
Identifying and Estimating Causal Effects under Weak Overlap by Generative Prognostic Model

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 As an important problem of causal inference, we discuss the identification and
2 estimation of treatment effects (TEs) under weak overlap, i.e., subjects with certain
3 features all belong to a single treatment group. We use a latent variable to model
4 a prognostic score (PGS), which is widely used in biostatistics and sufficient for
5 TEs, i.e., we build a generative prognostic model. We prove that the latent variable
6 recovers a PGS, and the model identifies individualized treatment effects. The
7 model is then learned as the Intact-VAE, a new type of variational autoencoder
8 (VAE). We derive counterfactual generalization bounds which motivate representa-
9 tion balanced for treatment groups conditioned on individualized features. The
10 proposed method is compared with recent methods using (semi-)synthetic datasets.

11 1 Introduction

12 Causal inference [21, 34], i.e, inferring causal effects of interventions, is a fundamental problem. In
13 this work, we focus on treatment effects (TEs), such as effects of public policies or a new drug, based
14 on a set of observations consisting of binary labels for treatment / control (non-treated), outcome, and
15 other covariates (e.g, patients' personal records). The fundamental difficulty of causal inference is
16 that we never observe *counterfactual* outcomes, which would have been if we had made the other
17 decision (treatment or control). While the ideal protocol for causal inference is randomized controlled
18 trials (RCTs), they often have ethical and practical issues, or suffer from prohibitive costs. Thus,
19 causal inference from observational data is indispensable. It introduces other challenges, however.
20 The most crucial one is *confounding*: there may be variables (called *confounders*) that causally affect
21 both the treatment and the outcome, and spurious correlation follows.

22 A large majority of works, including this work, rely on the *unconfoundedness*, which means that
23 appropriate covariates are collected so that the confounding can be controlled by conditioning on
24 covariates. That is, all the confounders are in essence observed. This is still challenging, due to
25 systematic *imbalance* (difference) of the distributions of the covariates between the treatment and
26 control groups, introducing bias in estimation. Among classical ways of dealing with imbalance
27 are matching and re-weighting [44, 35]. Machine learning methods are also exploited; there are
28 semi-parametric methods, e.g, [48, TMLE], which have better finite sample performance, and also
29 non-parametric, tree-based methods, e.g., [49, Causal Forests (CF)]. Notably, starting from [23], there
30 is a recent rise of interest in learning representation of covariates, which is independent of treatment
31 groups, i.e., *balanced* representation learning (BRL) .

32 The most serious form of imbalance is that sample points with certain values of covariate are all
33 belong to a single treatment group, which is called *weak overlap* of the covariate. Causal effects are
34 not directly estimable at non-overlapped covariate values. There are lines of work that give robustness
35 to weak overlap [3], trim non-overlapped sample points [52], or study convergence rate depending
36 on overlap [19]. Weak overlap is particularly relevant to machine learning methods exploiting rich
37 covariates, because, with higher-dimensional covariates, overlap is harder to satisfy and verify [10].

38 Our approach to the weak overlap issue is based on the prognostic score (PGS) [14], which is among
 39 the important concepts of sufficient scores. While the most well-known score is the propensity score
 40 (PPS) [36], PGSs have also long been known to improve methods using PPS [37, 5], and interests last
 41 in biostatistics [45, 2]. Prognostic modeling can benefit more from predictive systems and exploit
 42 richer literature than propensity modeling, particularly in Medicine and Health. A comparative study
 43 in [13] shows PGS-based methods perform better, or as well as, PPS methods. Thus, it is promising
 44 to combine the predictive powers of prognostic modeling and machine learning.

45 To solve the inverse problem of recovering PGSs, our method exploits also the recent advance
 46 of identifiable representation, particularly of VAE [26, iVAE]. *Identification* means parameters of
 47 interest (for us, representation function and causal effects) are uniquely determined and given by
 48 true *observational* distribution. Identification logically precedes estimation and inference. Without
 49 identification there is no hope of a consistent estimator, and a model would fail silently; it may fit
 50 perfectly but return an estimator that converges to the wrong one or does not converge [29, particularly
 51 Sec. 8]. Identification is even more important for causal inference, because, unlike usual (non-causal)
 52 model misspecification, causal assumptions are often unverifiable through observables [50]. Thus, it
 53 is critical to specify theoretical conditions for identification, and then the applicability of methods
 54 can be judged by knowledge of an application domain.

55 In this work, we study identification (Sec. 3) and estimation (Sec. 4) of TEs under weak overlap. We
 56 particularly discuss individualized treatment effects, conditioned on the covariates. Code and proofs
 57 are in Supplementary Materials. The main **contributions** of this paper are:

- 58 1) theory of TE identification under weak overlap of covariates, using PGS and identifiable model;
- 59 2) counterfactual generalization bounds on TE error, which motivates our *conditional* BRL;
- 60 3) a new regularized VAE to estimate TEs, with connections to identification and balancing;
- 61 4) experimental comparison to state-of-the-art methods on (semi-)synthetic datasets.

62 2 Setup and motivation

63 2.1 Counterfactuals, treatment effects, and identification

64 Following [21], we introduce *potential outcomes* (POs, or counterfactual outcomes) $\mathbf{y}(t) \in \mathbb{R}^d, t \in$
 65 $\{0, 1\}$. $\mathbf{y}(t)$ is the outcome that would have been observed, if treatment value $t = t$ had been applied.
 66 Formally, this is the *consistency of counterfactuals*: $\mathbf{y} = \mathbf{y}(t)$ if $t = t$, or simply $\mathbf{y} = \mathbf{y}(t)$. We see
 67 $\mathbf{y}(t)$ as the hidden variables that give *factual* \mathbf{y} under *factual assignment* $t = t$. The *fundamental*
 68 *problem of causal inference* is that, for a unit under research, we could observe only one of $\mathbf{y}(0)$ or
 69 $\mathbf{y}(1)$, corresponding the treatment value applied. That is, “factual” refers to \mathbf{y} or t that is in principle
 70 *observable* in data, or statistical entities (e.g, estimators) built on them. We also observe relevant
 71 covariate(s) $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^m$, which is associated with individuals, with distribution $\mathcal{D} \sim p(\mathbf{x}, \mathbf{y}, t)$.
 72 Note, we use Roman fonts for random variables (e.g., \mathbf{x}) and italic for realization (e.g., \mathbf{x}).

73 The expected PO is denoted by $\mu_t(\mathbf{x}) = \mathbb{E}(\mathbf{y}(t)|\mathbf{x} = \mathbf{x})$, conditioned on $\mathbf{x} = \mathbf{x}$. The estimands in
 74 this work are the Conditional Average TE (CATE) and Average TE (ATE), defined respectively by

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

75 CATE is an *individual-level*, personalized, treatment effect, given highly discriminative covariate.

76 Standard results [38][16, Ch. 3] give sufficient conditions for identification under general setting.
 77 They are *Exchangeability*: $\mathbf{y}(t) \perp\!\!\!\perp t | \mathbf{x}$, and *Overlap*: $p(t|\mathbf{x}) > 0$ for any $\mathbf{x} \in \mathcal{X}$. Both are required for
 78 $t \in \{0, 1\}$. When t appears in a statement without quantification, we always mean “for both t ”. Often,
 79 *Consistency* is also listed, but, as above, it is better known as well-definedness of counterfactuals.
 80 Exchangeability means, just as in RCTs but additionally given \mathbf{x} , that there is no correlation between
 81 factual treatment t and counterfactual outcomes $\mathbf{y}(t)$. Overlap means that the supports of $p(\mathbf{x}|t = 0)$
 82 and $p(\mathbf{x}|t = 1)$ should be the same, and this ensures it is valid to condition on any (\mathbf{x}, t) .

83 We relax overlapped covariate in Sec. 3.2, to allow some non-overlapped values \mathbf{x} , i.e., covariate \mathbf{x} is
 84 *weakly overlapped*. In Sec. 2.2, we introduce a condition which gives exchangeability (and PGSs). In
 85 this paper, we also discuss overlap of variables other than \mathbf{x} (e.g. PGSs), and we give a definition for
 86 any random variable \mathbf{v} with support \mathcal{V} as following.

87 **Definition 1.** *Overlap* of \mathbf{v} means $p(t|\mathbf{v}) > 0$ for all $t \in \{0, 1\}, \mathbf{v} \in \mathcal{V}$. We also say \mathbf{v} is *overlapped*.
 88 If the condition is violated at some value \mathbf{v} , then \mathbf{v} is *non-overlapped* and \mathbf{v} is *weakly overlapped*.

