

CAUSAL EFFECT ESTIMATION WITH MIXED LATENT CONFOUNDERS AND POST-TREATMENT VARIABLES

Anonymous authors

Paper under double-blind review

ABSTRACT

Causal inference from observational data has attracted considerable attention in recent years. One main obstacle is the handling of confounders. As the direct measure of confounders may not always be feasible, recent methods seek to address the confounding bias with proxy variables, which are covariates researchers postulate to be conducive to the inference of latent confounders. However, observed covariates may scramble both latent confounders and latent post-treatment variables in observational study, where existing methods risk biasing the estimation by unintentionally controlling for variables affected by the treatment. In this paper, we systematically investigate the bias due to latent post-treatment variables, i.e., *latent post-treatment bias*, in causal effect estimation. We first derive the bias of existing methods when selected proxies scramble both latent confounders and latent post-treatment variables, which we demonstrate can be arbitrarily bad. We then propose a novel Confounder-identifiable VAE (CiVAE) to address the bias. CiVAE is built upon the assumption that the prior of the latent variables belongs to a general exponential family with at least one invertible sufficient statistic in the factorized part. Based on this, we show that latent confounders and latent post-treatment variables can be individually identified up to simple bijective transformations. Finally, we prove that the true causal effects can be unbiasedly estimated with the transformed confounders inferred by CiVAE. Experiments on both simulated and real-world datasets demonstrate that CiVAE is significantly more robust to latent post-treatment bias than existing methods for causal effects estimation.

1 INTRODUCTION

Causal inference, which seeks to draw conclusions about cause-and-effect relationships among variables of interest, has gained increasing prominence in various fields, such as social science, economics, and public health (Glass et al., 2013; Johansson et al., 2016; Prospero et al., 2020). Traditional methods rely on randomized control trials (RCT) to draw valid causal conclusions from experimentation (Cook et al., 2002). Recently, more attention has been dedicated toward causal inference from observational datasets, which contain samples with passively observed past treatment, the associated outcome, and possibly features, and in which researchers have no control over the treatment assignment mechanism (Shalit et al., 2017; Shi et al., 2019; Wager & Athey, 2018).

One main obstacle to inferring causal relations from observational data is confounding bias, which occurs when past treatments were determined by variables that causally influence the outcome, i.e., confounders. In such cases, the difference in the average outcome between the treatment group and the non-treatment group cannot be attributed solely to the treatment, but may also be due to the systematic difference of samples in the two groups (Mickey & Greenland, 1989). If the confounders can be observed, a simple strategy to address such a bias is to control them via methods such as covariate adjustment (Pocock et al., 2002) or propensity score re-weighting (Li et al., 2018). However, confounders are not always measurable (Kuroki & Pearl, 2014). Therefore, recent methods seek to adjust for the influence of confounders based on their noisy proxies, which are generally covariates researchers postulated to be conducive for the inference of confounders (Miao et al., 2018; Yao et al., 2018; Madras et al., 2019). One exemplar work from this strain is the causal effect variational auto-encoder (CEVAE) (Louizos et al., 2017) (Fig. 1-(a)), which has demonstrated that confounding bias can be mitigated by controlling latent variables inferred from proxies of confounders.

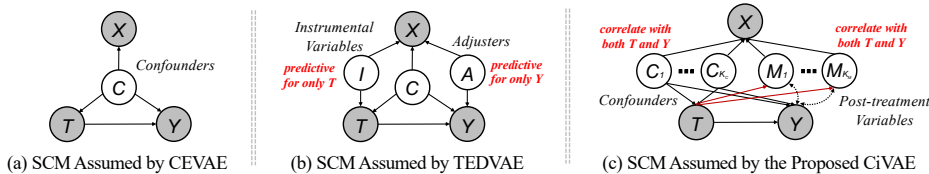


Figure 1: Structural causal model (SCM) assumed by CEVAE, TEDVAE, and CiVAE.

Although proxy-of-confounder-based methods have achieved substantial progress, we argue that these algorithms may risk controlling latent post-treatment variables (i.e., variables causally affected by the treatment) scrambled in the proxy variables, where **post-treatment bias** may be unintentionally introduced in the estimated treatment effect. Here, we note that the negative effects of controlling post-treatment variables have been investigated in prior research (Acharya et al., 2016; Elwert & Winship, 2014; King & Zeng, 2006). For example, Montgomery et al. (2018) found that more than 50% of the papers published in top journals of politics inadvertently control post-treatment variables in the experimental setting, although researchers have complete control over the treatment assignment mechanism and the covariates to control for. On this basis, we postulate that post-treatment bias could be even worse for proxy-based methods in the setting of observational study, as when treatments are passively recorded, it is difficult to determine which variables causally influence the treatment and which variables are influenced by it (as both confounders and post-treatment variables are correlated with the treatment and the outcome). In addition, the post-treatment variables can be latent, which may be scrambled into the observed covariates together with the latent confounders.

