our model identifiability is essentially the same as that of Theorem 1 in Khemakhem et al. (2020a), but now adapted to include the dependency on t . Here we give an outline of the proof, and the details can be easily filled by referring to Khemakhem et al. (2020a). In the proof, subscripts t are omitted for convenience.

Proof of Lemma 1. Using **(M1)** i) and ii), we transform $p_{\mathbf{f},\lambda}(\mathbf{y}|\mathbf{x}, t) = p_{\mathbf{f}',\lambda'}(\mathbf{y}|\mathbf{x}, t)$ into equality of noiseless distributions, that is,

$$q_{\mathbf{f}',\lambda'}(\mathbf{y}) = q_{\mathbf{f},\lambda}(\mathbf{y}) := p_{\lambda}(\mathbf{f}^{-1}(\mathbf{y})|\mathbf{x}, t)\text{vol}(\mathbf{J}_{\mathbf{f}^{-1}}(\mathbf{y}))\mathbb{I}_{\mathcal{Y}}(\mathbf{y}) \quad (14)$$

where p_{λ} is the Gaussian density function of the conditional prior defined in (3) and $\text{vol}(A) := \sqrt{\det AA^T}$. $q_{\mathbf{f}',\lambda'}$ is defined similarly to $q_{\mathbf{f},\lambda}$.

Then, apply model (3) to (14), plug the $2n + 1$ points from **(D1')** into it, and re-arrange the resulting $2n + 1$ equations in matrix form, we have

$$\mathcal{F}'(Y) = \mathcal{F}(Y) := \mathbf{L}^T \mathbf{t}(\mathbf{f}^{-1}(Y)) - \beta \quad (15)$$

where $\mathbf{t}(Z) := (Z, Z^2)^T$ is the sufficient statistics of factorized Gaussian, and $\beta_t := (\alpha_t(\mathbf{x}_1) - \alpha_t(\mathbf{x}_0), \dots, \alpha_t(\mathbf{x}_{2n}) - \alpha_t(\mathbf{x}_0))^T$ where $\alpha_t(X; \lambda_t)$ is the log-partition function of the conditional prior in (3). \mathcal{F}' is defined similarly to \mathcal{F} , but with $\mathbf{f}', \lambda', \alpha'$

Since \mathbf{L} is invertible, we have

$$\mathbf{t}(\mathbf{f}^{-1}(Y)) = \mathbf{A}\mathbf{t}(\mathbf{f}'^{-1}(Y)) + \mathbf{c} \quad (16)$$

where $\mathbf{A} = \mathbf{L}^{-T}\mathbf{L}'^T$ and $\mathbf{c} = \mathbf{L}^{-T}(\beta - \beta')$.

The final part of the proof is to show, by following the same reasoning as in Appendix B of Sorrenson et al. (2019), that \mathbf{A} is a sparse matrix such that

$$\mathbf{A} = \begin{pmatrix} \text{diag}(\mathbf{a}) & \mathbf{O} \\ \text{diag}(\mathbf{u}) & \text{diag}(\mathbf{a}^2) \end{pmatrix} \quad (17)$$

where \mathbf{A} is partitioned into four n -square matrices. Thus

$$\mathbf{f}^{-1}(Y) = \text{diag}(\mathbf{a})\mathbf{f}'^{-1}(Y) + \mathbf{b} \quad (18)$$

where \mathbf{b} is the first half of \mathbf{c} . \square

Proof of Proposition 2. Under **(G2)**, and **(M3)**, we have

$$\mathbb{E}_{p_{\theta}}(Y|X, T) = \mathbb{E}(Y|X, T) \implies \mathbf{f}_t \circ \mathbf{h}(\mathbf{x}) = \mathbf{j}_t \circ \mathbf{p}(\mathbf{x}) \text{ on } (\mathbf{x}, t) \text{ such that } p(t, \mathbf{x}) > 0. \quad (19)$$

We show the solution set of (19) on *overlapping* \mathbf{x} is

$$\{(\mathbf{f}, \mathbf{h}) | \mathbf{f}_t = \mathbf{j}_t \circ \Delta^{-1}, \mathbf{h} = \Delta \circ \mathbf{p}, \Delta : \mathcal{P} \rightarrow \mathbb{R}^n \text{ is injective}\}. \quad (20)$$

By **(G2)(M1)**, and with injective $\mathbf{f}_t, \mathbf{j}_t$ and $\dim(Z) = \dim(Y) \geq \dim(\mathbf{p})$, for any Δ above, there exists a functional parameter \mathbf{f}_t such that $\mathbf{j}_t = \mathbf{f}_t \circ \Delta$. Thus, set (20) is non-empty, and any element is indeed a solution because $\mathbf{f}_t \circ \mathbf{h} = \mathbf{j}_t \circ \Delta^{-1} \circ \Delta \circ \mathbf{p} = \mathbf{j}_t \circ \mathbf{p}$.

Any solution of (19) should be in (20). A solution should satisfy $\mathbf{h}(\mathbf{x}) = \mathbf{f}_t^{-1} \circ \mathbf{j}_t \circ \mathbf{p}(\mathbf{x})$ for both t since \mathbf{x} is overlapping. This means the *injective* function $\mathbf{f}_t^{-1} \circ \mathbf{j}_t$ should *not* depend on t , thus it is one of the Δ in (20).

We proved conclusion 1) with $\mathbf{v} := \Delta$. And, on overlapping \mathbf{x} , conclusion 2) is quickly seen from

$$\hat{\mu}_t(\mathbf{x}) = \mathbf{f}_t(\mathbf{h}(\mathbf{x})) = \mathbf{j}_t \circ \mathbf{v}^{-1}(\mathbf{v} \circ \mathbf{p}(\mathbf{x})) = \mathbf{j}_t(\mathbf{p}(\mathbf{x})) = \mu_t(\mathbf{x}). \quad (21)$$

We rely on overlapping \mathbf{p} to work for non-overlapping \mathbf{x} . For any \mathbf{x}_t with $p(1-t|\mathbf{x}_t) = 0$, to ensure $p(1-t|\mathbf{p}(\mathbf{x}_t)) > 0$, there should exist \mathbf{x}_{1-t} such that $\mathbf{p}(\mathbf{x}_{1-t}) = \mathbf{p}(\mathbf{x}_t)$ and $p(1-t|\mathbf{x}_{1-t}) > 0$. And we also have $\mathbf{h}(\mathbf{x}_{1-t}) = \mathbf{h}(\mathbf{x}_t)$ due to **(M2)**. Then, we have

$$\hat{\mu}_{1-t}(\mathbf{x}_t) = \mathbf{f}_{1-t}(\mathbf{h}(\mathbf{x}_t)) = \mathbf{f}_{1-t}(\mathbf{h}(\mathbf{x}_{1-t})) = \mathbf{j}_{1-t}(\mathbf{p}(\mathbf{x}_{1-t})) = \mathbf{j}_{1-t}(\mathbf{p}(\mathbf{x}_t)) = \mu_{1-t}(\mathbf{x}_t). \quad (22)$$

The third equality uses (19) on $(\mathbf{x}_{1-t}, 1-t)$. \square

Below we prove Theorem 1 with **(D2)** replaced by

(D2') (Spontaneous balance) there exist $2n+1$ points $\mathbf{x}_0, \dots, \mathbf{x}_{2n} \in \mathcal{X}$, $2n$ -square matrix \mathbf{C} , and $2n$ -vector \mathbf{d} , such that $\mathbf{L}_0^{-1}\mathbf{L}_1 = \mathbf{C}$ and $\beta_0 - \mathbf{C}^{-T}\beta_1 = \mathbf{d}/k$ for optimal λ_t (see below), where \mathbf{L}_t is defined in **(D1')**, $\beta_t := (\alpha_t(\mathbf{x}_1) - \alpha_t(\mathbf{x}_0), \dots, \alpha_t(\mathbf{x}_{2n}) - \alpha_t(\mathbf{x}_0))^T$, and $\alpha_t(X; \lambda_t)$ is the log-partition function of the prior in (3).

(D2') restricts the discrepancy between λ_0, λ_1 on $2n+1$ values of X , thus is relatively easy to satisfy with high-dimensional X . **(D2')** is general despite (or thanks to) the involved formulation. Let us see its generality even under a highly special case: $\mathbf{C} = c\mathbf{I}$ and $\mathbf{d} = \mathbf{0}$. Then, $\mathbf{L}_0^{-1}\mathbf{L}_1 = c\mathbf{I}$ requires that, $\mathbf{h}_1(\mathbf{x}_k) - c\mathbf{h}_0(\mathbf{x}_k)$ is the same for $2n+1$ points \mathbf{x}_k . This is easily satisfied except for $n \gg m$ where m is the dimension of X , which rarely happens in practice. And, $\beta_0 - \mathbf{C}^{-T}\beta_1 = \mathbf{d}$ becomes just $\beta_1 = c\beta_0$. This is equivalent to $\alpha_1(\mathbf{x}_k) - c\alpha_0(\mathbf{x}_k)$ same for $2n+1$ points, again fine in practice. However, the high generality comes with price. Verifying **(D2')** using data is challenging, particularly with high-dimensional covariate and latent variable. Although we believe fast algorithms for this purpose could be developed, the effort would be nontrivial. This is another motivation to use the extreme case $\lambda_0 = \lambda_1$ in Sec. 4.1, which corresponds to $\mathbf{C} = \mathbf{I}$ and $\mathbf{d} = \mathbf{0}$.

Proof of Theorem 1. By **(M1)** and **(G1')**, for any injective function $\Delta : \mathcal{P} \rightarrow \mathbb{R}^n$, there exists a functional parameter \mathbf{f}_t^* such that $\mathbf{j}_t = \mathbf{f}_t^* \circ \Delta$. Let $\mathbf{h}_t^* = \Delta \circ \mathbf{p}_t$, then, clearly from **(M3')**, such parameters $\theta^* = (\mathbf{f}^*, \mathbf{h}^*)$ are optimal: $p_{\theta^*}(\mathbf{y}|\mathbf{x}, t) = p(\mathbf{y}|\mathbf{x}, t)$.