89 **2.2 Prognostic scores**

90 Our method is motivated by PGSs [14], adapted as Pt-score and P-score in Definition 2 in this paper,
 91 related to balancing scores $\mathbf{b}(\mathbf{x})$, which is defined by $t \perp\!\!\!\perp \mathbf{x} | \mathbf{b}(\mathbf{x})$ [36]. The PPS $p(t = 1 | \mathbf{x})$ is a special
 92 case of this. Both are sufficient scores for identification; PGSs are sufficient statistics of *outcome*
 93 *predictors* and $\mathbf{b}(\mathbf{x})$ is for the treatment (see Appendix for details).

94 **Definition 2.** A *Pt-score* (PtS) is two functions $\mathbb{P}_t(\mathbf{x})$ ($t = 0, 1$) such that $\mathbf{y}(t) \perp\!\!\!\perp \mathbf{x} | \mathbb{P}_t(\mathbf{x})$. A PtS is
 95 called a *P-score* (PS) if $\mathbb{P}_0 = \mathbb{P}_1$.

96 Note that, a PtS is by definition two functions, thus overlapped $\mathbb{P}_t(\mathbf{x})$ means *both* $\mathbb{P}_0(\mathbf{x})$ and $\mathbb{P}_1(\mathbf{x})$
 97 are overlapped. **Why PtS (PGS)?** PtS is more applicable than balancing score $\mathbf{b}(\mathbf{x})$ under weak
 98 overlap. Overlapped $\mathbf{b}(\mathbf{x})$ implies overlapped \mathbf{x} , which in turn implies overlapped PGS [10]. Lower-
 99 dimensional than \mathbf{x} , PtS is likely more overlapped than \mathbf{x} , and, moreover, there is evidence that PtS
 100 *maximizes overlap* among all sufficient scores for ATE [9].

101 Below is a direct corollary of Proposition 5 in [14]. Both of PtS and PS give CATE, but, as we will
 102 see, PS is better as a *conditionally balanced* representation, since $\mathbb{P}_t(\mathbf{x}) \perp\!\!\!\perp t | \mathbf{x}$ only when $\mathbb{P}_0 = \mathbb{P}_1$.

103 **Proposition 1** (CATE by PtS). *If \mathbb{P}_t is a PtS, then CATE can be given by*

$$\mu_t(\mathbf{x}) = \mathbb{E}(\mathbb{E}(\mathbf{y}(t) | \mathbb{P}_t, \mathbf{x})) = \mathbb{E}(\mathbb{E}(\mathbf{y} | \mathbb{P}_t(\mathbf{x}), t = t)) = \int p(y | \mathbb{P}_t = \mathbb{P}_t(\mathbf{x}), t) y dy \quad (2)$$

104 With the knowledge of \mathbb{P}_t , we choose one of $\mathbb{P}_0, \mathbb{P}_1$ corresponding to the counterfactual outcome of
 105 interest. This ability of counterfactual assignment resolves the problem of non-overlap at \mathbf{x} .

106 PtSs exist under general settings when $\mathbf{y}(t)$ follows an additive noise model (ANM).

107 **(G1)** (ANM) the data generating process (DGP) for \mathbf{y} is $\mathbf{y} = \mathbf{f}^*(\mathbb{M}(\mathbf{x}), t) + \mathbf{e}$ where \mathbf{f}^*, \mathbb{M} are
 108 functions and \mathbf{e} is a zero-mean exogenous (external) noise.

109 The DGP defines $\mathbf{y}(t)$ by setting $t = t$ in the equation. And it also specifies how other variables
 110 causally affect \mathbf{y} . For example, \mathbf{x} affects \mathbf{y} through \mathbb{M} , so $\mathbb{M}(\mathbf{x})$ is the effect modifier [14], which is
 111 often components of \mathbf{x} affecting \mathbf{y} directly. Note **(G1)** also implies exchangeability given \mathbf{x} , through
 112 $\mathbf{y}(t) \perp\!\!\!\perp t | \mathbb{M}(\mathbf{x})$. ANMs are also commonly used in nonparametric regression methods for TEs [6].

113 Under **(G1)**, 1) $\mathbb{P}_t := \mathbf{f}_t^*(\mathbb{M}(\mathbf{x})) = \mu_t(\mathbf{x})^1$ is a PtS² but not PS, 2) \mathbb{M} is a PS (\mathbf{x} is a trivial PS), and
 114 3) $\mathbb{P} := (\mu_0(\mathbf{x}), \mu_1(\mathbf{x}))$ is a PS. We use the same symbol to denote a PtS and the random variable
 115 defined by it, when appropriate.

116 **3 Identification under generative prognostic model**

117 In Sec. 3.1, we introduce our generative prognostic model and VAE based
 118 on $p(\mathbf{y}, \mathbf{z} | \mathbf{x}, t)$ and prove identifiability of our model. In Sec. 3.2, we prove
 119 identification of CATEs, one of our main contributions. The theoretical
 120 analysis involves only our generative model (i.e., prior and decoder), but
 121 not encoder, because model identifiability is a property of *model*, and
 122 causal identification is about *DGP* and model. The encoder is involved
 123 in *estimation* which is studied in Sec. 4.

124 **3.1 Intact-VAE: model, architecture, and identifiability**

125 Generative models are useful to solve the inverse problem of recovering
 126 Pt-score. Our goal is to build a model that can be learned by VAE from
 127 observational data to obtain a PtS, or more ideally PS, via the latent
 128 variable \mathbf{z} . That is, a generative prognostic model.

129 With the above goal, the generative model of our VAE is built as

$$p(\mathbf{y}, \mathbf{z} | \mathbf{x}, t) = p(\mathbf{y} | \mathbf{z}, t) p(\mathbf{z} | \mathbf{x}, t). \quad (3)$$

130 The first factor is our decoder which models $p(\mathbf{y} | \mathbb{P}_t, t)$ in (2), and the
 131 second factor is our *conditional* prior which models $\mathbb{P}_t(\mathbf{x})$. Conditioning on \mathbf{x} in the joint model

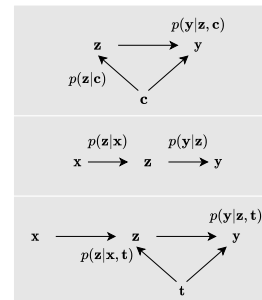


Figure 1: Graphical models of the decoders. From top: CVAE, iVAE, and Intact-VAE. The encoders are similar, taking all observables and build approximate posteriors, and thus are omitted.

¹We often write t of function argument in subscript, which indicates possible counterfactual assignment.

² μ_t is the most common PGS, to the extent that some call it *the* PGS (e.g. [39, 9, 47]), even without ANMs.

132 $p(\mathbf{y}, \mathbf{z}|\mathbf{x}, t)$ reflects that our estimand is CATE given \mathbf{x} . Modeling the score by a conditional distri-
 133 bution rather than a deterministic function is more flexible. We parameterize our model by ANM
 134 outcome and factorized Gaussian prior as

$$p_{\mathbf{f}}(\mathbf{y}|\mathbf{z}, t) = p_{\epsilon}(\mathbf{y} - \mathbf{f}_t(\mathbf{z})); p_{\lambda}(\mathbf{z}|\mathbf{x}, t) \sim \mathcal{N}(\mathbf{z}; \mathbf{h}_t(\mathbf{x}), \text{diag}(\mathbf{k}_t(\mathbf{x}))) \quad (4)$$

135 where $\theta = (\mathbf{f}, \mathbf{h}, \mathbf{k})$ are functional parameters and ϵ is a noise. $\lambda(\mathbf{x}) := \text{diag}(\mathbf{k}_t^{-1}(\mathbf{x}))(\mathbf{h}_t(\mathbf{x}), -\frac{1}{2})^T$
 136 is the natural parameter of the Gaussian prior, and we also use it as a shorthand for both \mathbf{h}, \mathbf{k} .

137 The ELBO of our model can be derived from standard variational lower bound

$$\begin{aligned} \log p(\mathbf{y}|\mathbf{x}, t) &\geq \log p(\mathbf{y}|\mathbf{x}, t) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}, \mathbf{y}, t) \| p(\mathbf{z}|\mathbf{x}, \mathbf{y}, t)) \\ &= \mathbb{E}_{\mathbf{z} \sim q} \log p(\mathbf{y}|\mathbf{z}, t) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}, \mathbf{y}, t) \| p(\mathbf{z}|\mathbf{x}, t)). \end{aligned} \quad (5)$$

138 Our encoder q , which conditions on all the observables, is standard, and we will see its importance
 139 later. We name this architecture *Intact-VAE* (*I*dentifiable *t*reatment-*c*onditional VAE).