Consider the following real-world example that researchers from the Company¹ have encountered when estimating the average causal effects of *switching a job from onsite to online mode to the statistics of the applicants* (e.g., average age, gender/geographical diversity, etc.). In this case, the Company collected a dataset of two groups of online (i.e., the treatment group) and onsite jobs (i.e., the non-treatment group), where for each job, the statistics of the applicants (i.e., the average age) are calculated as the outcome. Clearly, the seniority of the job is a confounder between the treatment and the outcome, as less senior jobs (e.g., internships) are more likely to permit online work, and applicants for these jobs tend to be younger on average. The seniority of a job can be difficult to measure. Therefore, the required skills of the job, which the recruiter must provide when publishing a job Ad in the Company, can be used as the proxy of the confounder "seniority". However, a **caveat** is that, switching to an online working mode may also alter the required skills of a job, thereby affecting the qualification of the applicants (where these altered skills are mediators). Consequently, directly using the required skills as the proxy of the confounder "seniority" could unintentionally control latent mediators, which introduces post-treatment bias in the causal effect estimation results.

Addressing the **latent post-treatment bias** faces multi-faceted challenges. First, there lacks a theoretical formulation of the bias when the selected proxies scramble latent post-treatment variables for proxy-of-confounder-based methods; the trade-off between deconfounding and introducing new post-treatment bias is not clear. In addition, it is difficult to distinguish confounders and post-treatment variables in the latent space. Existing covariate disentanglement-based methods, e.g., TEDVAE (Zhang et al., 2021), mainly focus on an easier task of disentangling latent confounders with latent adjusters and instrumental variables. This can be achieved by using their different predictive abilities w.r.t. the treatment and outcome (see Fig. 1-(b)). However, since latent confounders and post-treatment variables correlate with both the treatment and outcome, the two cannot be disentangled by these methods. One solution is to assume the proxy of latent post-treatment variables can be observed, from which post-treatment variables can be inferred and disentangled from the latent confounders. However, this assumption is **too strong**, as in the previous online/onsite job case, we can never know which skills are causally influenced by the work mode. Finally, even if latent confounders can be distinguished, since general latent variable models have no identifiability guarantee (Khemakhem et al., 2020), it is unclear whether controlling the inferred latent variables, which may be arbitrary transformations of the true confounders, can provide unbiased estimations of the causal effects.

To address the aforementioned challenges, we provide a systematic investigation of the latent post-treatment bias in causal inference. We first analyze the behavior of existing proxy-based causal inference methods when the selected proxies scramble both latent confounders and post-treatment variables, where we demonstrate that the estimated average causal effects can be arbitrarily biased.

¹Anonymized due to double-blind review policy.

We then propose the Confounder-identifiable VAE (CiVAE) to address such biases. Specifically, we show that based on a mild assumption that the prior distribution of latent variables (i.e., the latent confounders and post-treatment variables) belongs to a general exponential family with at least one invertible sufficient statistic in the factorized part, latent confounders and latent post-treatment variables can be *individually* identified up to *simple bijective transformations*. In addition, based on the causal relations among confounders, mediators, and treatment, we further demonstrate that the inferred confounders (which are actually transformed proxies of the true confounders) could be properly distinguished from the inferred latent post-treatment variables with pair-wise conditional independence tests. Finally, we prove that the true causal effects can be unbiasedly estimated based on transformed confounders inferred by CiVAE. Experiments on both simulated and real-world datasets demonstrate that CiVAE shows more robustness to latent post-treatment bias than existing methods.

2 PROBLEM FORMULATION AND ANALYSIS

2.1 PROBLEM FORMULATION

Throughout this paper, we assume the causal model in Fig. 1-(c), where the dashed lines denote indeterminate causal mechanisms that might vary in different cases. We use a binary random variable T to denote the treatment, a random vector $\mathbf{X} \in \mathbb{R}^{K_X}$ to denote the observed covariates, and a random scalar $Y \in \mathbb{R}$ to denote the outcome. Furthermore, observed covariates \mathbf{X} are assumed to be generated from K_C independent latent confounders $\mathbf{C} \triangleq [C_1, C_2, \dots, C_{K_C}]$ and K_M latent post-treatment variables $\mathbf{M} \triangleq [M_1, M_2, \dots, M_{K_M}]$ under the causal influence of treatment T . We use the random vector $\mathbf{Z} \triangleq [\mathbf{C} || \mathbf{M}] \in \mathbb{R}^{K_Z=K_C+K_M}$ to denote all latent factors. **Our aim** is to estimate the average causal effects of treatment T on outcome Y with auxiliary confounder information in \mathbf{X} , where the estimation should be devoid of both confounding bias and post-treatment bias.

2.2 ANALYSIS OF LATENT POST-TREATMENT BIAS

2.2.1 PRELIMINARIES AND ASSUMPTIONS

To achieve such a purpose, we first formally define the (conditional) average treatment effects (C/ATE) when covariates \mathbf{X} scramble both latent confounders \mathbf{C} and post-treatment variables \mathbf{M} . We then define the post-treatment bias when covariates \mathbf{X} are used directly as the proxy of confounders. To facilitate the analysis, we make the following assumption regarding the causal generative process.

Assumption 1. (Noisy-Injectivity). We assume $\mathbf{X} = f(\mathbf{C}, \mathbf{M}) + \epsilon$, where f is a deterministic function that combines latent confounders \mathbf{C} and latent post-treatment variables \mathbf{M} into observations \mathbf{X} and ϵ is random noise. In addition, we assume that the function f is **injective**; beyond injectivity, f can be arbitrarily nonlinear. We use $f^\dagger : \mathbf{X} \rightarrow [\mathbf{C} || \mathbf{M}]$ to denote its left inverse. We use $f_C^\dagger : \mathbf{X} \rightarrow \mathbf{C}$ and $f_M^\dagger : \mathbf{X} \rightarrow \mathbf{M}$ to denote the mapping from \mathbf{X} to \mathbf{C} , \mathbf{M} , respectively.