Since have all assumptions for Lemma 1, we have

$$\Delta \circ \mathbf{j}^{-1}(\mathbf{y}) = \mathbf{f}^* \circ \mathbf{p}^{-1}(\mathbf{y}) = \mathcal{A} \circ \mathbf{f}^{-1}(\mathbf{y})|_t, \text{ on } (\mathbf{y}, t) \in \{(\mathbf{j}_t \circ \mathbf{p}_t(\mathbf{x}), t) | p(t, \mathbf{x}) > 0\}, \quad (23)$$

where \mathbf{f} is any optimal parameter, and “ $|_t$ ” collects all subscripts t . Note, except for Δ , all the symbols should have subscript t .

Nevertheless, using **(D2')**, we can further prove $\mathcal{A}_0 = \mathcal{A}_1$.

We repeat the core quantities from Lemma 1 here: $\mathbf{A}_t = \mathbf{L}_t^{-T}\mathbf{L}'_t$ and $\mathbf{c}_t = \mathbf{L}_t^{-T}(\beta_t - \beta'_t)$.

From **(D2')**, we immediately have

$$\mathbf{L}_0^{-1}\mathbf{L}_1 = \mathbf{L}'_0^{-1}\mathbf{L}'_1 = \mathbf{C} \iff \mathbf{A}_0 = \mathbf{A}_1 \quad (24)$$

And also,

$$\begin{aligned} \mathbf{L}_0^{-1}\mathbf{L}_1 = \mathbf{C} &\iff \mathbf{L}_0^{-T}\mathbf{C}^{-T} = \mathbf{L}_1^{-T} \\ \beta_0 - \mathbf{C}^{-T}\beta_1 = \beta'_0 - \mathbf{C}^{-T}\beta'_1 = \mathbf{d}/k &\iff \mathbf{C}^T(\beta_0 - \beta'_0) = \beta_1 - \beta'_1 \end{aligned} \quad (25)$$

Multiply right hand sides of the two lines, we have $\mathbf{c}_0 = \mathbf{c}_1$. Now we have $\mathcal{A}_0 = \mathcal{A}_1 := \mathcal{A}$. Apply this to (23), we have

$$\mathbf{f}_t = \mathbf{j}_t \circ \mathbf{v}^{-1}, \quad \mathbf{v} := \mathcal{A}^{-1} \circ \Delta \quad (26)$$

for any optimal parameters $\theta = (\mathbf{f}, \mathbf{h})$. Again, from **(M3')**, we have

$$p_{\theta}(\mathbf{y}|\mathbf{x}, t) = p(\mathbf{y}|\mathbf{x}, t) \implies p_{\epsilon}(\mathbf{y} - \mathbf{f}_t(\mathbf{h}_t(\mathbf{x}))) = p_{\epsilon}(\mathbf{y} - \mathbf{j}_t(\mathbf{p}_t(\mathbf{x}))) \quad (27)$$

where $p_{\epsilon} = p_{\epsilon}$. And the above is only possible when $\mathbf{f}_t \circ \mathbf{h}_t = \mathbf{j}_t \circ \mathbf{p}_t$. Combined with $\mathbf{f}_t = \mathbf{j}_t \circ \mathbf{v}^{-1}$, we have conclusion 1).

And conclusion 2) follows from the same reasoning as Proposition 2, applied to both \mathbf{p}_0 and \mathbf{p}_1 . \square

Note, when multiplying the two lines of (25), the effects of $k \rightarrow 0$ cancel out, and \mathbf{c}_t is finite and well-defined. Also, it is apparent from above proof that **(D2')** is a necessary and sufficient condition for $\mathcal{A}_0 = \mathcal{A}_1$, if other conditions of Theorem 1 are given.

Below, we prove the results in Sec. 4.2. The definitions and results work for the prior; simply replace $q_t(\mathbf{x}|\mathbf{x})$ with $p_t(\mathbf{z}|\mathbf{x}) := p_\lambda(\mathbf{z}|\mathbf{x}, t)$ in definitions and statements, and the proofs below hold as the same. The dependence on \mathbf{f} prevail, and the superscripts are omitted. The arguments \mathbf{x} are sometimes also omitted.

Lemma 2 (Counterfactual risk bound). *Assume $|\mathcal{L}_f(\mathbf{z}, t)| \leq M$, we have*

$$\epsilon_{CF}(\mathbf{x}) \leq \sum_t q(1-t|\mathbf{x})\epsilon_{F,t}(\mathbf{x}) + M\mathbf{D}(\mathbf{x}) \quad (28)$$

where $\epsilon_{CF}(\mathbf{x}) := \sum_t p(1-t|\mathbf{x})\epsilon_{CF,t}(\mathbf{x})$, and $\mathbf{D}(\mathbf{x}) := \sum_t \sqrt{D_{\text{KL}}(q_t||q_{1-t})}/2$.

Proof of Lemma 2.

$$\begin{aligned} & \epsilon_{CF} - \sum_t p(1-t|\mathbf{x})\epsilon_{F,t} \\ &= p(0|\mathbf{x})(\epsilon_{CF,1} - \epsilon_{F,1}) + p(1|\mathbf{x})(\epsilon_{CF,0} - \epsilon_{F,0}) \\ &= p(0|\mathbf{x}) \int \mathcal{L}_f(\mathbf{z}, 1)(q_0(\mathbf{z}|\mathbf{x}) - q_1(\mathbf{z}|\mathbf{x}))d\mathbf{z} + p(1|\mathbf{x}) \int \mathcal{L}_f(\mathbf{z}, 0)(q_1(\mathbf{z}|\mathbf{x}) - q_0(\mathbf{z}|\mathbf{x}))d\mathbf{z} \\ &\leq 2M\mathbb{T}\mathbb{V}(q_1, q_0) \leq M\mathbf{D}. \end{aligned}$$

□

$\mathbb{T}\mathbb{V}(p, q) := \frac{1}{2}\mathbb{E}|p(\mathbf{z}) - q(\mathbf{z})|$ is the total variance distance between probability density p, q . The last inequality uses Pinsker's inequality $\mathbb{T}\mathbb{V}(p, q) \leq \sqrt{D_{\text{KL}}(p||q)}/2$ twice, to get the symmetric \mathbf{D} .

Theorem 2 is a direct corollary of Lemma 2 and the following.

Lemma 3. *Define $\epsilon_F = \sum_t p(t|\mathbf{x})\epsilon_{F,t}$. We have*

$$\epsilon_f \leq 2(G(\epsilon_F + \epsilon_{CF}) - \mathbf{V}_Y). \quad (29)$$

Simply bound ϵ_{CF} in (29) by Lemma 2, we have Theorem 2. To prove Lemma 3, we first examine a bias-variance decomposition of ϵ_F and ϵ_{CF} .

$$\begin{aligned} \epsilon_{CF,t} &= \mathbb{E}_{q_{1-t}(\mathbf{z}|\mathbf{x})} \mathbf{g}_t(\mathbf{z}) \mathbb{E}_{p_{Y(t)|p_t}(\mathbf{y}|\mathbf{z})} (\mathbf{y} - \mathbf{f}_t(\mathbf{z}))^2 \\ &\geq G \mathbb{E}_{q_{1-t}(\mathbf{z}|\mathbf{x})} \mathbb{E}_{p_{Y(t)|p_t}(\mathbf{y}|\mathbf{z})} (\mathbf{y} - \mathbf{f}_t(\mathbf{z}))^2 \\ &= G \mathbb{E}_{q_{1-t}(\mathbf{z}|\mathbf{x})} \mathbb{E}_{p_{Y(t)|p_t}(\mathbf{y}|\mathbf{z})} ((\mathbf{y} - \mathbf{j}_t(\mathbf{z}))^2 + (\mathbf{j}_t(\mathbf{z}) - \mathbf{f}_t(\mathbf{z}))^2) \end{aligned} \quad (30)$$

The second line uses $|\mathbf{g}_t(\mathbf{z})| \leq G$, and the third line is a bias-variance decomposition. Now we can define $\mathbf{V}_{CF,t}(\mathbf{x}) := \mathbb{E}_{q_{1-t}(\mathbf{z}|\mathbf{x})} \mathbb{E}_{p_{Y(t)|p_t}(\mathbf{y}|\mathbf{z})} (\mathbf{y} - \mathbf{j}_t(\mathbf{z}))^2$ and $\mathbb{B}_{CF,t}(\mathbf{x}) := \mathbb{E}_{q_{1-t}(\mathbf{z}|\mathbf{x})} (\mathbf{j}_t(\mathbf{z}) - \mathbf{f}_t(\mathbf{z}))^2$, and we have

$$\epsilon_{CF,t} \geq G(\mathbf{V}_{CF,t}(\mathbf{x}) + \mathbb{B}_{CF,t}(\mathbf{x})) \implies \epsilon_{CF} \geq G(\mathbf{V}_{CF}(\mathbf{x}) + \mathbb{B}_{CF}(\mathbf{x})) \quad (31)$$

where $\mathbf{V}_{CF} := \sum_t p(1-t|\mathbf{x})\mathbf{V}_{CF,t} = \sum_t \mathbb{E}_{q(\mathbf{z}, 1-t|\mathbf{x})} \mathbb{E}_{p_{Y(t)|p_t}(\mathbf{y}|\mathbf{z})} (\mathbf{y} - \mathbf{j}_t(\mathbf{z}))^2$ and similarly $\mathbb{B}_{CF} = \sum_t \mathbb{E}_{q(\mathbf{z}, 1-t|\mathbf{x})} (\mathbf{j}_t(\mathbf{z}) - \mathbf{f}_t(\mathbf{z}))^2$. Repeat the above derivation for ϵ_F , we have

$$\epsilon_F \geq G(\mathbf{V}_F(\mathbf{x}) + \mathbb{B}_F(\mathbf{x})) \quad (32)$$

where $\mathbf{V}_F = \sum_t \mathbb{E}_{q(\mathbf{z}, t|\mathbf{x})} \mathbb{E}_{p_{Y(t)|p_t}(\mathbf{y}|\mathbf{z})} (\mathbf{y} - \mathbf{j}_t(\mathbf{z}))^2$ and $\mathbb{B}_F = \sum_t \mathbb{E}_{q(\mathbf{z}, t|\mathbf{x})} (\mathbf{j}_t(\mathbf{z}) - \mathbf{f}_t(\mathbf{z}))^2$. Now, we are ready to prove Lemma 3.