140 We naturally have an identifiable conditional VAE (CVAE), as the name suggests. Note that (3) has a
 141 similar factorization with the generative model of iVAE [26], that is $p(\mathbf{y}, \mathbf{z}|\mathbf{x}) = p(\mathbf{y}|\mathbf{z})p(\mathbf{z}|\mathbf{x})$; the
 142 first factor does not depend on \mathbf{x} . Further, since we have the conditioning on t in both the factors of
 143 (3), our VAE architecture is a combination of iVAE and CVAE [42, 28], with t as the conditioning
 144 variable. See Figure 1 for the comparison in terms of graphical models. The core idea of iVAE is
 145 reflected in our model identifiability (Lemma 1 below). See Appendix for the basics of VAEs.

146 The following conditions on the model are used in theoretical analysis.

147 **(M1)** i) \mathbf{f}_t is injective, ii) \mathbf{f}_t is differentiable, and iii) $n := \dim(\mathbf{z}) = d (= \dim(\mathbf{y}))$.

148 **Lemma 1** (Model identifiability). *Given model (3) and (4) under (M1) i) and ii), for $t = t$, assume*

149 **(D1)** (Linear independence of λ) there exist $2n+1$ points $\mathbf{x}_0, \dots, \mathbf{x}_{2n} \in \mathcal{X}$ such that the $2n$ -square
 150 matrix $\mathbf{L}_t := [\gamma_{t,1}, \dots, \gamma_{t,2n}]$ is invertible, where $\gamma_{t,k} := \lambda_t(\mathbf{x}_k) - \lambda_t(\mathbf{x}_0)$.

151 *Then, given $t = t$, the family is identifiable up to an equivalence class. That is, if $p_{\theta}(\mathbf{y}|\mathbf{x}, t = t) =$
 152 $p_{\theta'}(\mathbf{y}|\mathbf{x}, t = t)$, we have the relation between parameters: for any \mathbf{y}_t in the image of \mathbf{f}_t ,*

$$\mathbf{f}_t^{-1}(\mathbf{y}_t) = \text{diag}(\mathbf{a})\mathbf{f}'_t^{-1}(\mathbf{y}_t) + \mathbf{b} =: \mathcal{A}(\mathbf{f}'_t^{-1}(\mathbf{y}_t)) \quad (6)$$

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

154 The conditions are inherited from iVAE. **(D1)** holds easily in practice, if the components of $\lambda_t(\mathbf{x})$
 155 are *linearly independent*; if **(D1)** fails, then the support of $\lambda_t(\mathbf{x})$ is in a $(2n - 1)$ -dimensional space.

156 The essence of the result is $\mathbf{f}'_t = \mathbf{f}_t \circ \mathcal{A}_t$, that is, \mathbf{f}_t can be identified (learned) up to an affine
 157 transformation defined by λ_t . This is achieved by combining the techniques from [26] and [43], and
 158 essentially the same results can be proved for other exponential family priors [43]. In this paper,
 159 symbol ' (prime) always indicates another parameter (variable, etc.).

160 3.2 Nonparametric identifications under weakly-overlapped covariate

161 In this subsection, we give two identification results based on (partial) recovery of PS or PtS,
 162 respectively. Since PtSs are functions of \mathbf{x} , the recovery is achieved by a noiseless prior, that is,
 163 $\mathbf{k}(\mathbf{x}) = \mathbf{0}$; the prior $\mathbf{z}_{\lambda,t} \sim p_{\lambda}(\mathbf{z}|\mathbf{x}, t = t)$ degenerates to deterministic function $\mathbf{h}_t(\mathbf{x})$.

164 PtSs with dimensionality lower than or equal to \mathbf{y} are essential to work for weak overlap of \mathbf{x} , to the
 165 extent that, from now on, we simply say PSs / PtSs when referring to this kind of *low-dimensional*
 166 *PSs / PtSs*, unless particularly indicated. **(M1)** iii), i.e. $n = d$, is not restrictive because μ_t is a PtS of
 167 the same dimension as \mathbf{y} under **(G1)**. Also, in practice, $n = d$ means that we seek a low-dimensional
 168 representation of \mathbf{x} . In fact, to make the dimensionality explicit in **(G1)**, we introduce an alternative
 169 **(G1')** which includes **(G1)** with $\mathbb{P}_t = \mu_t$ and \mathbf{j}_t is identity.

170 **(G1')** (Low-dimensional PtS) Under **(G1)**, $\mu_t(\mathbf{x}) = \mathbf{j}_t(\mathbb{P}_t(\mathbf{x}))$ for some \mathbb{P}_t and injective \mathbf{j}_t .

171 We use **(G1')** afterwards. Clearly, \mathbb{P}_t in **(G1')** is a PtS, and injectivity and $n = d$ ensure $n =$
 172 $\dim(\mathbf{y}) \geq \dim(\mathbb{P}_t)$. Similarly, the next **(G2)** reduces unverifiable $n \geq \dim(\mathbb{P})$ to $n = d$, for PS.

173 **(G2)** (Low-dimensional PS) Under **(G1)**, $\mu_t(\mathbf{x}) = \mathbf{j}_t(\mathbb{P}(\mathbf{x}))$ for some \mathbb{P} and injective \mathbf{j}_t .

174 **(G2)** means that CATEs are given by μ_0 and an invertible function $\mathbf{i} := \mathbf{j}_1 \circ \mathbf{j}_0^{-1}$. See Appendix for
 175 more discussion and a (closely related) real world example. In Sec. 4.1, we argue that *there often*
 176 *exist equivalent PSs under (G1')*, at least approximately.

177 With **(G1')** or **(G2)**, overlapped \mathbf{x} can be relaxed to overlapped P(t)S plus the following.

178 **(M2)** (Score partition preserving) For any $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$, $\mathbb{P}_t(\mathbf{x}) = \mathbb{P}_t(\mathbf{x}') \implies \mathbf{h}_t(\mathbf{x}) = \mathbf{h}_t(\mathbf{x}')$.

179 Note that **(M2)** is in fact required for optimal \mathbf{h} , in the sense specified in Proposition 1 and Theorem 1
 180 below. The intuition is that, \mathbb{P}_t maps non-overlapped \mathbf{x} to an overlapped value, and \mathbf{h}_t preserves this
 181 property, through learning. In fact, **(M2)** is trivially satisfied if \mathbb{P}_t and \mathbf{h}_t are *linear*, and this is still
 182 challenging and considered by many works [32, 9], or some with linear outcome models [11, 39].

183 Our first identification, Proposition 2, relies on **(G2)** and our generative model, *without* model
 184 identifiability (so differentiable \mathbf{f}_t is not needed). This is a nonparametric³ identification under shape
 185 restriction [7], because \mathbf{f}, \mathbf{h} are functional parameters, and injectivity is monotonicity if \mathbf{j}_t is on \mathbb{R} .

186 **Proposition 2** (Identification with PS). *Given **(G2)** and model (3) and (4) under **(M1)** i) and iii),
 187 and **(M3)** (PS matching) $\mathbf{h}_0(\mathbf{x}) = \mathbf{h}_1(\mathbf{x})$ and $\mathbf{k}(\mathbf{x}) = \mathbf{0}$. Then, if $\mathbb{E}_{p_\theta}(\mathbf{y}|\mathbf{x}, t) = \mathbb{E}(\mathbf{y}|\mathbf{x}, t)$, we have⁴*

- 188 1) (Recovery of PS) $\mathbf{z}_{\lambda, t} = \mathbf{h}_t(\mathbf{x}) = \mathbf{v}(\mathbb{P}(\mathbf{x}))$ on overlapped \mathbf{x} ,
 189 where $\mathbf{v} : \mathcal{P} \rightarrow \mathbb{R}^n$ is an injective function and $\mathcal{P} := \{\mathbb{P}(\mathbf{x})|\text{overlapped } \mathbf{x}\}$
 190 2) (Identification) if \mathbb{P} in **(G2)** is overlapped, and **(M2)** is satisfied, then $\mu_t(\mathbf{x}) = \hat{\mu}_t(\mathbf{x})$
 191 for any $t \in \{0, 1\}$, $\mathbf{x} \in \mathcal{X}$, where $\hat{\mu}_t(\mathbf{x}) := \mathbb{E}_{p_{\lambda}(\mathbf{z}|\mathbf{x}, t)}\mathbb{E}_{p_f}(\mathbf{y}|\mathbf{z}, t) = \mathbf{f}_t(\mathbf{h}_t(\mathbf{x}))$.

192 The essence is, i) the true DGP is identified up to an invertible mapping \mathbf{v} , so that $\mathbf{f}_t = \mathbf{j}_t \circ \mathbf{v}^{-1}$ and
 193 $\mathbf{h}_t = \mathbf{v} \circ \mathbb{P}_t$, and ii) \mathbb{P}_t is recovered up to \mathbf{v} and $\mathbf{y}(t) \perp\!\!\!\perp \mathbf{x} | \mathbb{P}_t$ is preserved, with *same* \mathbf{v} for both t .