Noisy-Injectivity is a common assumption made either explicitly or implicitly in most existing proxy-of-confounder-based causal inference algorithms. For example, if both \mathbf{X} and \mathbf{C} are categorical, Pearl (2012) assumes that \mathbf{X} has at least the same number of categories as \mathbf{C} , whereas the effect restoration algorithm (Rothman et al., 2008) assumes that the matrix of $p(\mathbf{C}, \mathbf{X})$ to be full-rank. Although CEVAE (Louizos et al., 2017) makes no explicit injectivity assumption between \mathbf{C} and \mathbf{X} , it requires that the joint distribution $p(\mathbf{C}, \mathbf{X}, T, Y)$ can be fully recovered from the observations (\mathbf{X}, T, Y) . The literature shows that some of the possible identification criteria are **1)** multiple independent views of \mathbf{C} in \mathbf{X} (Edwards et al., 2015), and **2)** \mathbf{C} is categorical and \mathbf{X} is a mixture of Gaussian components determined by \mathbf{C} (that is, \mathbf{X} is generated by bijective mapping of \mathbf{C} to the mean of the corresponding component with added Gaussian noise) (Anandkumar et al., 2014).

In the following part of this section, we omit the noise ϵ to gain better intuition of latent post-treatment bias (but all the exact conclusions will still hold in the posterior sense). In Section 3, we assume that noise exists and demonstrate that our method can still adequately identify latent confounders.

2.2.2 CAUSAL ESTIMAND AND THE TRUE ATE

Based on Assumption 1, we are ready to define the estimated average treatment effect (ATE) by controlling the covariates \mathbf{X} , as well as the true (conditional) average treatment effects.

Definition 1. We define the *Difference in Conditional Expected Values (DCEV)* as

$$DCEV(\mathbf{x}) = \mathbb{E}[Y|T = 1, \mathbf{X} = \mathbf{x}] - \mathbb{E}[Y|T = 0, \mathbf{X} = \mathbf{x}], \quad (1)$$

which is the difference of the expected value of the outcome Y for units with variable $\mathbf{X} = \mathbf{x}$ in the treatment group and the non-treatment group. Based on $DCEV(\mathbf{x})$, we define the *Difference in Expected Value (DEV)*, i.e., $DEV(\mathbf{X}) = \mathbb{E}_{p(\mathbf{X})}[DCEV(\mathbf{X})]$ as the expected value DCEV.

$DEV(\mathbf{X})$ denotes the ATE estimand by controlling covariates \mathbf{X} . If $\mathbf{X} = \emptyset$, $DEV(\emptyset)$ represents the *naive estimator* that directly calculates the expected difference of Y between the treatment group and the non-treatment group. With the causal estimand $DEV(\mathbf{X})$ introduced, we then define the true causal effects (i.e., C/ATE) when covariates \mathbf{X} scramble both latent confounders and post-treatment variables according to the generative process described in Assumption 1. The main issue that hinders a direct definition of C/ATE with $DCEV(\mathbf{x})$ and $DEV(\mathbf{X})$ is that, since \mathbf{X} contains latent post-treatment variables \mathbf{M} , conditional on \mathbf{X} , the strong ignorability assumption (Imbens & Rubin, 2015) widely used for the identification of causal effects **does not hold**². Accordingly, we have:

Definition 2. Under Assumption 1, we define the *Conditional Average Treatment Effect (CATE)* for individuals with observed covariates $\mathbf{X} = \mathbf{x}$ as follows:

$$CATE(\mathbf{x}) = \mathbb{E}[Y|T = 1, \mathbf{C} = f_C^\dagger(\mathbf{x})] - \mathbb{E}[Y|T = 0, \mathbf{C} = f_C^\dagger(\mathbf{x})], \quad (2)$$

with the *Average Treatment Effect (ATE)* of treatment T defined as

$$ATE = \mathbb{E}[Y|do(T = 1)] - \mathbb{E}[Y|do(T = 0)] = \mathbb{E}_{p(\mathbf{C})}[\mathbb{E}[Y|T = 1, \mathbf{C}] - \mathbb{E}[Y|T = 0, \mathbf{C}]]. \quad (3)$$

In Definition 2, we only consider the latent confounder component of \mathbf{X} for CATE in Eq. (2), as the causal relationship between the post-treatment variables \mathbf{M} and the outcome Y is indeterminate (see Fig. 1-(c)). However, if the specific relationship between \mathbf{M} and Y can be further established by the researcher (e.g., all elements of \mathbf{M} are latent mediators), more precise forms of CATE can be derived with path-specific counterfactual analysis (Imai et al., 2010; Cheng et al., 2022).

2.2.3 LATENT POST-TREATMENT BIAS

With $DEV(\mathbf{X})$ (the ATE estimator that controls the covariates \mathbf{X}), CATE, and ATE defined in Section 2.2.2, in this section, we analyze the *latent post-treatment bias* of existing proxy-of-confounder-based causal inference methods, such as CEVAE (Louizos et al., 2017), that control latent variables inferred from the covariates \mathbf{X} to estimate the ATE of T on Y , when \mathbf{X} scrambles both latent confounders and post-treatment variables. In our analysis, Lemma 2.1 will be frequently used.