Proof of Lemma 3.

$$\begin{aligned} \epsilon_f &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} ((\mathbf{f}_1 - \mathbf{f}_0) - (\mathbf{j}_1 - \mathbf{j}_0))^2 \\ &= \mathbb{E}_q((\mathbf{f}_1 - \mathbf{j}_1) + (\mathbf{j}_0 - \mathbf{f}_0))^2 \\ &\leq 2\mathbb{E}_q((\mathbf{f}_1 - \mathbf{j}_1)^2 + (\mathbf{j}_0 - \mathbf{f}_0)^2) \\ &= 2 \int [(\mathbf{f}_1 - \mathbf{j}_1)^2 q(\mathbf{z}, 1|\mathbf{x}) + (\mathbf{j}_0 - \mathbf{f}_0)^2 q(\mathbf{z}, 0|\mathbf{x}) + \\ &\quad (\mathbf{f}_1 - \mathbf{j}_1)^2 q(\mathbf{z}, 0|\mathbf{x}) + (\mathbf{j}_0 - \mathbf{f}_0)^2 q(\mathbf{z}, 1|\mathbf{x})] d\mathbf{z} \\ &= 2(\mathbb{B}_F + \mathbb{B}_{CF}) \leq 2(G(\epsilon_F + \epsilon_{CF}) - \mathbf{V}_Y) \end{aligned}$$

□

The first inequality uses $(a + b)^2 \leq 2(a^2 + b^2)$. The next equality splits $q(\mathbf{z}|\mathbf{x})$ into $q(\mathbf{z}, 0|\mathbf{x})$ and $q(\mathbf{z}, 1|\mathbf{x})$ and rearranges to get \mathbb{B}_F and \mathbb{B}_{CF} . The last inequality uses the two bias-variance decompositions, and $\mathbf{V}_Y = \mathbf{V}_F + \mathbf{V}_{CF}$.

B ADDITIONAL BACKGROUNDS

B.1 PROGNOSTIC SCORE AND BALANCING SCORE

In the fundamental work of (Hansen, 2008), prognostic score is defined equivalently to our \mathbf{p}_0 (P0-score), but it in addition requires no effect modification to work for $Y(1)$. Thus, a useful prognostic score corresponds to our PGS. We give main properties of PGS as following.

Proposition 3. *If V gives exchangeability, and $\mathbf{p}_t(V)$ is a PGS, then $Y(t) \perp\!\!\!\perp V, T | \mathbf{p}_t$.*

The following three properties of conditional independence will be used repeatedly in proofs.

Proposition 4 (Properties of conditional independence). (Pearl, 2009, Sec. 1.1.55) *For random variables W, X, Y, Z . We have:*

$$\begin{aligned} X \perp\!\!\!\perp Y | Z \wedge X \perp\!\!\!\perp W | Y, Z &\implies X \perp\!\!\!\perp W, Y | Z \text{ (Contraction)}. \\ X \perp\!\!\!\perp W, Y | Z &\implies X \perp\!\!\!\perp Y | W, Z \text{ (Weak union)}. \\ X \perp\!\!\!\perp W, Y | Z &\implies X \perp\!\!\!\perp Y | Z \text{ (Decomposition)}. \end{aligned}$$

Proof of Proposition 3. From $Y(t) \perp\!\!\!\perp T | V$ (exchangeability of V), and since \mathbf{p}_t is a function of V , we have $Y(t) \perp\!\!\!\perp T | \mathbf{p}_t, V(1)$.

From (1) and $Y(t) \perp\!\!\!\perp V | \mathbf{p}_t(V)$ (definition of Pt-score), using contraction rule, we have $Y(t) \perp\!\!\!\perp T, V | \mathbf{p}_t$ for both t . \square

Prognostic scores are closely related to the important concept of balancing score (Rosenbaum & Rubin, 1983). Note particularly, the proposition implies $Y(t) \perp\!\!\!\perp T | \mathbf{p}_t$ (using decomposition rule). Thus, if $\mathbf{p}(V)$ is a P-score, then \mathbf{p} also gives weak ignorability (exchangeability and overlap), which is a nice property shared with balancing score, as we will see immediately.

Definition 4 (Balancing score). $\mathbf{b}(V)$, a function of random variable V , is a balancing score if $T \perp\!\!\!\perp V | \mathbf{b}(V)$.

Proposition 5. *Let $\mathbf{b}(V)$ be a function of random variable V . $\mathbf{b}(V)$ is a balancing score if and only if $f(\mathbf{b}(V)) = p(T = 1 | V) := e(V)$ for some function f (or more formally, $e(V)$ is $\mathbf{b}(V)$ -measurable). Assume further that V gives weak ignorability, then so does $\mathbf{b}(V)$.*

Obviously, the propensity score $e(V) := p(T = 1 | V)$, the propensity of assigning the treatment given V , is a balancing score (with f be the identity function). Also, given any invertible function v , the composition $v \circ \mathbf{b}$ is also a balancing score since $f \circ v^{-1}(v \circ \mathbf{b}(V)) = f(\mathbf{b}(V)) = e(V)$.

Compare the definition of balancing score and prognostic score, we can say balancing score is sufficient for the treatment T ($T \perp\!\!\!\perp V | \mathbf{b}(V)$), while prognostic score (Pt-score) is sufficient for the potential outcomes $Y(t)$ ($Y(t) \perp\!\!\!\perp V | \mathbf{p}_t(V)$). They complement each other; conditioning on either deconfounds the potential outcomes from treatment, with the former focuses on the treatment side, the latter on the outcomes side.

B.2 VAE, CONDITIONAL VAE, AND IVAE

VAEs (Kingma et al., 2019) are a class of latent variable models with latent variable Z , and observable Y is generated by the decoder $p_\theta(\mathbf{y}|\mathbf{z})$. In the standard formulation (Kingma & Welling, 2013), the variational lower bound $\mathcal{L}(\mathbf{y}; \theta, \phi)$ of the log-likelihood is derived as:

$$\begin{aligned} \log p(\mathbf{y}) &\geq \log p(\mathbf{y}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{y}) \| p(\mathbf{z}|\mathbf{y})) \\ &= \mathbb{E}_{\mathbf{z} \sim q} \log p_\theta(\mathbf{y}|\mathbf{z}) - D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{y}) \| p(\mathbf{z})), \end{aligned} \tag{33}$$

where D_{KL} denotes KL divergence and the encoder $q_\phi(\mathbf{z}|\mathbf{y})$ is introduced to approximate the true posterior $p(\mathbf{z}|\mathbf{y})$. The decoder p_θ and encoder q_ϕ are usually parametrized by NNs. We will omit the parameters θ, ϕ in notations when appropriate.

The parameters of the VAE can be learned with stochastic gradient variational Bayes. With Gaussian latent variables, the KL term of \mathcal{L} has closed form, while the first term can be evaluated by drawing samples from the approximate posterior q_ϕ using the reparameterization trick (Kingma & Welling, 2013), then, optimizing the evidence lower bound (ELBO) $\mathbb{E}_{\mathbf{y} \sim \mathcal{D}}(\mathcal{L}(\mathbf{y}))$ with data \mathcal{D} , we train the VAE efficiently.

Conditional VAE (CVAE) (Sohn et al., 2015; Kingma et al., 2014) adds a conditioning variable C , usually a class label, to standard VAE (See Figure 1). With the conditioning variable, CVAE can give better reconstruction of each class. The variational lower bound is

$$\log p(\mathbf{y}|\mathbf{c}) \geq \mathbb{E}_{\mathbf{z} \sim q} \log p(\mathbf{y}|\mathbf{z}, \mathbf{c}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{y}, \mathbf{c})||p(\mathbf{z}|\mathbf{c})). \quad (34)$$

The conditioning on C in the prior is usually omitted (Doersch, 2016), i.e., the prior becomes $Z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ as in standard VAE, since the dependence between C and the latent representation is also modeled in the encoder q . Moreover, unconditional prior in fact gives better reconstruction because it encourages learning representation independent of class, similarly to the idea of beta-VAE (Higgins et al., 2017).

As mentioned, *identifiable* VAE (iVAE) (Khemakhem et al., 2020a) provides the first identifiability result for VAE, using auxiliary variable X . It assumes $Y \perp\!\!\!\perp X|Z$, that is, $p(\mathbf{y}|\mathbf{z}, \mathbf{x}) = p(\mathbf{y}|\mathbf{z})$. The variational lower bound is

$$\begin{aligned} \log p(\mathbf{y}|\mathbf{x}) &\geq \log p(\mathbf{y}|\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{y}, \mathbf{x})||p(\mathbf{z}|\mathbf{y}, \mathbf{x})) \\ &= \mathbb{E}_{\mathbf{z} \sim q} \log p_{\mathbf{f}}(\mathbf{y}|\mathbf{z}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{y}, \mathbf{x})||p_{\mathbf{T}, \lambda}(\mathbf{z}|\mathbf{x})), \end{aligned} \quad (35)$$

where $Y = \mathbf{f}(Z) + \epsilon$, ϵ is additive noise, and Z has exponential family distribution with sufficient statistics \mathbf{T} and parameter $\lambda(X)$. Note that, unlike CVAE, the decoder does *not* depend on X due to the independence assumption.