194 PS is preferred since it satisfies overlap more easily and **(M2)** than PtS which refers to two functions.
 195 However, the existence of low-dimensional PS is uncertain in practice when our knowledge of the
 196 DGP is limited. Thus, we need Theorem 1 to work under PtS which generally exists.

197 **Theorem 1** (Identification with PtS). *Given the DGP **(G1')** and model (3)&(4) under **(M1)** and
 198 **(M3')** (Noise matching) $p_\epsilon = p_e$ and $\mathbf{k}(\mathbf{x}) = k\mathbf{k}'(\mathbf{x}), k \rightarrow 0$, assume **(D1)** and*

- 199 **(D2)** (Spontaneous balance) There exist $2n + 1$ points $\mathbf{x}_0, \dots, \mathbf{x}_{2n} \in \mathcal{X}$, $2n$ -square matrix \mathbf{C} ,
 200 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
 201 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
 202 $\alpha_t(\mathbf{x}; \lambda_t)$ is the log-partition function of the prior in (4).

203 *Then, if $p_\theta(\mathbf{y}|\mathbf{x}, t) = p(\mathbf{y}|\mathbf{x}, t)$, conclusions 1) and 2) in Proposition 2 hold with \mathbb{P} replaced with \mathbb{P}_t
 204 in **(G1')**, and the domain of \mathbf{v} becomes $\mathcal{P} := \mathcal{P}_0 \cup \mathcal{P}_1, \mathcal{P}_t := \{\mathbb{P}_t(\mathbf{x})|\text{overlapped } \mathbf{x}\}$.*

205 Theorem 1 also achieves the two essential points, but in different and complementary ways. Proposi-
 206 tion 2 starts from the prior by $\mathbb{P}_0 = \mathbb{P}_1$ and setting $\mathbf{h}_0 = \mathbf{h}_1$. Conversely, Theorem 1 starts from the
 207 decoder with $p_\epsilon = p_e$ and strengthens model identifiability (6) by **(D2)**. **(D2)** restricts the discrepancy
 208 between λ_0, λ_1 on $2n + 1$ points of \mathbf{x} , thus is relatively easy to satisfy with high-dimensional \mathbf{x} .

209 We see more reasons to prefer PS. In general, to identify the mean function $\mu_t(\mathbf{x})$, a regression is
 210 enough, and $p_\epsilon = p_e$ is unnecessary as in Proposition 2. Also, **(D2)** is trivial if we have PS and set
 211 $\lambda_0 = \lambda_1$. See Appendix for more on the complementarity between the two identifications.

212 4 Estimation by β -Intact-VAE

213 4.1 Prior as PS, posterior as PtS, and β as regularization strength

214 In this subsection, we discuss our focus on balanced PtSs and give an estimation method to realize it.

215 In learning the Intact VAE with data, we assume that there is a PtS and the decomposition of **(G1')**
 216 holds. Such a decomposition is not unique in general, however. Among possible PtSs, we wish to
 217 learn a balanced PtS, which is close to PS. This is based on the observations in Sec. 3.2: we saw
 218 that existence of PS is preferable in identifying the true DGP up to equivalent expression. Here, we
 219 introduce the notion of balanced PtS in a non-rigorous way: a PtS \mathbb{P}_t is called *balanced* if the value of
 220 some measure for the conditional independence $\mathbb{P}_t(\mathbf{x}) \perp\!\!\!\perp t | \mathbf{x}$ is small. The idea is common in practice.
 221 For example, in a real-world nutrition study, [20] reduces 11 covariates to a 1-dimensional linear PS.

³Some references (e.g. [16]) only refer to identification without models as “nonparametric”. However, with recent advances, we have identification under nonparametric models [29], also the case in this paper.

⁴Same results hold *without* “under **(G1)**” in **(G2)** (or **(G1')** for Theorem 1), except that $\mathbf{z}_{\lambda, t}$ is *not* necessarily a PS (or PtS for Theorem 1). This is because **(G1)** is required to ensure μ_t is a PtS.

222 Assuming that there is a balanced PtS, we consider two ways for estimating it with Intact VAE. One
 223 is to exploit a prior that does not depend on t . Namely, we set $\lambda_0 = \lambda_1 =: \lambda$ in (4). The other is
 224 to introduce a hyperparameter β in the ELBO as in β -VAE [17]. More specifically, with the prior
 225 $p_\lambda(\mathbf{z}|\mathbf{x})$, we use factorized Gaussian for the decoder and encoder:

$$p_{f,g}(\mathbf{y}|\mathbf{z}, t) \sim \mathcal{N}(\mathbf{y}; \mathbf{f}_t(\mathbf{z}), \text{diag}(\mathbf{g}_t(\mathbf{z}))); \quad q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}, t) \sim \mathcal{N}(\mathbf{z}; \mathbf{r}_t(\mathbf{x}, \mathbf{y}), \text{diag}(\mathbf{s}_t(\mathbf{x}, \mathbf{y}))). \quad (7)$$

226 The modified ELBO with β , up to additive constant, is derived as

$$\mathbb{E}_{\mathcal{D}}\{-\beta D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}, t)||p_\lambda(\mathbf{z}|\mathbf{x})) - \mathbb{E}_{\mathbf{z}\sim q}[(\mathbf{y} - \mathbf{f}_t(\mathbf{z}))^2/2\mathbf{g}_t^2(\mathbf{z})] - \mathbb{E}_{\mathbf{z}\sim q} \log |g_t(\mathbf{z})|\}. \quad (8)$$

227 Here, for convenience, we omit the summation (also in \mathcal{L}_f in Sec. 4.2), as if \mathbf{y} was univariate. The
 228 approximate posterior (or encoder) $q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}, t)$ depends on t , which can realize a PtS. With β , we
 229 control the trade-off between the first and second term: the former is the divergence of the posterior
 230 from the balanced prior, and the latter is the reconstruction of the outcome. By choosing β in an
 231 appropriate way, such as by validation, the ELBO can find a solution that explains the outcome while
 232 keeping the balancedness of the posterior. Also note that, the parameters \mathbf{g} and \mathbf{k} , which models the
 233 outcome noise and expresses uncertainty of the prior, respectively, are both learned by the ELBO.
 234 This deviates from the theoretical conditions in Sec. 3.2, but is more practical and gives better results
 235 in experiments. See Appendix for much more on ideas and connections behind the ELBO.

236 Once the encoder q_ϕ is learned⁵, the estimate of the expected POs is given by

$$\hat{\mu}_{\hat{t}}(\mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x}=\mathbf{x})} \mathbf{f}_{\hat{t}}(\mathbf{z}) = \mathbb{E}_{\mathcal{D}|\mathbf{x}\sim p(\mathbf{y}, t|\mathbf{x})} \mathbb{E}_{\mathbf{z}} \mathbf{f}_{\hat{t}}(\mathbf{z}) q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}, t), \hat{t} \in \{0, 1\}, \quad (9)$$

237 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 and $\mathcal{D} \sim p(\mathbf{x}, \mathbf{y}, t)$. In estimation,
 238 we consider the case where \mathbf{x} is observed in the data, and the sample of (\mathbf{y}, t) are taken from the
 239 data given $\mathbf{x} = \mathbf{x}$ (when \mathbf{x} is not in the data, we replace q_ϕ with p_λ in (9), see Appendix for details
 240 and results). Note that \hat{t} in (9) indicates counterfactual assignment which may not be the same as
 241 factual $t = t$ in the data. That is, we set $t = \hat{t}$ in the decoder. The assignment is not applied to the
 242 encoder, but it is learned from factual \mathbf{x}, \mathbf{y} (see also Sec. 4.2, the explanation for $\epsilon_{CF,t}$). The overall
 243 **algorithm** steps are i) we train VAE by (8), ii) infer CATE $\hat{\tau}(\mathbf{x}) = \hat{\mu}_1(\mathbf{x}) - \hat{\mu}_0(\mathbf{x})$ by (9).

244 4.2 Conditional balanced representation learning

245 We formally justify our ELBO (8) from the viewpoint of BRL. Usually, particularly in ATE estimation,
 246 balance means covariate balance, i.e., $\mathbf{x} \perp\!\!\!\perp t$ [44]. Influenced by this, most BRL methods learn balanced
 247 covariate representation \mathbf{z} such that $\mathbf{z} \perp\!\!\!\perp t$ [40, 31] and usually \mathbf{z} is a function of \mathbf{x} . From Sec. 4.1, we
 248 understand that larger β in ELBO (8) encourages $\mathbf{z} \perp\!\!\!\perp t|\mathbf{x}$ which is given by the prior, corresponding
 249 to $\mathbb{P}_t(\mathbf{x}) \perp\!\!\!\perp t|\mathbf{x}$ for a balanced PtS. Here, we show that, this *conditional balance* of representation \mathbf{z}
 250 is natural for CATE estimation, and CATE error due to bad recovery of $\hat{\mathbf{j}}_t$ in $(\mathbf{G1}')$ is controlled by
 251 ELBO (8). In Appendix, we detail novel implications of our bounds, compared to those in [40, 31].