Lemma 2.1. For an injective function g , $\mathbb{E}[Y|\mathbf{X} = \mathbf{x}] = \mathbb{E}[Y|g(\mathbf{X}) = g(\mathbf{x})]$ holds.

The proof when g is differentiable *a.e.* can be referred to in Appendix A.1. Since the latent variable models used in existing methods (such as VAE with factorized Gaussian prior in CEVAE) lack identifiability guarantee (i.e., the recovery of the exact latent variables), we assume that these models can recover the true latent space $\mathbf{Z} = [\mathbf{C}, \mathbf{M}]$ up to invertible transformations \tilde{f} , where the inference process can be represented as $\tilde{\mathbf{Z}} = \tilde{f}(\mathbf{X}) = \tilde{f} \circ f^\dagger(\mathbf{X})$. With such an assumption, we have the following theorem regarding the latent post-treatment bias when \mathbf{X} mixes post-treatment variables.

Theorem 2.2. If the observed covariates \mathbf{X} are generated from latent confounders \mathbf{C} and latent post-treatment variables \mathbf{M} according to Assumption 1, the latent post-treatment bias of a proxy-of-confounder-based causal inference algorithm that controls latent variables $\tilde{\mathbf{Z}}$ inferred from \mathbf{X} via $\tilde{f} = \tilde{f} \circ f^\dagger : \mathbb{R}^{K_X} \rightarrow \mathbb{R}^{K_C + K_M}$ to estimate the ATE can be formulated as follows:

$$\begin{aligned} Bias(\mathbf{X}) &= ATE - DEV(\tilde{f}(\mathbf{X})) = ATE - \mathbb{E}[\mathbb{E}[Y|T = 1, \tilde{f}(\mathbf{X})] - \mathbb{E}[Y|T = 0, \tilde{f}(\mathbf{X})]] \\ &= ATE - \mathbb{E}[\mathbb{E}[Y|T = 1, \tilde{f} \circ f^\dagger(f(\mathbf{C}, \mathbf{M}))] - \mathbb{E}[Y|T = 0, \tilde{f} \circ f^\dagger(f(\mathbf{C}, \mathbf{M}))]] \\ &= \mathbb{E}[\mathbb{E}[Y|T = 1, \mathbf{C}] - \mathbb{E}[Y|T = 0, \mathbf{C}]] - \mathbb{E}[\mathbb{E}[Y|T = 1, \mathbf{C}, \mathbf{M}] - \mathbb{E}[Y|T = 0, \mathbf{C}, \mathbf{M}]], \end{aligned} \quad (4)$$

which can be arbitrarily bad. Therefore, the estimator of existing proxy-of-confounder-based methods, i.e., $DEV(\tilde{f}(\mathbf{X}))$, is an arbitrarily biased estimator of the ATE, when the selected proxy of confounders \mathbf{X} accidentally mixes in latent post-treatment variables \mathbf{M} .

²Equivalently, we could say that given covariates \mathbf{X} , the **backdoor criteria** between T and Y does not hold, which requires the conditional set of variables contains no descendants of the treatment T (Glymour et al., 2016).

The final step of Eq. (4) can be proved since f is injective and \bar{f} bijective, the composite $\bar{f} \circ f^\dagger \circ f : [\mathbf{C}, \mathbf{M}] \rightarrow \hat{\mathbf{Z}}$ is bijective, so we can use Lemma 2.1 to remove $\bar{f} \circ f^\dagger \circ f$ in the condition.

2.2.4 EXAMPLES IN THE LINEAR CASES

Generally, the latent post-treatment bias defined in Eq. (4) cannot be simplified because **1**) the causal relationship between \mathbf{M} and Y is indeterminate, and **2**) the causal influence of \mathbf{C} , \mathbf{M} , and T on Y can be arbitrary. However, for linear structural causal models with causal relationships determined between \mathbf{M} and Y (e.g., \mathbf{M} are mediators, which are post-treatment variables that have causal influences on the outcomes), stronger conclusions can be drawn as follows:

Corollary 2.3. (*MixedMediator*) For the linear Structural Causal Model (SCM) defined as:

$$\begin{aligned} T &\leftarrow \mathbf{1}(\alpha_T + \sum \beta_i \cdot C_i > a) \\ M_j &\leftarrow \alpha_M + \gamma_j \cdot T \\ \mathbf{X} &\leftarrow \boldsymbol{\alpha}_X + \mathbf{A}[\mathbf{M}|\mathbf{C}] \\ Y &\leftarrow \alpha_Y + \tau \cdot T + \sum \theta_j \cdot M_j + \sum \kappa_i \cdot C_i, \end{aligned} \quad (5)$$

where the mixture function $f = \mathbf{A} \in \mathbb{R}^{K_X \times (K_C + K_M)}$ is a full column-rank matrix, the CATE, ATE, and the bias of proxy-of-confounder-based causal inference model that controls the latent variables $\hat{\mathbf{Z}}$ inferred via $\hat{\mathbf{Z}} = \tilde{f}(\mathbf{X}) = \mathbf{B}^T \mathbf{X}$ can be formulated as follows:

$$\begin{aligned} ATE &= CATE = \tau + \sum \gamma_j \cdot \theta_j \\ DEV(\hat{\mathbf{Z}}) &= \mathbb{E}[DCEV(\hat{\mathbf{Z}})] = DCEV(\hat{\mathbf{Z}}) = \tau \\ Bias(\hat{\mathbf{Z}}) &= ATE - DEV(\hat{\mathbf{Z}}) = \sum \gamma_j \cdot \theta_j, \end{aligned} \quad (6)$$

where $\mathbf{B} \in \mathbb{R}^{K_X \times (K_C + K_M)}$ is another full column-rank matrix. Since $\sum \gamma_j \cdot \theta_j$ is arbitrary, the estimator $DEV(\hat{\mathbf{Z}}) = \mathbb{E}[DCEV(\mathbf{B}^T \mathbf{X})]$ is arbitrarily biased for ATE estimation.