Here, *identifiability of the model* means that the functional *parameters* $(\mathbf{f}, \mathbf{T}, \lambda)$ can be identified (learned) up to certain simple transformation. Further, in the limit of $\epsilon \rightarrow \mathbf{0}$, iVAE solves the nonlinear ICA problem of recovering $Z = \mathbf{f}^{-1}(Y)$.

C EXPOSITIONS

The order of subsections below follows that they are referred in the main text.

C.1 LIST OF ASSUMPTIONS

The following is a list of assumptions required by our identification theory, with comments on their roles and subtleties.

(G1) additive noise model is needed to ensure the existence of PtSs. **(G1')** is equivalent to **(G1)**, and is introduced for better presentation, e.g., it connects to **(G2)** and **(M1)** through injectivity.

(M1) and **(D1)** are inherited from iVAE and are required for model (parameter) identifiability (identifying \mathbf{f}_t up to affine mapping), which does not imply CATE identification in general. Arguably here the most important is that the mapping \mathbf{f}_t from latent Z to outcome Y is injective, or else some information of Z is in principle unrecoverable. These two conditions are not required by Proposition 2 which does not need model identifiability.

(M2), together with overlapping PtSs, is important to address limited overlap of X and can be seen as a weak form of OOD generalization.

(M3') means 1) we need to know or learn the distribution of hidden noise ϵ and 2) noiseless prior. This simplifies the proof of identification, but when implementing the VAE as an estimation method, both noises are learned.

(D2), or in fact **(D2')**, strengthens the model identifiability to determine both \mathbf{f}_0 and \mathbf{f}_1 up to the *same* affine mapping, which replaces the balance of PS.

(G2) is required by Proposition 2 but not Theorem 1. It is no less important than **(G1')**, because the core intuition of our method is that **(G2)** should hold approximately. Sec. C.3 contains several detailed real-world examples on **(G2)**.

C.2 DETAILS AND EXPLANATIONS ON INTACT-VAE

Our goal is to build a model that can be learned by VAE from observational data to obtain a PGS, or more ideally bPGS, via the latent variable Z . That is, a generative prognostic model. Generative models are useful to solve the inverse problem of recovering PGSs.

With the above goal, the generative model of our VAE is built as (3). Conditioning on X in the joint model $p(\mathbf{y}, \mathbf{z}|\mathbf{x}, t)$ reflects that our estimand is CATE given X . Modeling the score by a conditional distribution rather than a deterministic function is more flexible.

The ELBO of our model can be derived from standard variational lower bound as following:

$$\begin{aligned} \log p(\mathbf{y}|\mathbf{x}, t) &\geq \log p(\mathbf{y}|\mathbf{x}, t) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}, \mathbf{y}, t)||p(\mathbf{z}|\mathbf{x}, \mathbf{y}, t)) \\ &= \mathbb{E}_{\mathbf{z}\sim q} \log p(\mathbf{y}|\mathbf{z}, t) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}, \mathbf{y}, t)||p(\mathbf{z}|\mathbf{x}, t)). \end{aligned} \quad (36)$$

We naturally have an identifiable conditional VAE (CVAE), as the name suggests. Note that (3) has a similar factorization with the generative model of iVAE (Khemakhem et al., 2020a), that is $p(\mathbf{y}, \mathbf{z}|\mathbf{x}) = p(\mathbf{y}|\mathbf{z})p(\mathbf{z}|\mathbf{x})$; the first factor does not depend on X . Further, since we have the conditioning on T in both the factors of (3), our VAE architecture is a combination of iVAE and CVAE (Sohn et al., 2015; Kingma et al., 2014), with T as the conditioning variable. See Figure 1 for the comparison in terms of graphical models. The core idea of iVAE is reflected in our model identifiability (see Lemma 1).

Please do not confuse the DGP (**G1**) and the generative model (3) of Intact-VAE. The former is the causal model, but the latter is not (at least before we show the TE identifications in Sec. 3.2). In our case, the generative model is built as a way to learn the scores through the correspondence to (2).

In particular, note that conditionally balanced representation $Z \perp\!\!\!\perp T|X$ is possible under the generative model. This requires a violation of *causal faithfulness*, so that there are other conditional independence relations, which are not generally implied by the graphical model. Our method, based on iVAE, which achieves ICA, performs nonlinear ICA to recover the scores. In fact, ICA procedures often violate causal faithfulness, because it requires finding causes from effects. Also, the violation of causal faithfulness is not caused by the generative model (which is shown in Figure 1), because the representation is learned by the encoder, and $Z \perp\!\!\!\perp T|X$ is enforced by β .

C.3 DISCUSSIONS AND EXAMPLES OF (G2)

We focus on univariate outcome on \mathbb{R} which is the most practical case and the intuitions apply to more general types of outcomes. Then, i , the mapping between μ_0 and μ_1 , is monotone, i.e, either increasing or decreasing. The increasing i means, if a change of the value of X increases (decreases) the outcome in the treatment group, then it is also the case for the controlled group. This is often true because the treatment does *not* change the mechanism how the covariates affect the outcome, under the principle of “independence of causal mechanisms (ICM)” (Janzing & Scholkopf, 2010). The decreasing i corresponds to another common interpretation when ICM does not hold. Now, the treatment does change the way covariates affect Y , but in a *global* manner: it acts like a “switch” on the mechanism: the same change of X always has *opposite* effects on the two treatment groups.

We support the above reasoning by real world examples. First we give two examples where μ_0 and μ_1 are both monotone increasing. This, and also that both μ_t are monotone decreasing, are natural and sufficient conditions for increasing i , though not necessary. The first example is form Health. (Starling et al., 2019) mentions that gestational age (length of pregnancy) has a monotone increasing effect on babies’ birth weight, regardless of many other covariates. Thus, if we intervene on one of the other binary covariates (say, t = receive healthcare program or not), both μ_t should be monotone increasing in gestational age. The next example is from economics. (Gan & Li, 2016) shows that job-matching probability is monotone increasing in market size. Then, we can imagine that, with t = receive training in job finding or not, the monotonicity is not changed. Intuitively, the examples corresponds to two common scenarios: the causal effects are accumulated though time (the first example), or the link between a covariate and the outcome is direct and/or strong (the second example).

Examples for decreasing i are rarer and the following is a bit deliberate. This example is also about babies’ birth weight as the outcome. (Abrevaya et al., 2015) shows that, with t = mother smokes

or not and X = mother’s age, the CATE $\tau(x)$ is monotone decreasing for $20 < x < 26$ (smoking decreases birth weight, and the absolute causal effect is larger for older mother). On the other hand, it is shown that birth weight slightly increases (by about 100g) in the same age range in a surveyed population (Wang et al., 2020). Thus, it is convince that, smoking changes the the tendency of birth weight w.r.t mother’s age from increasing to decreasing, and gives the large decreasing of birth weight (by about 300g) as its causal effect. This could be understood: the negative effects of smoking on mother’s heath and in turn on birth weight are accumulated during the many years of smoking.

C.4 COMPLEMENTARITY BETWEEN THE TWO IDENTIFICATIONS

We examine the complementarity between the two identifications more closely. The conditions **(M3)** / **(M3’)** and **(G2)** / **(D2’)** form two pairs, and are complementary inside each pair. The first pair matches model and truth, while the second pair restricts the discrepancy between the treatment groups. In Theorem 1, **(G2)** ($p_0 = p_1$) is replaced by **(D2’)** which instead makes $\mathcal{A}_0 = \mathcal{A}_1 := \mathcal{A}$ in (5). And **(D2’)** is easily satisfied with high-dimensional X , even if the possible values of C, d are restricted to $C = cI$ and $d = \mathbf{0}$ (see below). On the other hand, $p_\epsilon = p_e$ in **(M3’)** is impractical, but it ensures that $p_\theta(y|x, t) = p(y|x, t)$ so that (5) can be used. In Sec. 4.1, we consider practical estimation method and introduce the *regularization* that encourages learning a PGS similar to bPGS so that $p_\epsilon = p_e$ can be relaxed.

C.5 IDEAS AND CONNECTIONS BEHIND THE ELBO (7)

Bayesian approach is favorable to express the prior belief that bPGSs exist and the preference for them, and to still have reasonable posterior estimation when the belief fails and learning general PGS is necessary. This is the causal importance of VAE as an estimation method for us. By the unconditional but still flexible Λ , and also the identifications, the ELBO encourages the recovery of an approximate bPGS as the posterior, which still learns the dependence on T if necessary. Moreover, β expresses our additional knowledge (or, inductive bias) about whether or not there exist approximate bPGSs (e.g., from domain expertise).

In fact, β connects our VAE to β -VAE (Higgins et al., 2017), which is closely related to noise and variance control (Doersch, 2016, Sec. 2.4)(Mathieu et al., 2019).

Considerations on noise modeling. In Theorem 1, with large and mismatched *noises* (then **(M3’)** is easily violated), the identification of outcome model $f_t = j_t \circ v^{-1}$ would fail, and, in turn, the prior would learn confounding bias, by confusing the causal effect of T on p_T and the correlation between T and X . This is another reason to prefer $\lambda_0 = \lambda_1$, besides balancing. On the other hand, the posterior conditioning on Y provides information of noise e , and it is shown in (Bonhomme & Weidner, 2021) that posterior effect estimation has *minimum worst-case error* under model misspecification (of the noise and prior, in our case).

Under large e , a relatively small β implicitly encourages g *smaller* than the scale of e , through stressing the third term in ELBO (7). And the the model as a whole would still learn $p(y|x, t)$ well, because the uncertainty of e can be moved to and modeled by the prior. This is why k is *not* set to zero because learnable prior noise (variance) allows us to implicitly control g via β . Intuitively, smaller g strengthens the correlation between Y and Z in our model, and this naturally reflects that posterior conditioning on Y is more important under larger e . Hopefully, precise learning of outcome noise **(M3’)** is not required, as in Proposition 2.