252 Using (9) to estimate CATE, $\hat{\tau}_f(\mathbf{z}) = \mathbf{f}_1(\mathbf{z}) - \mathbf{f}_0(\mathbf{z})$ is marginalized on $q(\mathbf{z}|\mathbf{x})$. The bounds below
 253 motivate both prior and posterior balancing. Let us first consider errors defined by the aggregated
 254 prior $p(\mathbf{z}|\mathbf{x}) := \sum_t p(t|\mathbf{x}) p_t(\mathbf{z}|\mathbf{x})$ (denote $p_t(\mathbf{z}|\mathbf{x}) := p_\lambda(\mathbf{z}|\mathbf{x}, t = t)$), and $q(\mathbf{z}|\mathbf{x})$ afterwards. We
 255 introduce the objective we bound. The *true CATE*, given covariate $\mathbf{x} = \mathbf{x}$ or score \mathbf{z} , is

$$\tau(\mathbf{x}) = m_1(\mathbb{P}_1(\mathbf{x})) - m_0(\mathbb{P}_0(\mathbf{x})); \quad \tau_m(\mathbf{z}) = m_1(\mathbf{z}) - m_0(\mathbf{z}) \quad (10)$$

256 where $m_t(\mathbf{z}) := \mathbb{E}(\mathbf{y}(t)|\mathbb{P}_t = \mathbf{z})$ and \mathbb{P}_t is a balanced PtS in $(\mathbf{G1}')$. Accordingly, given \mathbf{x} , the *error*
 257 *of prior CATE*, with or without knowing \mathbb{P}_t , is naturally defined as

$$\epsilon_f^{*,p}(\mathbf{x}) := \mathbb{E}_{p(\mathbf{z}|\mathbf{x})} (\hat{\tau}_f(\mathbf{z}) - \tau(\mathbf{x}))^2; \quad \epsilon_f^p(\mathbf{x}) := \mathbb{E}_{p(\mathbf{z}|\mathbf{x})} (\hat{\tau}_f(\mathbf{z}) - \tau_m(\mathbf{z}))^2. \quad (11)$$

258 We bound ϵ_f^p instead of $\epsilon_f^{*,p}$ because the error between $\tau(\mathbf{x})$ and $\tau_m(\mathbf{z})$ is small if balanced \mathbb{P}_t is
 259 recovered (then $\mathbf{z} \approx \mathbb{P}_0(\mathbf{x}) \approx \mathbb{P}_1(\mathbf{x})$ in (10), see Appendix for details). Instead, we consider the error
 260 between $\hat{\tau}_f$ and τ_m below. We define the risks of outcome regression, into which ϵ^p is decomposed.

261 **Definition 3 (PO Risks).** The *expected loss of PO* at (\mathbf{z}, t) , *factual risk*, and *counterfactual risk* are

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

⁵As usual, we expect variational inference and optimization procedure are (near) optimal, i.e., Consistency of VAE (see Appendix for formal statement). Consistent estimation using the prior is a direct corollary of consistent VAE. Under Gaussian models, it is possible to prove consistency of posterior estimation, as shown in [4].

262 With $\mathbf{y}(t)$ involved, \mathcal{L}_f is a PO error of f weighted by g . Factual and counterfactual counterparts are
 263 defined accordingly, w.r.t factual p_t learned from data. Note, in $\epsilon_{F,t}$, unit $\mathbf{u} = (\mathbf{x}, \mathbf{y}, t)$ involves in
 264 the learning of $p_t(\mathbf{z}|\mathbf{x})$ (and $q_t(\mathbf{z}|\mathbf{x})$ afterwards), and also in $\mathcal{L}_f(\mathbf{z}, t)$ since $\mathbf{y}(t) = \mathbf{y}$ for the unit. In
 265 $\epsilon_{CF,t}$, however, $\mathbf{y}(t) \neq \mathbf{y}' = \mathbf{y}(1-t)$ for $\mathbf{u}' = (\mathbf{x}, \mathbf{y}', 1-t)$ (particularly relevant for the posterior).

266 Thus, *the regression error (second) term in ELBO (8) controls $\epsilon_{F,t}^p$ via factual data*. On the other
 267 hand, $\epsilon_{CF,t}^p$ is not estimable due to unobservable $\mathbf{y}(1-t)$, but is bounded as below.

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

$$\epsilon_{CF}^p(\mathbf{x}) \leq \sum_t p(1-t|\mathbf{x})\epsilon_{F,t}^p(\mathbf{x}) + M\mathbb{D}^p(\mathbf{x}) \quad (12)$$

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

270 $\epsilon_{CF}^p(\mathbf{x})$ is bounded by $\epsilon_{F,t}^p$ plus $M\mathbb{D}^p(\mathbf{x})$, which measures the imbalance between $p_t(\mathbf{z}|\mathbf{x})$ and is
 271 symmetric for t . We can implicitly control ϵ_{CF}^p by making \mathbb{D}^p small. Again, this means PS is preferred
 272 as a conditional balanced representation, and justifies our balanced prior $p_\lambda(\mathbf{z}|\mathbf{x})$. Moreover, *the*
 273 *results (including Theorem 2 below) also hold for posterior estimation*, that is, replace $p_t(\mathbf{z}|\mathbf{x})$ with
 274 $q_t(\mathbf{z}|\mathbf{x}) := q(\mathbf{z}|\mathbf{x}, t) = \mathbb{E}_{p(\mathbf{y}|\mathbf{x}, t)} q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}, t)$, the results and proofs hold as it was. This implies
 275 that the imbalance between q_t should also be controlled. Correspondingly, *the symmetric KL term in*
 276 *ELBO (8) balances $q_t(\mathbf{z}|\mathbf{x})$ by encouraging $\mathbf{z} \perp\!\!\!\perp t|\mathbf{x}$ for the posterior.*

277 Theorem 2 in turn bounds ϵ_f^p , by decomposing it to $\epsilon_{F,t}^p$, $\epsilon_{CF,t}^p$, and \mathbb{V}_y^p .

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

$$\epsilon_f(\mathbf{x}) \leq 2(G^2(\epsilon_{F,0}(\mathbf{x}) + \epsilon_{F,1}(\mathbf{x})) + M\mathbb{D}(\mathbf{x}) - \mathbb{V}_y(\mathbf{x}))p \quad (13)$$

279 where $\mathbb{V}_y^p(\mathbf{x}) := \mathbb{E}_{p(\mathbf{z}|\mathbf{x})} \sum_t \mathbb{E}_{p(\mathbf{y}(t)|\mathbb{P}_t=\mathbf{z})} (\mathbf{y}(t) - m_t(\mathbf{z}))^2$, and “ $|^p$ ” collects all superscripts q .

280 The new term $\mathbb{V}_y^p(\mathbf{x})$ reflects the intrinsic variance in the DGP and can not be controlled, and it is
 281 negative because ϵ_f^p is defined by mean functions f_t and m_t , not $\mathbf{y}(t)$. The other two terms, as we
 282 indicated, is estimated by our ELBO.

283 Estimating G , M is nontrivial. Instead, similarly to [40], we rely on β in the ELBO to weight the
 284 two terms in (13). We do not need two hyperparameters since G is *implicitly controlled by the third*
 285 *term in ELBO (8)*, which is a norm constraint. As in matching methods, β is a trade-off between
 286 conditional balance of learned PtS (affected by f_t) and precision / effective sample size of outcome
 287 regression, and can be seen as the probabilistic counterpart of [47, 25].

288 Finally, we note the bounds do not directly address non-overlap; in Lemma 2, when $p(1-t|\mathbf{x}) = 0$,
 289 $\epsilon_{F,1-t}^p$ in the r.h.s is unbounded since $p_{1-t}(\mathbf{z}|\mathbf{x})$ can not be learned from data. However, as we argued
 290 in Sec. 3.2, with more balanced \mathbb{P}_t recovered as representation, overlap is more easily satisfied.

291 5 Related work

292 **Weak overlap.** Under (respective versions of) weak overlap, [32] estimates ATE by reducing
 293 covariates to a linear PGS, [11] estimates homogeneous (constant) TE under partial linear outcome
 294 model, and [9] studies identification of ATE by a general class of scores, given (linear) PPS and PGS.
 295 In machine learning, current focus is on finding overlap regions [33, 8], or indicating possible failure
 296 under weak overlap [22], but not remedies. An exception is [24] which provides bounds without
 297 overlap. [40, 31] are in line of [24] and have similar bounds to ours, without relating to overlap. Our
 298 method is the *first* in machine learning that gives identification without overlap.