The proof of Eq. (6) is provided in Appendix A.2. In addition, we show that the post-treatment variables \mathbf{M} DO NOT necessarily need to have direct causal effects on the outcome Y to incur arbitrary bias in ATE estimation. In Appendix A.3, we provide another example (i.e., *MixedCorrelator*) in the linear case where \mathbf{M} is correlated with Y through unobserved confounders \mathbf{U} in Corollary A.1.

3 METHODOLOGY

In this section, we introduce the proposed Confounder-identifiable Variational Auto-Encoder (CiVAE) to address latent post-treatment bias. Specifically, we first prove that if the prior distribution of the true latent variables $\mathbf{Z} = [\mathbf{C}, \mathbf{M}]$ satisfies certain weak assumptions, identifiability criterion holds, and each dimension of the inferred latent variables $\hat{\mathbf{Z}}$, i.e., \hat{Z}_i , corresponds to the invertible transformation of **either** a true confounder C_j **or** a true post-treatment variable M_k . Then, utilizing the causal relations between \mathbf{C} , \mathbf{M} , and T , we novelly transform the challenging confounder-identifiability problem into a tractable pair-wise conditional independence test problem, which can be effectively solved with kernel-based methods. Finally, we demonstrate that controlling the transformed confounders inferred by CiVAE can yield an unbiased estimation of the true ATE.

3.1 GENERATIVE PROCESS

The fundamental work of deep variational inference with identifiability guarantee, i.e., the identifiable VAE (iVAE) (Khemakhem et al., 2020), makes a strict assumption that the prior of true latent variables \mathbf{Z} (i.e., $[\mathbf{C}, \mathbf{M}]$ in our case) is conditionally factorized given the available covariates (i.e., the treatment T and the outcome Y in our case). However, since both latent confounders \mathbf{C} and latent post-treatment variables \mathbf{M} form fork structures with the outcome Y (see Fig. 1-(c)) (Koller & Friedman, 2009), $C_i, C_j, M_i,$ and M_j are not independent given Y . Recently, Non-Factorized iVAE (NF-iVAE) (Lu et al., 2021) was proposed that allows arbitrary dependence among the true latent variables \mathbf{Z} in the conditional priors, where \mathbf{Z} can be identified up to arbitrary non-linear

transformations, However, the transformation are not necessarily invertible, which is risky for causal inference, as multiple values of the confounders may collapse, leading to bias when estimating the ATE by averaging the *DCEV* calculated in each stratum of the inferred confounders.

The proposed NF-iVAE guarantees the identifiability of \mathbf{Z} by putting a general exponential family distribution with at least one invertible sufficient statistic in the factorized part as its prior when conditioning on treatment T and outcome Y , which can be formulated as follows.

Assumption 2. Let $\mathbf{Z} = [\mathbf{C}||\mathbf{M}]$ be the random vector for latent variables that causally generate the observed covariates \mathbf{X} according to Assumption 1. We assume that the conditional prior of \mathbf{Z} given the outcome Y and the treatment T belongs to a general exponential family with parameter vector $\boldsymbol{\lambda}(Y, T)$ and sufficient statistics $\mathbf{S}(\mathbf{Z}) = [\mathbf{S}_f(\mathbf{Z})^T, \mathbf{S}_{n_f}(\mathbf{Z})^T]^T$. Specifically, $\mathbf{S}(\mathbf{Z})$ is composed of (i) the sufficient statistics of a factorized exponential family, i.e., $\mathbf{S}_f(\mathbf{Z}) = [\mathbf{S}_1(Z_1)^T, \dots, \mathbf{S}_{K_Z}(Z_{K_Z})^T]^T$, where all components $\mathbf{S}_i(Z_i)$ have dimension larger than or equal to 2 and **each \mathbf{S}_i has at least one invertible dimension**, and (ii) $\mathbf{S}_{n_f}(\mathbf{Z})$, where \mathbf{S}_{n_f} is a neural network with ReLU activation. The density of the conditional prior can be formulated as:

$$p_{\mathbf{S}, \boldsymbol{\lambda}}(\mathbf{Z} | Y, T) = \mathcal{Q}(\mathbf{Z}) / \mathcal{C}(Y, T) \exp[\mathbf{S}(\mathbf{Z})^T \boldsymbol{\lambda}(Y, T)], \quad (7)$$

where $\mathcal{Q}(\mathbf{Z})$ is the base measure and $\mathcal{C}(Y, T)$ not dependent on \mathbf{Z} is the normalizing constant.