Now, it is clear that β naturally controls at the same time noise scale and balancing. And the regularization can also be understood as an interpolation between Proposition 2 and Theorem 1: relying on bPGS, or on model identifiability; learning loosely, or precisely, the outcome regression. When the noise scale is different from truth, there would be error due to imperfect recovery of j . Sec. 4.2 shows that this error and balancing form a trade-off, which is adjusted by β .

Importance of balancing from misspecification view. If we must learn an unapproximate bPGS, we have larger misspecification under a balanced prior and rely more on Y in the posterior. Both are bad because it is shown in (Bonhomme & Weidner, 2021) that posterior only helps under bounded (small) misspecification, and posterior estimator has higher variance than prior estimator (see below

for an extreme case). Again, we want a regularizer to encourage learning of bPGS, so that we can explore the *middle ground*: relatively low-dimensional \mathbf{p} , or relatively small \mathbf{e} .

Example. Assume the true outcome noise is (near) zero. By setting $\epsilon \rightarrow \mathbf{0}$ in our model, the posterior $p_{\theta}(\mathbf{z}|\mathbf{x}, \mathbf{y}, t) = p_{\theta}(\mathbf{y}, \mathbf{z}|\mathbf{x}, t)/p_{\theta}(\mathbf{y}|\mathbf{x}, t)$ degenerates to $\mathbf{f}_T^{-1}(Y) = \mathbf{f}_T^{-1}(j_T(\mathbf{p}_T)) = \mathbf{v}^{-1}(\mathbf{p}_T)$, a *factual* PGS. However, $\mathbf{f}_{1-T}^{-1}(Y) = \mathbf{f}_{1-T}^{-1}(j_T(\mathbf{p}_T)) = \mathbf{v}^{-1}(j_{1-T}^{-1} \circ j_T(\mathbf{p}_T)) \neq \mathbf{v}^{-1}(\mathbf{p}_{1-T})$, the score recovered by posterior does not work for counterfactual assignment! The problem is, unlike X , the outcome $Y = Y(T)$ is affected by T , and, the degenerated posterior disregards the information of X from the prior and depends exclusively on factual (Y, T) .

C.6 CONSISTENCY OF VAE AND PRIOR ESTIMATION

The following is a refined version of Theorem 4 in Khemakhem et al. (2020a). The result is proved by assuming: i) our VAE is flexible enough to ensure the ELBO is tight (equals to the true log likelihood) for some parameters; ii) the optimization algorithm can achieve the *global* maximum of ELBO (again equals to the log likelihood).

Proposition 6 (Consistency of Intact-VAE). *Given model (3)&(6), and let $p^*(\mathbf{x}, \mathbf{y}, t)$ be the true observational distribution, assume*

- i) *there exists $(\bar{\theta}, \bar{\phi})$ such that $p_{\bar{\theta}}(\mathbf{y}|\mathbf{x}, t) = p^*(\mathbf{y}|\mathbf{x}, t)$ and $p_{\bar{\theta}}(\mathbf{z}|\mathbf{x}, \mathbf{y}, t) = q_{\bar{\phi}}(\mathbf{z}|\mathbf{x}, \mathbf{y}, t)$;*
- ii) *the ELBO $\mathbb{E}_{\mathcal{D} \sim p^*}(\mathcal{L}(\mathbf{x}, \mathbf{y}, t; \theta, \phi))$ (4) can be optimized to its global maximum at (θ', ϕ') ;*

Then, in the limit of infinite data, $p_{\theta'}(\mathbf{y}|\mathbf{x}, t) = p^(\mathbf{y}|\mathbf{x}, t)$ and $p_{\theta'}(\mathbf{z}|\mathbf{x}, \mathbf{y}, t) = q_{\phi'}(\mathbf{z}|\mathbf{x}, \mathbf{y}, t)$.*

Proof. From i), we have $\mathcal{L}(\mathbf{x}, \mathbf{y}, t; \bar{\theta}, \bar{\phi}) = \log p^*(\mathbf{y}|\mathbf{x}, t)$. But we know \mathcal{L} is upper-bounded by $\log p^*(\mathbf{y}|\mathbf{x}, t)$. So, $\mathbb{E}_{\mathcal{D} \sim p^*}(\log p^*(\mathbf{y}|\mathbf{x}, t))$ should be the global maximum of the ELBO (even if the data is finite).

Moreover, note that, for any (θ, ϕ) , we have $D_{\text{KL}}(p_{\theta}(\mathbf{z}|\mathbf{x}, \mathbf{y}, t) \| q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y}, t)) \geq 0$ and, in the limit of infinite data, $\mathbb{E}_{\mathcal{D} \sim p^*}(\log p_{\theta}(\mathbf{y}|\mathbf{x}, t)) \leq \mathbb{E}_{\mathcal{D} \sim p^*}(\log p^*(\mathbf{y}|\mathbf{x}, t))$. Thus, the global maximum of ELBO is achieved *only* when $p_{\theta}(\mathbf{y}|\mathbf{x}, t) = p^*(\mathbf{y}|\mathbf{x}, t)$ and $p_{\theta}(\mathbf{z}|\mathbf{x}, \mathbf{y}, t) = q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y}, t)$. \square

Consistent prior estimation of CATE follows directly from the identifications. The following is a corollary of Theorem 1.

Corollary 1. *Under the conditions of Theorem 1, further require the consistency of Intact-VAE. Then, in the limit of infinite data, we have $\mu_t(X) = \mathbf{f}_t(\mathbf{h}_t(X))$ where \mathbf{f}, \mathbf{h} are the optimal parameters learned by the VAE.*

C.7 PRE/POST-TREATMENT PREDICTION

Sampling posterior requires *post-treatment* observation (\mathbf{y}, t) . Often, it is desirable that we can also have *pre-treatment* prediction for a new subject, with only the observation of its covariate $X = \mathbf{x}$. To this end, we use the prior as a pre-treatment predictor for Z : replace q_{ϕ} with p_{Λ} in (8) and get rid of the outer average taken on \mathcal{D} ; all the others remain the same. We also have sensible pre-treatment prediction even without true low-dimensional PSSs, because p_{Λ} gives the best balanced approximation of the target PGS. The results of pre-treatment prediction are given in the experimental section E.

D MORE ON RELATED WORK

D.1 CFR AND CEVAE

CFR and CEVAE are well-known machine learning methods for CATE estimation. Here we make detailed comparisons to them.

D.1.1 COMPARISONS WITH CFR

Our method is related to CFR in two ways. Theoretically, our bounds in Sec. 4.2 resemble those in Shalit et al. (2017). But we bound CATE error, while CFR bounds PEHE; thus, our bounds give

conditional balancing while CFR only has unconditional balancing. See Sec. D.3 for more on the bounds. Conceptually, CFR is loosely related to our method because it also learns a representation as an outcome predictor, as mentioned in the follow-up Johansson et al. (2020). However, CFR does not have a generative model, so their representation is not formally related to PGSs. Moreover, CFR does not account the outcome noise, while the uncertainty due to the noise is accounted by our VAE.

D.1.2 COMPARISONS WITH AND CRITICISMS OF CEVAE

Motivation CEVAE is motivated by exploiting proxy variables, and its intuition is that the hidden confounder U can be recovered by VAE from proxy variables.

Our method is motivated by prognostic scores (Hansen, 2008), and our model is directly based on equations (2) which identifies CATE. There is no need to recover the hidden confounder in our framework.

Architecture Our model is naturally based on (2), particularly the independence properties of PGS. And as a consequence, our VAE architecture is a natural combination of iVAE and CVAE (see Figure 1). Our ELBO (4) is derived by the standard variational lower bound.

On the other hand, the architecture of CEVAE is more ad hoc and complex. Its decoder follows the graphical model of descendant proxy mentioned above, but adds an ad hoc component to mimic TARnet (Shalit et al., 2017): it uses separated NNs for the two potential outcomes. We tried this idea on the IHDP dataset, and, as we show in Sec. 5.2, it has basically no merits for our method, because we have a principled way for balancing.

The encoder of CEVAE is even more complex. To have post-treatment estimation, $q(T|X)$ and $q(Y|X, T)$ are added into the encoder. As a result, the ELBO of CEVAE has two additional likelihood terms corresponding to the two distributions. However, in our Intact-VAE, post-treatment estimation is given naturally by our standard encoder, thanks to the correspondence between our model and (2).

Justification We have given the identifications and bounds of our method in this paper. Moreover, we carefully distinguish assumptions on the DGP and assumptions on our model, and identify the assumptions that are important for causality. There are few theoretical justifications for CEVAE. Their Theorem 1 directly assumes the joint distribution $p(\mathbf{x}, \mathbf{y}, t, \mathbf{u})$ including hidden confounder U is recovered, then identification is trivial by using the standard adjustment equation.

However, the challenge is exactly that the confounder is hidden, unobserved. Many years of work have been done in causal inference to derive conditions under which hidden confounder can be (partially) recovered (Greenland, 1980; Kuroki & Pearl, 2014; Miao et al., 2018). In particular, Miao et al. (2018) gives the most recent identification result for proxy setting, which requires very specific two proxies structure, and other completeness assumptions on distributions. Thus, it is unreasonable to believe that VAE, with simple descendant proxies, can recover the hidden confounder. Indeed, Rissanen & Marttinen (2021) recently give evidence that the method often fails.

Moreover, the identifiability of VAE itself is a challenging problem. As mentioned in Introduction, Khemakhem et al. (2020a) is the first identifiability result for VAE, but it only identifies an equivalence class, not a unique representation function. Thus, it is also unconvincing that VAE can learn a unique latent distribution, without certain assumptions. As we show in Sec. 5.1, for relatively simple synthetic datasets, CEVAE can not robustly recover the hidden confounder, even only up to transformation, while our method can (though, again, this is not needed for our method).

D.2 INJECTIVITY, INVERTIBILITY, MONOTONICITY, AND OVERLAP

Let us note that *any injective mapping defines an invertible mapping*, by restrict the domain of the inverse function to the range of the injective mapping. Also note that injectivity is weaker than monotonicity; a monotone mapping can be defined by an injective and *order-preserving* mapping between ordered sets. Particularly, *an injective and continuous mapping on \mathbb{R} is monotone*, and many works in econometrics give examples of this case.