299 **Prognostic scores** are recently combined with machine learning, mainly in biostatistics. For example,
 300 [39] trains a flexible PGS and fits a linear regression on the PGS among others, for constant TEs, and
 301 [12] models PGS in its Bayesian regression tree for CATE. More related, [20] estimates CATE by
 302 reducing covariates to a linear score that is a joint PPS and PGS, and [47] uses SVM to minimize
 303 the worst-case bias due to PGS imbalance. However, in machine learning, few methods consider
 304 PGSs; [55, 15] learn outcome predictors, without connection to PGS, while [24] conceptually, but
 305 not formally, connects BRL to PGS. Our work follows the recent boom in biostatistics and is the *first*
 306 to formally connect generative learning, PGS, and BRL (see below on BRL) for TE estimation.

307 **Identifiable representation.** Recently, independent component analysis (ICA) and representation
 308 learning, both ill-posed inverse problems, meet together to give nonlinear ICA and identifiable

309 representation, e.g., using VAEs [26], and energy models [27]. The results are exploited in causal
 310 discovery [51] and out-of-distribution generalization [46]. This work is the *first* to explore identifiable
 311 representations in TE identification.

312 **BRL and related methods** amount to a major direction. Early BRL methods are BLR/BNN [23]
 313 and TARnet/CFR [40]. Adding to this, [53] also exploits the local similarity of between data points.
 314 [41] uses similar architecture to TARnet, considering the importance of treatment probability. There
 315 are also methods using GAN [54, GANITE] and Gaussian process [1]. Our method shares the idea of
 316 BRL, and further extends to conditional balancing more suitable for CATE.

317 Our work hopefully lays conceptual and theoretical foundations of **VAE methods for TEs** (e.g.,
 318 [30, 31]), under unconfoundedness. Also, **monotonicity**, which is injectivity on \mathbb{R} , is important in
 319 causal inference, and some works consider it together with overlap [24, 56]. See Appendix for details.

320 6 Experiments

321 We compare the proposed method with existing methods on three types of datasets. Here we present
 322 two experiments, and the rest one, on the Pokec social network dataset, can be found in Appendix. As
 323 in previous works [40, 30], we report the absolute error of ATE $\epsilon_{ATE} := |\mathbb{E}_{\mathcal{D}}(y(1) - y(0)) - \mathbb{E}_{\mathcal{D}}\hat{\tau}(\mathbf{x})|$,
 324 and, as a surrogate of CATE, the empirical PEHE [18] $\epsilon_{PEHE} := \mathbb{E}_{\mathcal{D}}((y(1) - y(0)) - \hat{\tau}(\mathbf{x}))^2$.

325 Unless otherwise indicated, for each function f, h, k, r, s in (4)(7), we use a multilayer perceptron
 326 (MLP) that has 3×200 hidden units with ReLU activation, and $\lambda = (h, k)$ depends only on \mathbf{x} . We fix
 327 $g(\mathbf{x}) = 1$ because the datasets have fixed noise scale, and results with learned g on synthetic dataset
 328 with dependent noise is in Appendix. The Adam optimizer with initial learning rate 10^{-4} and batch
 329 size 100 is employed. All experiments use early-stopping of training by evaluating the ELBO on a
 330 validation set, and results are reported on a testing set. Each set of running on synthetic dataset (a line
 331 in the figure) is within 1 hour on an 8-CPU machine, and it is within a day for IHDP. More details on
 332 hyper-parameters and settings are given in each experiment and Appendix.

333 6.1 Synthetic dataset

334 We generate synthetic datasets following (14). Both \mathbf{x}, \mathbf{w} are factorized Gaussians. μ, σ are randomly
 335 sampled. The functions h, k, l are linear. Outcome models f_0, f_1 are built by NNs with invertible
 336 activations. y is univariate, $\dim(\mathbf{x}) = 30$, and $\dim(\mathbf{w})$ ranges from 1 to 5. \mathbf{w} is a PS, but is *not*
 337 *low-dimensional* when $\dim(\mathbf{w}) > 1$. We *control overlap* by ω which multiplies the logit value, and
 338 have 5 different overlap levels from strong overlap to very weak overlap. See Appendix for details.

$$\mathbf{x} \sim \mathcal{N}(\mu, \sigma); \mathbf{w} | \mathbf{x} \sim \mathcal{N}(h(\mathbf{x}), k(\mathbf{x})); t | \mathbf{x} \sim \text{Bern}(\text{Logit}(\omega l(\mathbf{x}))); y | \mathbf{w}, t \sim \mathcal{N}(f_t(\mathbf{w}), 1). \quad (14)$$

339 With the same $(\dim(\mathbf{w}), \omega)$, we evaluate our method and
 340 CFR on 10 random DGPs, with different sets of functions
 341 f, h, k, l in (14). For each DGP, we sample 1500 data
 342 points, and split them into 3 equal sets for training, val-
 343 idation, and testing. We show our results for different
 344 hyperparameter β . For CFR, we try different balancing
 345 parameters and present the best results (see Appendix for
 346 details). We report ϵ_{PEHE} , see Appendix for ATE results.

347 In each panel of Figure 2, we adjust one of $\omega, \dim(\mathbf{w})$
 348 respectively, with the other fixed to the lowest. In the left
 349 panel, both our method and CFR are quite robust to overlap
 350 level, supporting respective theories ([24] gives bounds for
 351 CFR under weak overlap). Too large β seems to worsen
 352 the results, possibly because we already have an apparent
 353 PS (\mathbf{w} with $\dim(\mathbf{w}) = 1$) and large β incurs sub-optimal ELBO with no gain.

354 In the right panel, when $\dim(\mathbf{w}) > 1$, f_t in (14) is non-injective and learning of PtS is necessary.
 355 Thus, larger β has a negative effect, and particularly, $\beta = 1$ is significantly better than $\beta = 3$. The
 356 drop of error for $\dim(\mathbf{w}) > 3$ is due to the randomness of f in (14). In Sec. 2.2, we saw that the
 357 2-dimensional PS $\mathbb{P} := (\mu_0(\mathbf{x}), \mu_1(\mathbf{x}))$ always exists under ANMs. Thus, when $\dim(\mathbf{w}) > 2$, our
 358 method tries to recover that \mathbb{P} , and generally performs not worse than under $\dim(\mathbf{w}) = 2$, but still
 359 not better than under $\dim(\mathbf{w}) = 1$.

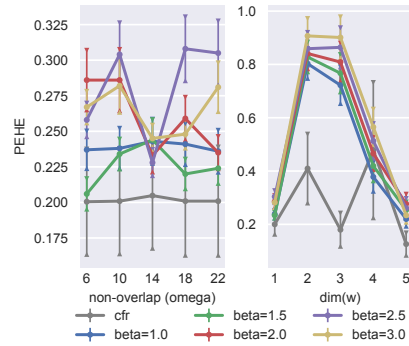


Figure 2: $\sqrt{\epsilon_{PEHE}}$ on synthetic dataset. Error bar on 10 random DGPs.

360 Our method is much more robust against different DGPs than
 361 CFR (see error bars), though it is worse than CFR when
 362 $\dim(\mathbf{w}) > 1$. This is unsurprising because our model
 363 has *1-dimensional* \mathbf{z} , while CFR uses 200-dimensional
 364 representation. Thus, the results already show the power of
 365 identification and recovery of scores (see Figure 3 also). In
 366 fact, we observed that our method outperforms or matches
 367 CFR with higher-dimensional \mathbf{z} (see Appendix). Thus, we
 368 believe the performance gap with $\dim(\mathbf{z}) = 1$ is due to the

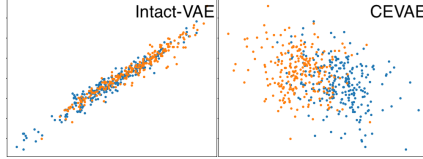


Figure 3: Plots of recovered - true latent. Blue: $t = 0$, Orange: $t = 1$.

capacity of NNs in Intact-VAE.
 369 When $\dim(\mathbf{w}) = 1$, there are no better PSs than \mathbf{w} , because f_t is invertible and no information
 370 can be dropped from \mathbf{w} . Thus, as shown in Figure 3, our method learns \mathbf{z} as an approximate affine
 371 transformation of the true \mathbf{w} , showing identification. For comparison, we run [30, CEVAE] which
 372 is also based on VAE but without identification, and it shows much lower quality of recovery. As
 373 expected, both recovery and estimation are better with balanced prior $p_\lambda(\mathbf{z}|\mathbf{x})$, and we can see an
 374 example of bad recovery using $p_\lambda(\mathbf{z}|\mathbf{x}, t)$ in Appendix. More latent plots can also be found there.