We justify that assumption 2 is weak and practical as follows. **1)** Neural networks with ReLU activation have universal approximation ability of distributions (Lu & Lu, 2020). Therefore, Eq. (7) can model arbitrary dependence between true latent confounders \mathbf{C} and true post-treatment variables \mathbf{M} conditional on T and Y . **2)** Although CiVAE makes an extra assumption that $\forall i$, at least one dimension of \mathbf{S}_i is invertible, this can be easily satisfied as most commonly used exponential family distributions, such as Gaussian, Bernoulli, etc., has at least one invertible sufficient statistics³.

The reason why we use ReLU as the activation is that, the identifiability of iVAE relies on the condition that the sufficient statistics \mathbf{S} have zero second-order cross-derivative. The factorized part, i.e., \mathbf{S}_f , satisfies it trivially since all cross-derivatives of \mathbf{S}_f are zero. In addition, since the ReLU neural networks are linear *a.e.*, all second-order derivatives of \mathbf{S}_{n_f} are zero. Therefore, identifiability holds after adding \mathbf{S}_{n_f} in the prior that allows the capturing of arbitrary dependence among \mathbf{Z} .

3.2 OPTIMIZATION OBJECTIVE

Combining Assumptions 1 and 2, the generative process of CiVAE can be formulated as follows:

$$p_{\boldsymbol{\theta}}(\mathbf{X}, \mathbf{Z} | Y, T) = p_f(\mathbf{X} | \mathbf{Z}) p_{\mathbf{S}, \boldsymbol{\lambda}}(\mathbf{Z} | Y, T), \quad (8)$$

$$p_f(\mathbf{X} | \mathbf{Z}) = p_{\epsilon}(\mathbf{X} - f(\mathbf{Z})). \quad (9)$$

where $\boldsymbol{\theta} = (f, \boldsymbol{\lambda}, \mathbf{S}) \in \Theta$ are the parameters of the generative distribution⁴. Since the generative process of CiVAE is parameterized by deep neural networks, the posterior distribution of \mathbf{Z} , i.e., $p_{\boldsymbol{\theta}}(\mathbf{Z} | \mathbf{X}, Y, T)$, is intractable. Therefore, we resort to variational inference (Blei et al., 2017), where we introduce approximate posterior $q_{\phi}(\mathbf{Z} | \mathbf{X}, Y, T)$ parameterized by deep neural network with trainable parameter ϕ , and in $q_{\phi}(\mathbf{Z} | \cdot)$ finds the one closes to $p_{\boldsymbol{\theta}}(\mathbf{Z} | \cdot)$ measured by KL divergence. Minimization of the KL is equivalent to maximization of the evidence lower bound (ELBO) as:

$$\mathcal{L}(\boldsymbol{\theta}, \phi) := \mathbb{E}_{q_{\phi}(\mathbf{Z} | \mathbf{X}, Y, T)} \left[\log p_f(\mathbf{X} | \mathbf{Z}) + \underbrace{\log p_{\mathbf{S}, \boldsymbol{\lambda}}(\mathbf{Z} | Y, T) - \log q_{\phi}(\mathbf{Z} | \mathbf{X}, Y, T)}_{\text{KL of posterior with prior}} \right]. \quad (10)$$

Since the normalization constant \mathcal{C} in Eq. (7) is generally intractable, it is infeasible to directly learn $\mathbf{S}, \boldsymbol{\lambda}$ by optimizing Eq. (10). Therefore, we substitute the KL term in Eq. (10) with the widely-used score matching (Hyvärinen & Dayan, 2005) to learn unnormalized densities instead as follows:

$$\begin{aligned} \mathcal{L}(\mathbf{S}, \boldsymbol{\lambda}, \phi) &:= \mathbb{E}_{q_{\phi}(\mathbf{Z} | \mathbf{X}, Y, T)} \left[\|\nabla_{\mathbf{Z}} \log q_{\phi}(\mathbf{Z} | \mathbf{X}, Y, T) - \nabla_{\mathbf{Z}} \log p_{\mathbf{S}, \boldsymbol{\lambda}}(\mathbf{Z} | Y, T)\|^2 \right] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{Z} | \mathbf{X}, Y, T)} \left[\sum_{j=1}^{K_Z} \left[\frac{\partial^2 p_{\mathbf{S}, \boldsymbol{\lambda}}(\mathbf{Z} | Y, T)}{\partial Z_j^2} + \frac{1}{2} \left(\frac{\partial p_{\mathbf{S}, \boldsymbol{\lambda}}(\mathbf{Z} | Y, T)}{\partial Z_j} \right)^2 \right] \right] + \text{const.} \end{aligned} \quad (11)$$

³There are a few exponential family with no invertible sufficient statistics, e.g., Weibull distribution when shape parameter k is even.

⁴Note that although f is a function, we include it in the parameter set to be consistent with the iVAE paper.

3.3 IDENTIFIABILITY OF CiVAE

With the generative process and optimization objective of CiVAE introduced in the previous subsections, we are ready to introduce the final assumption of CiVAE, which, combined with Assumptions 1 and 2, leads to the main theorem of this paper, which states the identifiability of CiVAE.