Many classical and recent works (with many real world applications, see C.1) in econometrics are based on monotonicity. Particularly, there is a long line of work based on *monotonicity of treatment* (Huber & Wüthrich, 2018). More related to our method is another line of work based on *monotonicity of outcome*, see (Chernozhukov & Hansen, 2013) and references therein for early results. Some recent works apply monotonicity of outcome to nonparametric IV regression (NPIV) (Freyberger & Horowitz, 2015; Li et al., 2017; Chetverikov & Wilhelm, 2017), where the structural equation of the outcome is assumed to be $Y = f(T) + \epsilon$, and f is monotone and T (the treatment) is often continuous. Particularly, (Chetverikov & Wilhelm, 2017) combines monotonicity of both treatment and outcome, and (Freyberger & Horowitz, 2015) considers *discrete* treatment (note continuity or differentiability is not necessary for monotonicity). NPIV with monotone f is closely related to our method, but the difference is that T is replaced by a PGS in our method, and the PGS is recovered from observables. Finally, as we mentioned in Sec. 3.2, monotonicity is a kind of shape restriction which also includes, e.g., concavity and symmetry and attracts recent interests (Chetverikov et al., 2018). However, most of NPIV works focus on identifying f but not directly on TEs, and we do not know any works that use monotonicity to address limited overlap.

Recently in machine learning, (Johansson et al., 2019; Zhang et al., 2020b; Johansson et al., 2020) note the relationship between invertibility and overlap. As mentioned, (Johansson et al., 2020) gives bounds without overlap, but the relationship between invertibility and overlap is not explicit in their theory. (Johansson et al., 2019) explicitly discuss overlap and invertibility, but does not focus on TEs. (Zhang et al., 2020b) assumes overlap so that identification is given, and then focuses on learning overlapping representation that preserves the overlapping the covariate. However, it does not relate invertibility and overlap, but uses invertible representation function to *preserve exchangeability given the covariate*, and linear outcome regression to simplify the model. Related, our identifications required (M2), of which linearity of PGS and representation function is a sufficient condition, and our outcome model is injective, to *preserve the exchangeability given the PGS*. Thus, our method works under more general setting, and arguably under weaker conditions.

D.3 ADDITIONAL NOTES ON NOVELTIES OF THE BOUNDS IN SEC. 4.2

We give details and additional points regarding the novelties. Lu et al. (2020) also use a VAE and derive bounds most related to ours. Still, our method strengthens Lu et al. (2020), in a simpler and principled way: we distinguish true score and latent Z and show that identification is the link; considering both prior and posterior, we show the symmetric nature of the balancing term and relate it to our KL term in (7), without ad hoc regularization; moreover, we consider outcome noise modeling which is a strength of VAE and relate it to hyperparameter β . Particularly, in (Lu et al., 2020), latent variable Z is confused with the true representation (p_t up to invertible mapping in our case). *Without* identification, the method in fact has unbounded error. Note that Shalit et al. (2017) do not consider connection to identification and noise modeling as well. The error between $\hat{\tau}_f$ and τ_j , which we bound, is due to the unknown outcome noise that is not accounted by our Theorem 1; thus, the theory in Sec. 4.2 is complementary to that in Sec. 3.2. Finally, β is a trade-off between the conditional balance of learned PGS (affected by f_t), and precision/effective sample size of outcome regression—and can be seen as the probabilistic counterpart of Tarr & Imai (2021) and Kallus et al. (2018).

E DETAILS AND ADDITIONS OF EXPERIMENTS

We evaluate the post-treatment performance on training and validation set jointly (This is non-trivial. Recall the fundamental problem of causal inference). The treatment and (factual) outcome should not be observed for pre-treatment predictions, so we report them on a testing set. See also Sec. C.7 the pre/post-treatment distinction.

E.1 SYNTHETIC DATA

We detail how the random parameters in the DGPs are sampled. μ_i and σ_i are uniformly sampled in range $(-0.2, 0.2)$ and $(0, 0.2)$, respectively. The weights of linear functions h, k, l are sampled from standard normal distributions. The NNs f_0, f_1 use leaky ReLU activation with $\alpha = 0.5$ and are of 3 to 8 layers randomly, and the weights of each layer are sampled from $(-1.1, -0.9)$. To have a large but still reasonable outcome variance, the output of f_t is divided by $C_t := \text{Var}_{\{D|T=t\}}(f_t(Z))$. When generating DGPs with dependent noise, the variance parameter g_t for the outcome is generated by adding a softplus layer after respective f_t , and then normalized to range $(0, 2)$.

We use the original implementation of CFR³. Very possibly due to bugs in implementation, the CFR version using Wasserstein distance has error of TensorFlow type mismatch on our synthetic dataset, and the CFR version using MMD diverges with very large loss value on one or two of the 10 random DGPs. We use MMD version, and, when the divergence of training happens, report the results from trained models before divergence, which still give reasonable results. We search the balancing parameter alpha in $[0.16, 0.32, 0.64, 0.8, 1.28]$, and fix other hyperparameters as they were in the default config file.

We characterize the degree of limited overlap by examining the percentage of observed values x that give probability less than 0.001 for one of $p(t|x)$. The threshold is chosen so that all sample points near those values x almost certainly belong to a single group since we have 500 sample point in total. If we regard a DGP as very limited-overlapping when the above percentage is larger than 50%, then, as shown in Figure 4, non (all) of the 10 DGPs are very limited-overlapping with $\omega = 6$ ($\omega = 22$).

For diversity of the datasets, we set $g_t(W) = 1$ in DGPs in Appendix. Figure 5 shows, with $\text{dim}(Z) = 200$, our method works better than CFR under $\text{dim}(W) = 1$ and as well as CFR under $\text{dim}(W) > 1$. As mentioned in Conclusion, this indicates that the theoretical requirement of injective f_t in our model might be relaxed. Interestingly, larger β seems to give better results here, this is understandable because β controls the trade-off between fitting and balancing, and the fitting capacity of our decoder is much increased with $\text{dim}(Z) = 200$. Note that the above observations on $\text{dim}(Z)$ are not caused by fixing $g_t(W) = 1$ (compare Figure 5 with Figure 6 below).

Figure 6 shows the importance of noise modeling. Compared to Figure 2 in the main text, where $g_t(W)$ in DGPs is not fixed, our method works worse here, particularly for large β , because now noise modeling (g, k in the ELBO) only adds unnecessary complexity. The changes of performance w.r.t different ω should be unrelated to overlap levels, but to the complexity of random DGPs; compare to Figure 5, with larger NNs in our VAE, the changes become much insignificant. The drop of error for $\text{dim}(W) > 3$ is due to the randomness of f in (36). In Sec. 2.2, we saw that the 2-dimnsional bPGS $p := (\mu_0(X), \mu_1(X))$ always exists under additive noise models. Thus, when $\text{dim}(W) > 2$, our method tries to recover that p , and generally performs not worse than under $\text{dim}(W) = 2$, but still not better than under $\text{dim}(W) = 1$.

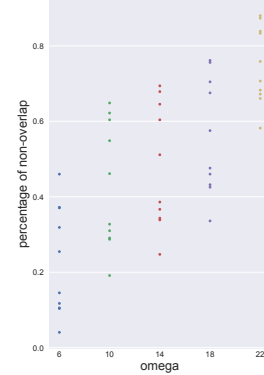


Figure 4: Degree of limited overlap w.r.t ω .

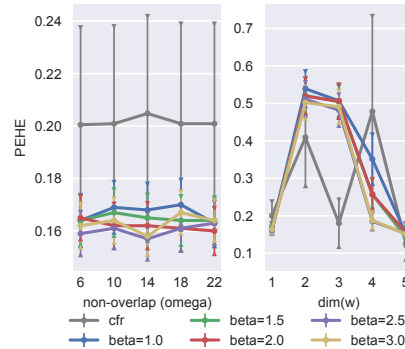


Figure 5: $\sqrt{\epsilon_{pehe}}$ on synthetic dataset, with $g_t(W) = 1$ in DGPs, and $\text{dim}(Z) = 200$ in our model. Error bar on 10 random DGPs.

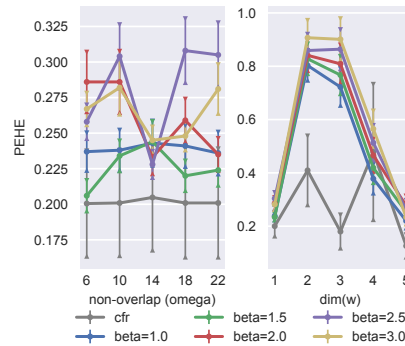


Figure 6: $\sqrt{\epsilon_{pehe}}$ on synthetic dataset, with $g_t(W) = 1$ in DGPs. Error bar on 10 random DGPs.

³<https://github.com/clinicalml/cfrnet>

Figure 7 shows results of ATE estimation. Notably, CFR drops performance w.r.t degree of limited overlap. Our method does not show this tendency except for very large β ($\beta = 3$). This might be another evidence that CFR and its unconditional balancing overfit to PEHE (see Sec. 5.2). Also note that, under $\dim(W) = 1$, $\beta = 3$ gives the best results for ATE although it does not work well for PEHE, and we do not know if this generalizes to the conclusion that large β gives better ATE estimation under the existence of bPGS, but leave this for future investigation.

Figure 8 shows results of pre-treatment prediction. In left panel, both our method and CFR perform only slightly worse than post-treatment. This is reasonable because here we have bPGS W with $\dim(W) = 1$, there is no need to learn PGS. In the right panel, we also do not see significant drop of performance compared to post-treatment. This might be due to the hardness of learning approximate bPGS in this dataset, and posterior estimation does not give much improvements.

You can find more plots for latent recovery at the end of the paper.