375 6.2 IHDP benchmark dataset

376 This experiment shows our conditional balancing matches state-of-the-art BRL methods, *and does not*
 377 *overfit to PEHE*. The IHDP dataset [18] is widely used to evaluate machine learning based methods,
 378 e.g. [40, 41]. It is also used in [24] which considers weak overlap, because *the covariates are weakly*
 379 *overlapped* due to their correlation to the artificial treatment assignment. Finally, there is a linear PS
 380 (linear combination of the covariates). See Appendix for details.

381 Note, most of covariates are binary, so the support of the PS is often on small and separated intervals
 382 and is possibly discrete. Thus, Gaussian latent \mathbf{z} is misspecified. We use multivariate \mathbf{z} in model
 383 to address this, similarly to [30]. We set $\beta = 1$ since it works well on synthetic dataset with weak
 384 overlap. To see our balancing property clearly, we modify our method and add two components for
 385 unconditional balancing from [40] (see Appendix), and compare this modified version to the original.

Table 1: Errors on IHDP over 1000 random DGPs. We report results with $\dim(\mathbf{z}) = 10$. **Bold** indicates method(s) that are *significantly* better. The results are taken from [40], except GANITE [54] and CEVAE [30].

Method	TMLE	BNN	CFR	CF	CEVAE	GANITE	Ours	Ours Mod.
ϵ_{ATE}	.30 \pm .01	.37 \pm .03	.25 \pm .01	.18 \pm .01	.34 \pm .01	.43 \pm .05	.178 \pm .006	.167 \pm .005
$\sqrt{\epsilon_{PEHE}}$	5.0 \pm .2	2.2 \pm .1	.71 \pm .02	3.8 \pm .2	2.7 \pm .1	1.9 \pm .4	.859 \pm .033	.777 \pm .026

386 As shown in Table 1, Intact-VAE outperforms or matches the state-of-the-art methods. Particularly,
 387 our method has the *best* ATE estimation, and is slightly worse than CFR for PEHE. This is possibly
 388 due to the fitting capacity (recall Sec. 6.1), and also we do *not* tune β . Notably, our method
 389 outperforms other generative models (CEVAE and GANITE) by large margins. The modified
 390 version is slightly improved, but we should note that the improvement for ϵ_{ATE} is barely significant.
 391 This indicates overfitting to PEHE. In fact, PEHE estimates the *marginalized* error $\mathbb{E}\epsilon(\mathbf{x})$ where
 392 $\epsilon(\mathbf{x}) = (\tau(\mathbf{x}) - \hat{\tau}(\mathbf{x}))^2$, and, compared with ϵ_{ATE} , it focuses on values \mathbf{x} with high probability and
 393 / or large $\epsilon(\mathbf{x})$. The balancing in [40] is based on bounding $\mathbb{E}\epsilon(\mathbf{x})$, and thus tends to overly focus on
 394 the above values of \mathbf{x} , resulting in sub-optimal estimation of CATE and even of ATE. This tendency
 395 is more apparent with sub-optimal hyperparameter for the unconditional balancing (see Appendix).

396 7 Conclusion

397 In this work, we proposed a method for CATE estimation, under weak overlap. Our method exploits
 398 identifiable VAE, a recent advance in generative models, and is fully motivated and theoretically
 399 justified by causal considerations: identification, PGS, and balancing. We show that VAEs are
 400 suitable for *principled* causal inference thanks to its probabilistic nature, if not compromised by ad
 401 hoc heuristics. We believe it is possible to extend the bounds in Sec. 4.2 to weak overlap, just as
 402 [24] extends [40] to weak overlap, and leave this for future. Experiments show evidence that the
 403 injectivity of \mathbf{f} in our model is possibly unnecessary because $\dim(\mathbf{z}) > \dim(\mathbf{y})$ often gives better
 404 results. Theoretical study of this is an interesting future direction. To avoid potential negative societal
 405 impact (e.g. bad prescriptions), practitioners should judge the conditions of the proposed method by
 406 their domain expertise, and careful trials are always recommended.

407 **References**

- 408 [1] Ahmed M Alaa and Mihaela van der Schaar. Bayesian inference of individualized treatment
409 effects using multi-task gaussian processes. In *Advances in Neural Information Processing*
410 *Systems*, pages 3424–3432, 2017.
- 411 [2] Joseph Antonelli, Matthew Cefalu, Nathan Palmer, and Denis Agniel. Doubly robust matching
412 estimators for high dimensional confounding adjustment. *Biometrics*, 74(4):1171–1179, 2018.
- 413 [3] Timothy B. Armstrong and Michal Kolesár. Finite-sample optimal estimation and inference on
414 average treatment effects under unconfoundedness. *arXiv preprint arXiv:1712.04594v5*, 2021.
- 415 [4] Stéphane Bonhomme and Martin Weidner. Posterior average effects. *arXiv preprint*
416 *arXiv:1906.06360v5*, 2021.
- 417 [5] M Alan Brookhart, Sebastian Schneeweiss, Kenneth J Rothman, Robert J Glynn, Jerry Avorn,
418 and Til Stürmer. Variable selection for propensity score models. *American journal of epidemi-*
419 *ology*, 163(12):1149–1156, 2006.
- 420 [6] Alberto Caron, Ioanna Manolopoulou, and Gianluca Baio. Estimating individual treatment
421 effects using non-parametric regression models: a review. *arXiv preprint arXiv:2009.06472*,
422 2020.
- 423 [7] Denis Chetverikov, Andres Santos, and Azeem M Shaikh. The econometrics of shape restrictions.
424 *Annual Review of Economics*, 10:31–63, 2018.
- 425 [8] Wangzhi Dai and Collin M Stultz. Quantifying common support between multiple treatment
426 groups using a contrastive-vae. In *Machine Learning for Health*, pages 41–52. PMLR, 2020.
- 427 [9] Alexander D’Amour and Alexander Franks. Deconfounding scores: Feature representations for
428 causal effect estimation with weak overlap. *arXiv preprint arXiv:2104.05762*, 2021.
- 429 [10] Alexander D’Amour, Peng Ding, Avi Feller, Lihua Lei, and Jasjeet Sekhon. Overlap in
430 observational studies with high-dimensional covariates. *Journal of Econometrics*, 2020.
- 431 [11] Max H Farrell. Robust inference on average treatment effects with possibly more covariates
432 than observations. *Journal of Econometrics*, 189(1):1–23, 2015.
- 433 [12] P Richard Hahn, Jared S Murray, Carlos M Carvalho, et al. Bayesian regression tree models
434 for causal inference: Regularization, confounding, and heterogeneous effects (with discussion).
435 *Bayesian Analysis*, 15(3):965–1056, 2020.
- 436 [13] David Hajage, Yann De Rycke, Guillaume Chauvet, and Florence Tubach. Estimation of
437 conditional and marginal odds ratios using the prognostic score. *Statistics in medicine*, 36(4):687–
438 716, 2017.
- 439 [14] Ben B Hansen. The prognostic analogue of the propensity score. *Biometrika*, 95(2):481–488,
440 2008.
- 441 [15] Negar Hassanpour and Russell Greiner. Learning disentangled representations for counterfactual
442 regression. In *International Conference on Learning Representations*, 2019.
- 443 [16] Miguel A. Hernan and James M. Robins. *Causal Inference: What If*. CRC Press, 1st edition,
444 2020.
- 445 [17] Irina Higgins, Loïc Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick,
446 Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a con-
447 strained variational framework. In *5th International Conference on Learning Representations*,
448 2017.
- 449 [18] Jennifer L Hill. Bayesian nonparametric modeling for causal inference. *Journal of Computa-*
450 *tional and Graphical Statistics*, 20(1):217–240, 2011.
- 451 [19] Han Hong, Michael P Leung, and Jessie Li. Inference on finite-population treatment effects
452 under limited overlap. *The Econometrics Journal*, 23(1):32–47, 2020.