Assumption 3. Assume the following: (i) The set $\{\mathbf{X} \in \mathcal{X} | \phi(\mathbf{X}) = 0\}$ has measure zero, where ϕ is the characteristic function of the density p_f in Eq. (9). (ii) The sufficient statistics, \mathbf{S}_i in \mathbf{S}_f are all twice differentiable. (iii) The mixture function f in Eq. (9) has all second-order cross derivatives. (iv) There exist $k + 1$ distinct points $(Y, T)_0, \dots, (Y, T)_k$ such that the matrix $\mathbf{L} = [\boldsymbol{\lambda}((Y, T)_1) - \boldsymbol{\lambda}((Y, T)_0), \dots, \boldsymbol{\lambda}((Y, T)_k) - \boldsymbol{\lambda}((Y, T)_0)]$ of size $k \times k$ is invertible, where $k = \text{Dim}(\mathbf{S})$.

(i) - (iii) are trivial for neural networks. (iv) denotes that independent samples of (Y, T) are required to identify \mathbf{C} and \mathbf{M} . The identifiability theorem of CiVAE can be formulated as follows.

Theorem 3.1. If Assumptions 1, 2, and 3 hold, and if $\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}} \in \Theta \rightarrow p_{\boldsymbol{\theta}}(\mathbf{X} | Y, T) = p_{\tilde{\boldsymbol{\theta}}}(\mathbf{X} | Y, T)$, the true latent variables \mathbf{Z} are identifiable up to **permutation and element-wise bijective transformation**. Furthermore, in the case of **variational inference**, if we denote the true parameter that generates the data as $\boldsymbol{\theta}^*$, if (i) the distribution family $q_{\phi}(\mathbf{Z} | \mathbf{X}, Y, T)$ contains the posterior $p_{\boldsymbol{\theta}}(\mathbf{Z} | \mathbf{X}, Y, T)$, and $q_{\phi}(\mathbf{Z} | \mathbf{X}, Y, T) > 0$, (ii) we optimize Eq. (4) w.r.t. both $\boldsymbol{\theta}, \phi$, then in the limit of infinite data, true parameters $\boldsymbol{\theta}^*$ can be learned up to a permutation and bijective transformation of \mathbf{Z} .

The proof of Theorem 3.1 is based on the NF-iVAE paper (Lu et al., 2021), with the new assumption introduced in CiVAE that each \mathbf{S}_i has at least one invertible dimension incorporated to ensure that the transformation of each Z_i is bijective. The detailed proof is provided in Appendix A.4.

3.4 IDENTIFICATION OF LATENT CONFOUNDERS

Theorem 3.1 ensures that latent variables $\hat{\mathbf{Z}}$ inferred by CiVAE cannot **1)** mix confounders and post-treatment variables in each dimension, or **2)** collapse different values of the latent confounders into the same value. To further determine the dimensions of confounder and post-treatment variable in $\hat{\mathbf{Z}}$, we rely on the causal relations between latent variables $\mathbf{Z} = [\mathbf{C}, \mathbf{M}]$ and treatment T and the associated marginal/conditional dependence properties. These are discussed as follows.

- **Case 1. Intra-Confounders.** Latent confounders C_i, C_j and the treatment T form the *V-structure* $C_i \rightarrow T \leftarrow C_j$. Therefore, C_i and C_j are marginally **independent**, whereas they become **dependent** when conditioning on the assigned treatment T .
- **Case 2. Intra-Post Treatment Variables.** Latent post-treatment variables M_i, M_j and the treatment T form a *fork-structure* $M_i \leftarrow T \rightarrow M_j$, where M_i, M_j are marginally **dependent**, but they become **independent** after conditioning on the assigned treatment T .
- **Case 3. Cross-Confounder and Post-Treatment Variables.** Latent confounder C_i , latent post-treatment variable M_j , and the treatment T forms a *chain structure* $C_i \rightarrow T \rightarrow M_j$, where C_i, M_j are marginally dependent, and they become **independent** after conditioning on T .

From the above analysis we can find that, the dependence between two latent variables Z_i and Z_j **increases** after conditioning on the treatment T ONLY in the case of *intra-confounders*. Therefore, if more than one latent confounders exist, which is highly probable when covariates \mathbf{X} are high-dimensional, we can conduct independence test $\text{Ind}(\hat{Z}_i, \hat{Z}_j)$ and $\text{CInd}(\hat{Z}_i, \hat{Z}_j | T)$ for all pairs of inferred latent variables, which can be implemented via kernel-based methods as (Zhang et al., 2012), and select the pairs where p-value of CInd is larger than that of Ind as latent confounders.

Here, we note that the kernel-based (conditional) independence test incurs $N^2 \times K_{\hat{\mathbf{Z}}}^2$ complexity in the training phase. However, once the dimensions of the confounders in $\hat{\mathbf{Z}}$ are determined, CiVAE **has the same complexity as CEVAE** for the estimation of CATE and ATE in the test phase. Therefore, we argue that the additional complexity of model training is worthy due to the substantially increased robustness toward latent post-treatment bias (which will be demonstrated in Section 4).

3.5 ATE ESTIMATOR WITH TRANSFORMED CONFOUNDERS

Finally, we show that controlling transformed confounders $\hat{\mathbf{C}}$ inferred by CiVAE provides an unbiased estimation of ATE. Although assumptions weaker than Assumption 2, e.g., inferred confounders have

the same propensity score as the true confounders (i.e., \hat{C} does not have to be bijective transformation of C), could lead to the same unbiasedness results (Imbens & Rubin, 2015), since our main purpose is to analyze the latent post-treatment bias and propose a viable solution accordingly, this introduces unnecessary complexity, which could be explored as a direction for future study.