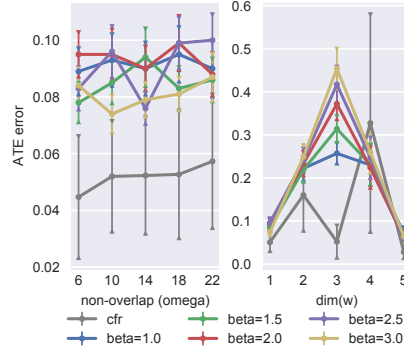


Figure 7: ϵ_{ate} on synthetic dataset, with $g_t(W) = 1$ in DGPs. Error bar on 10 random DGPs.

E.2 IHDP

IHDP is based on an RCT where each data point represents a baby with 25 features (6 continuous, 19 binary) about their birth and mothers. Race is introduced as a confounder by artificially removing all treated children with nonwhite mothers. There are 747 subjects left in the dataset. The outcome is synthesized by taking the covariates (features excluding Race) as input, hence *unconfoundedness* holds given the covariates. Following previous work, we split the dataset by 63:27:10 for training, validation, and testing. Note, there is no ethical concerns here, because the treatment assignment mechanism is artificial by processing the data. Also our results are only quantitative and we make no ethical conclusions.

The generating process is as following (Hill, 2011, Sec. 4.1).

$$Y(0) \sim \mathcal{N}(e^{\mathbf{a}^T(X+\mathbf{b})}, 1), \quad Y(1) \sim \mathcal{N}(\mathbf{a}^T X - c, 1), \tag{37}$$

where \mathbf{a} is a random coefficient, \mathbf{b} is a constant bias with all elements equal to 0.5, and c is a random parameter adjusting degree of overlapping between the treatment groups. As we can see, $\mathbf{a}^T X$ is a true bPGS. As mentioned in the main text, the bPGS might be discrete. Thus, this experiment also shows the importance of VAE, even if an apparent bPGS exists. Under *discrete* PSs, training an regression based on Proposition 2 is hard, but our VAE works well.

The two added components in the modified version of our method are as following. First, we build the two outcome functions $f_t(Z)$, $t = 0, 1$ in our learning model (3), using two separate NNs. Second, we add to our ELBO (4) a regularization term, which is the Wasserstein distance (Cuturi, 2013) between $\mathbb{E}_{\mathcal{D} \sim p(X|T=t)} p_{\Lambda}(Z|X)$, $t \in \{0, 1\}$. As shown in Table 2, best unconditional balancing parameter is 0.1. Larger parameters gives much worse PEHE and does not improve ATE estimation. Smaller parameters are more reasonable but still do not improve the results. The overall tendency is clear. Compared to ours, CFR with its unconditional balancing does not improve ATE estimation, it may improve PEHE results with fine tuned parameter, but possibly at the price of worse ATE estimation.

Table 3 shows pre-treatment results, All methods gives reasonable results.

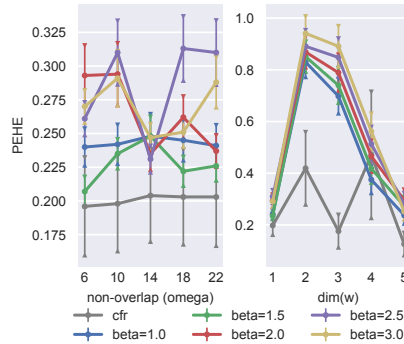


Figure 8: *Pre-treatment* $\sqrt{\epsilon_{pehe}}$ on synthetic dataset. Error bar on 10 random DGPs.

Table 2: Performance of modified version with different unconditional balancing parameter, the values of which are shown after ‘‘Mod.’’.

Method	Ours	Mod. 1	Mod. 0.2	Mod. 0.1	Mod. 0.05	Mod. 0.01	CFR
ϵ_{ate}	.177 \pm .007	.196 \pm .008	.177 \pm .007	.167 \pm .005	.177 \pm .006	.179 \pm .006	.25 \pm .01
$\sqrt{\epsilon_{pehe}}$.843 \pm .030	1.979 \pm .082	1.116 \pm .046	.777 \pm .026	.894 \pm .039	.841 \pm .029	.71 \pm .02

Table 3: *Pre-treatment* Errors on IHDP over 1000 random DGPs. We report results with $\dim(Z) = 10$. **Bold** indicates method(s) which is *significantly* better. The results are taken from Shalit et al. (2017), except GANITE (Yoon et al., 2018) and CEVAE (Louizos et al., 2017).

Method	TMLE	BNN	CFR	CF	CEVAE	GANITE	Ours
pre- ϵ_{ate}	NA	.42 \pm .03	.27 \pm .01	.40 \pm .03	.46 \pm .02	.49 \pm .05	.211\pm.011
pre- $\sqrt{\epsilon_{pehe}}$	NA	2.1 \pm .1	.76\pm.02	3.8 \pm .2	2.6 \pm .1	2.4 \pm .4	.946 \pm .048

E.3 POKEC SOCIAL NETWORK DATASET

This experiment shows our method is the best compared with the methods specialized for networked deconfounding, a challenging problem in its own right. Thus, our method has the potential to work under *unobserved confounding*, but we leave detailed experimental and theoretical investigation to future.

Pokec (Leskovec & Krevl, 2014) is a real world social network dataset. We experiment on a semi-synthetic dataset based on Pokec, which was introduced in (Veitch et al., 2019), and use exactly the same pre-processing and generating procedure. The pre-processed network has about 79,000 vertexes (users) connected by 1.3×10^6 undirected edges. The subset of users used here are restricted to three living districts which are within the same region. The network structure is expressed by binary adjacency matrix G . Following (Veitch et al., 2019), we split the users into 10 folds, test on each fold and report the mean and std of pre-treatment ATE predictions. We further separate the rest of users (in the other 9 folds) by 6 : 3, for training and validation.

Each user has 12 attributes, among which `district`, `age`, or `join date` is used as a confounder U to build 3 different datasets, with remaining 11 attributes used as covariate X . Treatment T and outcome Y are synthesised as following:

$$T \sim \text{Bern}(g(U)), \quad Y = T + 10(g(U) - 0.5) + \epsilon, \quad (38)$$

where ϵ is standard normal. Note that `district` is of 3 categories; `age` and `join date` are also discretized into three bins. $g(U)$, which is a bPGS, maps these three categories and values to $\{0.15, 0.5, 0.85\}$.

β -Intact-VAE is expected to learn a bPGS from G, X , if we can exploit the network structure effectively. Given the huge network structure, most users can practically be identified by their attributes and neighborhood structure, which means U can be roughly seen as a deterministic function of G, X . This idea is comparable to Assumptions 2 and 4 in (Veitch et al., 2019), which postulate directly that a balancing score can be learned in the limit of infinite large network. To extract information from the network structure, we use Graph Convolutional Network (GCN) (Kipf & Welling, 2017) in conditional prior and encoder of β -Intact-VAE. The implementation details are given at the end of this subsection.

Table 4 shows the results. The pre-treatment $\sqrt{\epsilon_{pehe}}$ for `Age`, `District`, and `Join date` confounders are 1.085, 0.686, and 0.699 respectively, practically the same as the ATE errors. Note that, Veitch et al. (2019) does not give individual-level prediction.

To extract information from the network structure, we use Graph Convolutional Network (GCN) (Kipf & Welling, 2017) in conditional prior and encoder of β -Intact-VAE. A difficulty is that, the network G and covariates X of *all* users are always needed by GCN, regardless of whether it is in training, validation, or testing phase. However, the separation can still make sense if we take care that the treatment and outcome are used only in the respective phase, e.g., (y_m, t_m) of a testing user m is only used in testing.

Table 4: Pre-treatment ATE on Pokec. Ground truth ATE is 1, as we can see in (38). “Unadjusted” estimates ATE by $\mathbb{E}_{\mathcal{D}}(y_1) - \mathbb{E}_{\mathcal{D}}(y_0)$. “Parametric” is a stochastic block model for networked data (Gopalan & Blei, 2013). “Embed-” denotes the best alternatives given by (Veitch et al., 2019). **Bold** indicates method(s) which is *significantly* better than all the others. We report results with 20-dimensional latent Z . The results of the other methods are taken from (Veitch et al., 2019).

	Age	District	Join Date
Unadjusted	4.34 ± 0.05	4.51 ± 0.05	4.03 ± 0.06
Parametric	4.06 ± 0.01	3.22 ± 0.01	3.73 ± 0.01
Embedding-Reg.	2.77 ± 0.35	1.75 ± 0.20	2.41 ± 0.45
Embedding-IPW	3.12 ± 0.06	1.66 ± 0.07	3.10 ± 0.07
Ours	2.08 ± 0.32	1.68 ± 0.10	1.70 ± 0.13

GCN takes the network matrix \mathbf{G} and the *whole* covariates matrix $\mathbf{X} := (\mathbf{x}_1^T, \dots, \mathbf{x}_M^T)^T$, where M is user number, and outputs a representation matrix \mathbf{R} , again for all users. During training, we *select* the rows in \mathbf{R} that correspond to users in training set. Then, treat this *training representation matrix* as if it is the covariates matrix for a non-networked dataset, that is, the downstream networks in conditional prior and encoder are the same as in the other two experiments, but take $(\mathbf{R}_{m,:})^T$ where \mathbf{x}_m was expected as input. And we have respective selection operations for validation and testing. We can still train β -Intact-VAE including GCN by Adam, simply setting the gradients of non-selected rows of \mathbf{R} to 0.

Note that GCN cannot be trained using mini-batch, instead, we perform batch gradient decent using full dataset for each iteration, with initial learning rate 10^{-2} . We use dropout (Srivastava et al., 2014) with rate 0.1 to prevent overfitting.

E.4 EMPIRICAL VALIDATION OF THE BOUNDS IN SEC. 4.2

Here we focus on the $\mathbf{D}(X)$ term in Theorem 2 because it is directly related to conditional balance.