- 453 [20] Ming-Yueh Huang and Kwun Chuen Gary Chan. Joint sufficient dimension reduction and
454 estimation of conditional and average treatment effects. *Biometrika*, 104(3):583–596, 2017.
- 455 [21] Guido W Imbens and Donald B Rubin. *Causal inference in statistics, social, and biomedical*
456 *sciences*. Cambridge University Press, 2015.
- 457 [22] Andrew Jesson, Sören Mindermann, Uri Shalit, and Yarin Gal. Identifying causal-effect
458 inference failure with uncertainty-aware models. *Advances in Neural Information Processing*
459 *Systems*, 33, 2020.
- 460 [23] Fredrik Johansson, Uri Shalit, and David Sontag. Learning representations for counterfactual
461 inference. In *International conference on machine learning*, pages 3020–3029, 2016.
- 462 [24] Fredrik D Johansson, Uri Shalit, Nathan Kallus, and David Sontag. Generalization bounds and
463 representation learning for estimation of potential outcomes and causal effects. *arXiv preprint*
464 *arXiv:2001.07426*, 2020.
- 465 [25] Nathan Kallus, Brenton Pennicooke, and Michele Santacatterina. More robust estimation of
466 sample average treatment effects using kernel optimal matching in an observational study of
467 spine surgical interventions. *arXiv preprint arXiv:1811.04274*, 2018.
- 468 [26] Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational au-
469 toencoders and nonlinear ica: A unifying framework. In *International Conference on Artificial*
470 *Intelligence and Statistics*, pages 2207–2217, 2020.
- 471 [27] Ilyes Khemakhem, Ricardo Monti, Diederik Kingma, and Aapo Hyvarinen. Ice-beem: Ident-
472 ifiable conditional energy-based deep models based on nonlinear ica. *Advances in Neural*
473 *Information Processing Systems*, 33, 2020.
- 474 [28] Durk P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised
475 learning with deep generative models. In *Advances in neural information processing systems*,
476 pages 3581–3589, 2014.
- 477 [29] Arthur Lewbel. The identification zoo: Meanings of identification in econometrics. *Journal of*
478 *Economic Literature*, 57(4):835–903, 2019.
- 479 [30] Christos Louizos, Uri Shalit, Joris M Mooij, David Sontag, Richard Zemel, and Max Welling.
480 Causal effect inference with deep latent-variable models. In *Advances in Neural Information*
481 *Processing Systems*, pages 6446–6456, 2017.
- 482 [31] Danni Lu, Chenyang Tao, Junya Chen, Fan Li, Feng Guo, and Lawrence Carin. Reconsidering
483 generative objectives for counterfactual reasoning. *Advances in Neural Information Processing*
484 *Systems*, 33, 2020.
- 485 [32] Wei Luo, Yeying Zhu, and Debashis Ghosh. On estimating regression-based causal effects
486 using sufficient dimension reduction. *Biometrika*, 104(1):51–65, 2017.
- 487 [33] Michael Oberst, Fredrik Johansson, Dennis Wei, Tian Gao, Gabriel Brat, David Sontag, and
488 Kush Varshney. Characterization of overlap in observational studies. In *International Conference*
489 *on Artificial Intelligence and Statistics*, pages 788–798. PMLR, 2020.
- 490 [34] Judea Pearl. *Causality: models, reasoning and inference*. Cambridge University Press, 2009.
- 491 [35] Paul R Rosenbaum. Modern algorithms for matching in observational studies. *Annual Review*
492 *of Statistics and Its Application*, 7:143–176, 2020.
- 493 [36] Paul R Rosenbaum and Donald B Rubin. The central role of the propensity score in observational
494 studies for causal effects. *Biometrika*, 70(1):41–55, 1983.
- 495 [37] Donald B Rubin. Estimating causal effects from large data sets using propensity scores. *Annals*
496 *of internal medicine*, 127(8_Part_2):757–763, 1997.
- 497 [38] Donald B Rubin. Causal inference using potential outcomes: Design, modeling, decisions.
498 *Journal of the American Statistical Association*, 100(469):322–331, 2005.

- 499 [39] Alejandro Schuler, David Walsh, Diana Hall, Jon Walsh, and Charles Fisher. Increasing the
500 efficiency of randomized trial estimates via linear adjustment for a prognostic score. *arXiv*
501 *preprint arXiv:2012.09935*, 2020.
- 502 [40] Uri Shalit, Fredrik D Johansson, and David Sontag. Estimating individual treatment effect:
503 generalization bounds and algorithms. In *International Conference on Machine Learning*, pages
504 3076–3085. PMLR, 2017.
- 505 [41] Claudia Shi, David Blei, and Victor Veitch. Adapting neural networks for the estimation of
506 treatment effects. In *Advances in Neural Information Processing Systems*, pages 2507–2517,
507 2019.
- 508 [42] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using
509 deep conditional generative models. In *Advances in neural information processing systems*,
510 pages 3483–3491, 2015.
- 511 [43] Peter Sorrenson, Carsten Rother, and Ullrich Köthe. Disentanglement by nonlinear ica with
512 general incompressible-flow networks (gin). In *International Conference on Learning Repre-*
513 *sentations*, 2019.
- 514 [44] Elizabeth A. Stuart. Matching Methods for Causal Inference: A Review and a Look Forward.
515 *Statistical Science*, 25(1):1 – 21, 2010.
- 516 [45] Elizabeth A Stuart, Brian K Lee, and Finbarr P Leacy. Prognostic score–based balance measures
517 can be a useful diagnostic for propensity score methods in comparative effectiveness research.
518 *Journal of clinical epidemiology*, 66(8):S84–S90, 2013.
- 519 [46] Xinwei Sun, Botong Wu, Chang Liu, Xiangyu Zheng, Wei Chen, Tao Qin, and Tie-yan Liu.
520 Latent causal invariant model. *arXiv preprint arXiv:2011.02203*, 2020.
- 521 [47] Alexander Tarr and Kosuke Imai. Estimating average treatment effects with support vector
522 machines. *arXiv preprint arXiv:2102.11926*, 2021.
- 523 [48] Mark J Van der Laan and Sherri Rose. *Targeted learning: causal inference for observational*
524 *and experimental data*. Springer Science & Business Media, 2011.
- 525 [49] Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects
526 using random forests. *Journal of the American Statistical Association*, 113(523):1228–1242,
527 2018.
- 528 [50] Halbert White and Karim Chalak. Identification and identification failure for treatment effects
529 using structural systems. *Econometric Reviews*, 32(3):273–317, 2013.
- 530 [51] Pengzhou Wu and Kenji Fukumizu. Causal mosaic: Cause-effect inference via nonlinear ica
531 and ensemble method. In *International Conference on Artificial Intelligence and Statistics*,
532 pages 1157–1167. PMLR, 2020.
- 533 [52] S Yang and P Ding. Asymptotic inference of causal effects with observational studies trimmed
534 by the estimated propensity scores. *Biometrika*, 105(2):487–493, 03 2018.
- 535 [53] Liuyi Yao, Sheng Li, Yaliang Li, Mengdi Huai, Jing Gao, and Aidong Zhang. Representa-
536 tion learning for treatment effect estimation from observational data. In *Advances in Neural*
537 *Information Processing Systems*, pages 2633–2643, 2018.
- 538 [54] Jinsung Yoon, James Jordon, and Mihaela van der Schaar. GANITE: Estimation of individ-
539 ualized treatment effects using generative adversarial nets. In *International Conference on*
540 *Learning Representations*, 2018.
- 541 [55] Weijia Zhang, Lin Liu, and Jiuyong Li. Treatment effect estimation with disentangled latent
542 factors. *arXiv preprint arXiv:2001.10652*, 2020.
- 543 [56] Yao Zhang, Alexis Bellot, and Mihaela Schaar. Learning overlapping representations for
544 the estimation of individualized treatment effects. In *International Conference on Artificial*
545 *Intelligence and Statistics*, pages 1005–1014. PMLR, 2020.

546 **Checklist**

- 547 1. For all authors...
- 548 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
- 549 contributions and scope? [Yes]
- 550 (b) Did you describe the limitations of your work? [Yes] Particularly in Conclusion.
- 551 (c) Did you discuss any potential negative societal impacts of your work? [Yes] We men-
- 552 tioned possible misuse of our method and suggestions for practitioners in Introduction
- 553 and Conclusion.
- 554 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
- 555 them? [Yes]
- 556 2. If you are including theoretical results...
- 557 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 558 (b) Did you include complete proofs of all theoretical results? [Yes] In Appendix.
- 559 3. If you ran experiments...
- 560 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 561 mental results (either in the supplemental material or as a URL)? [Yes] In supplemental
- 562 material.
- 563 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 564 were chosen)? [Yes] Possibly in Appendix.
- 565 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 566 ments multiple times)? [Yes]
- 567 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 568 of GPUs, internal cluster, or cloud provider)? [Yes]
- 569 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 570 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 571 (b) Did you mention the license of the assets? [Yes] In code project.
- 572 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 573 In supplemental material.
- 574 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 575 using/curating? [Yes] In Appendix. And we refer to the original paper where this is
- 576 discussed in details.
- 577 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 578 information or offensive content? [Yes] Refer to the original paper. And we believe
- 579 there are no privacy issues.
- 580 5. If you used crowdsourcing or conducted research with human subjects...
- 581 (a) Did you include the full text of instructions given to participants and screenshots, if
- 582 applicable? [N/A]
- 583 (b) Did you describe any potential participant risks, with links to Institutional Review
- 584 Board (IRB) approvals, if applicable? [N/A]
- 585 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 586 spent on participant compensation? [N/A]