In the Figure attached at the end of the paper, the rows correspond to 3 overlap levels from strong to weak ($\omega = 6, 14, 22$ respectively). The first column shows the histograms of correlation coefficients between $\mathbf{D}(X)$ and $\epsilon_f(X)$ on 100 random DGPs. The vertical bars in the histograms are 5, 25, 50, 75, 95 percentiles (the values are shown in the table below). The other 10 columns show the plots of distributions of $(\mathbf{D}(X), \epsilon_f(X))$ for the first 10 DGPs. The correlation coefficient for each DGP is shown as `corrcoef=*` above each histogram. The plots are in log-log scale, because both \mathbf{D} and ϵ_f are single-sided, and most data points concentrate near $(0, 0)$, making the plots bad-looking.

We have two important observations from the histograms: 1) on the majority of DGPs, there are positive correlations between \mathbf{D} and ϵ_f ; 2) the positive correlation is stronger with weaker overlap (the portion of large correlation increases, and the mean `corrcoef` are 0.100, 0.110, and 0.121, respectively).

Thus, our bounds and conditional balance have significance. Not all DGPs have positive correlations, and this is reasonable because our bound (11) has three other terms which can obscure the relation between \mathbf{D} and ϵ_f . The DGPs 1, 3, 6, 8, 10 show typical situations when there are positive correlations.

Table 5: Percentiles of correlation coefficients between $\mathbf{D}(X)$ and $\epsilon_f(X)$ on 100 random DGPs.

Percentile	5	25	50	75	95
$\omega = 6$	-0.289	-0.086	0.069	0.299	0.609
$\omega = 14$	-0.328	-0.124	0.055	0.337	0.636
$\omega = 22$	-0.274	-0.128	0.067	0.341	0.634

E.5 ADDITIONAL PLOTS ON SYNTHETIC DATASETS

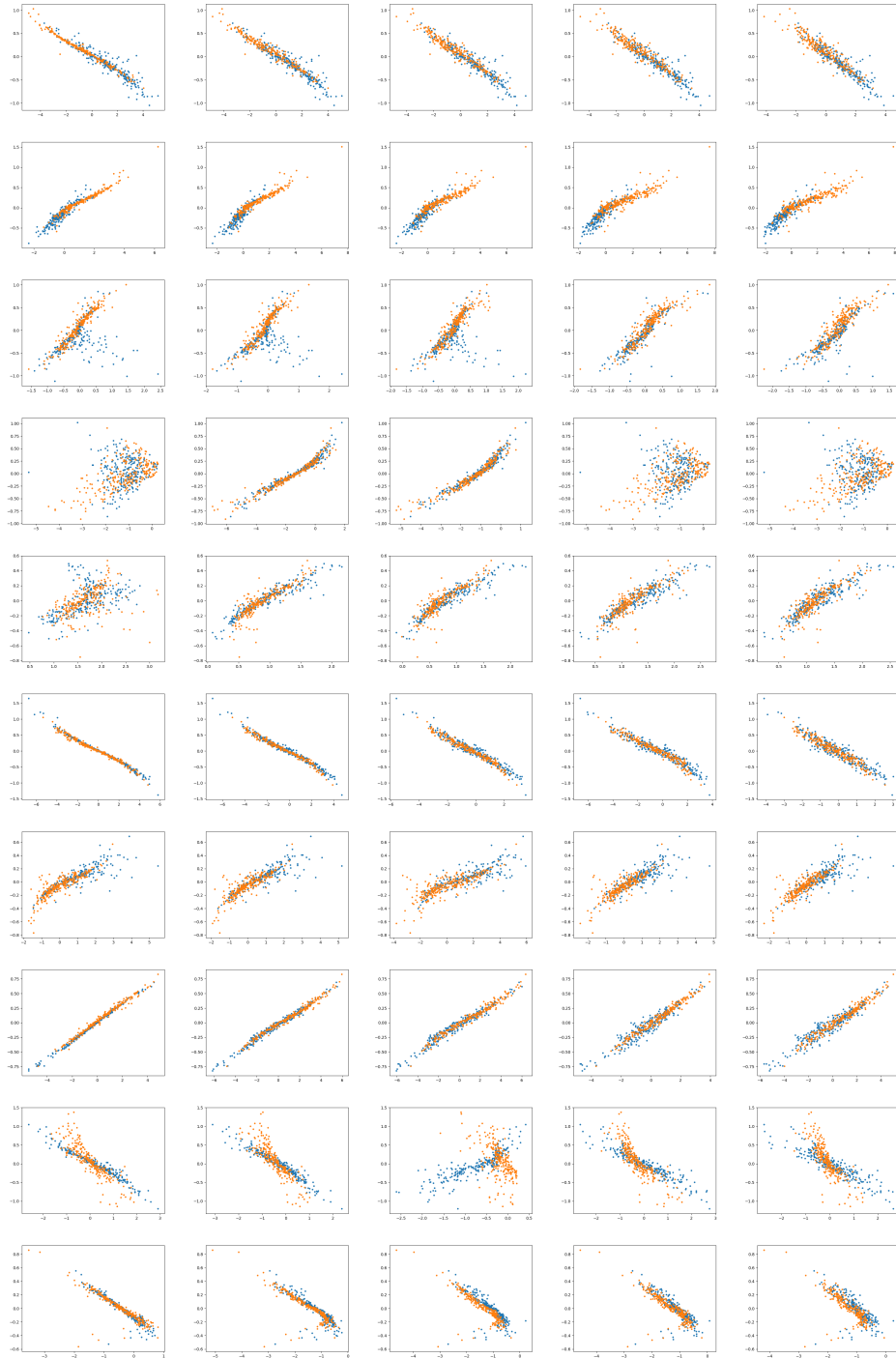


Figure 9: Plots of recovered-true latent. Rows: first 10 nonlinear random models, columns: outcome noise level.

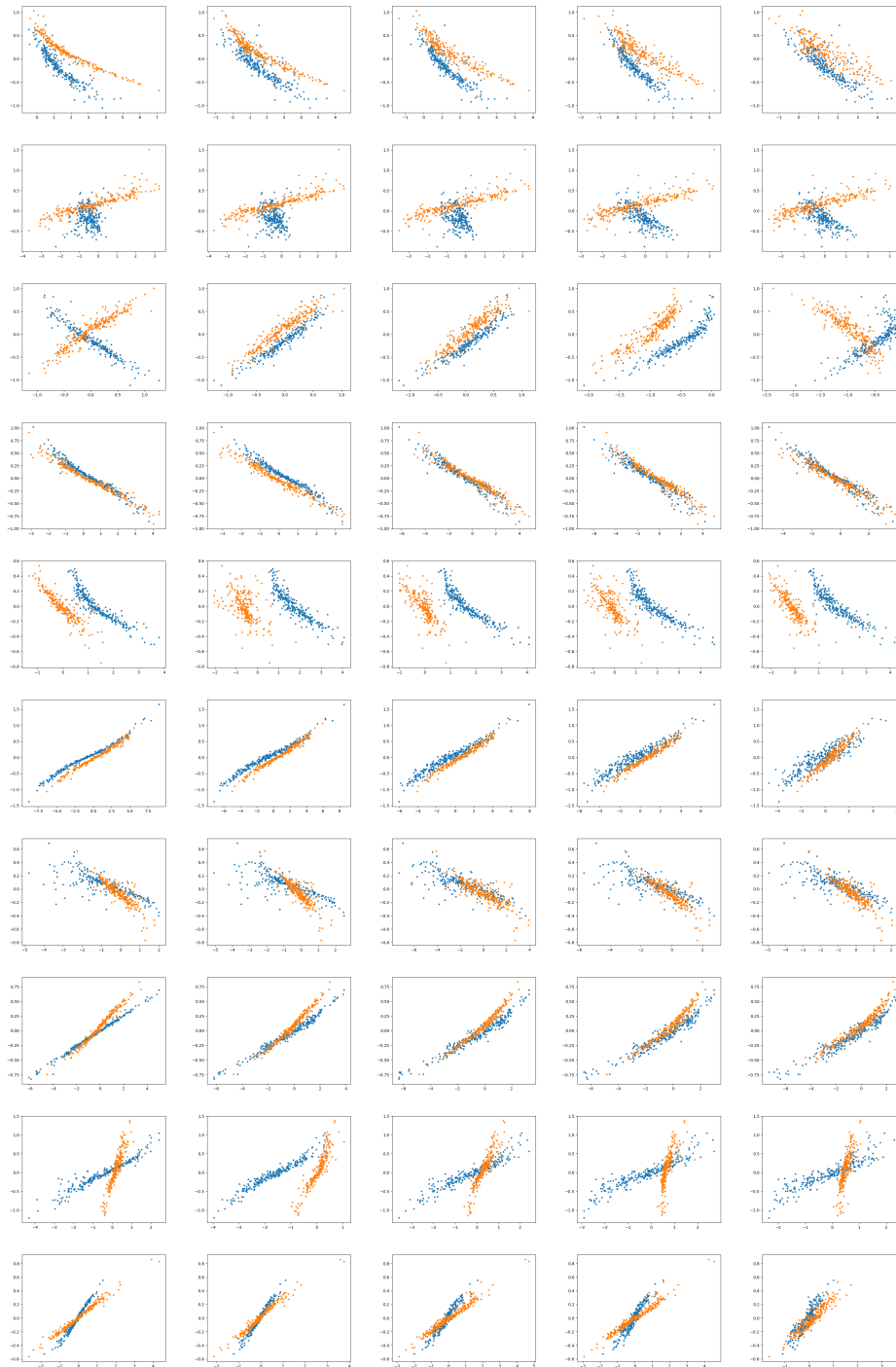


Figure 10: Plots of recovered-true latent. Conditional prior *depends* on t . Rows: first 10 nonlinear random models, columns: outcome noise level. Compare to the previous figure, we can see the transformations for $t = 0, 1$ are *not* the same, confirming the importance of balanced prior.