Theorem 3.2. *Controlling bijective of confounders is equivalent to controlling true confounders in ATE estimation, i.e., $DEV(\hat{C}) = DEV(g(C)) = ATE$, if transformation function g is bijective.*

The proof of Theorem 3.2 for discrete C is trivial (where $\hat{C} = g(C)$ represents a simple relabeling of the stratum that we calculate the $DCEV$ and take the expectation). The proof in the continuous case where g is differentiable is provided in Appendix A.5. With Theorem 3.2, we can control the identified latent confounders as true confounders, providing an unbiased estimate of ATE.

4 EMPIRICAL STUDY

4.1 DATASETS

We establish two simulated datasets, i.e., `MixedMediator` and `MixedCorrelator`, that consider two types of post-treatment variables, i.e., **1**) mediators and **2**) variables that are correlated with the outcome Y via latent confounders U . The generative process of the two datasets can be referred to in Corollary 2.3 and Corollary A.1, respectively, where the latent confounders C are generated from Gaussian as $C \sim \text{Gaussian}(0, \mathbf{I}_{K_C})$. For `MixedMediator`, γ is set as $[-1, -1, -1]$, θ is set as $[1, 1, 1]$, and τ is set as 2, which results in $ATE = -1$. For `MixedCorrelator`, we set the same γ and θ as `MixedMediator`, where parameters $\phi = 1$ and $\tau = 1$, which results in $ATE = 1$.

In addition, we build a real-world dataset based on the job Ads data from the Company, aiming to estimate the ATE of *switching a job from onsite to online working mode to the statistics of the applicants* (here we choose the average age as the outcome). In the dataset, treatment T represents the working mode of the job, where $T = 1$ represents the job is online, whereas $T = 0$ represents the job is onsite, Y is the standardized age, and $\mathbf{X} \in \{0, 1\}^{K_X}$ indicates the required skills of the job. We select 3,228 jobs from Bay Area, where a primary study shows that $DEV(\emptyset) \approx -2$ years⁵ (i.e., online job applicants are two years younger than onsite job applicants). To simulate the latent confounder C and post-treatment variables M , we first learn a generative model as follows:

$$\mathbf{Z} \sim \text{Gaussian}(\mathbf{0}, \mathbf{I}_{K_Z}), \mathbf{X} \sim \text{Multi}(NN_f(\mathbf{Z})), Y \sim \text{Gaussian}(\mathbf{w} \odot \mathbf{Z}, 1) \quad (12)$$

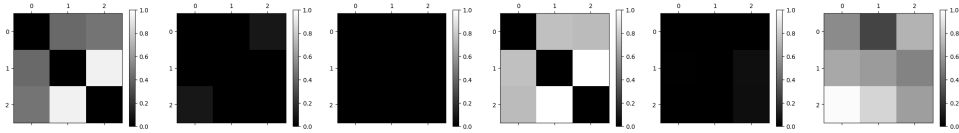
where Multi represents multinomial distribution, NN_f is a neural network with softmax activation, $\mathbf{Z}, \mathbf{w} \in \mathbb{R}^{K_Z}$, $K_Z = 6$, and \odot represents the element-wise product operator, respectively. We then treat the first $K_C = 3$ dimensions of \mathbf{Z} as the latent confounders C and the remaining $K_M = K_Z - K_C$ dimensions as the latent mediators M . After learning NN_f and \mathbf{w} according to Eq. (12), we draw latent confounders $C \in \text{Gaussian}(0, \mathbf{I})$, latent mediators $M = T \cdot \gamma$, and set the outcome $Y = \mathbf{w} \odot [C || M] + \tau \cdot T$, where the true ATE can be calculated as $\text{sum}(\gamma \odot \mathbf{w}_{-K_M}) + \tau$.

4.2 COMPARISONS WITH THE STATE-OF-THE-ART

The baselines we include for comparisons can be categorized into three classes. **1) Unawareness**, where no information in \mathbf{X} is used for ATE estimation. We implement the naive LR0 estimator, which regresses Y on T and uses the coefficient to estimate the ATE (Imbens & Rubin, 2015) (LR0 equals to $DEV(\emptyset)$, i.e., the difference of average outcome between the treatment and non-treatment group). **2) Control- \mathbf{X}** , which directly controls the covariates \mathbf{X} . In this class, LR1 regresses Y on T and \mathbf{X} , whereas TarNet uses a two-branch neural network to estimate the $DEV(\mathbf{X})$ **3) Control- \mathbf{Z}** , which controls latent variables \mathbf{Z} learned from the covariates \mathbf{X} . Methods from this class include the CEVAE (Louizos et al., 2017) and covariate disentanglement methods (see Fig. 1-(b)), such as DR-CFR Hassanpour & Greiner (2020) and TEDVAE (Zhang et al., 2021).

The comparisons are summarized in Table 1. From Table 1, we can empirically verify the correctness of Theorem 2.2 that post-treatment bias indeed poses a serious issue for proxy-of-confounder-based methods, because for the `MixedMediator` and `MixedConfounder` datasets, CEVAE is worse than the naive LR0 estimator that directly calculates the difference of mean outcome between the

⁵which leads to -0.178 after standardization. Code demo see <https://anonymous.4open.science/r/CiVAE-demo-54B9>.



(a) Case 1: Intra-Confounder (b) Case 2: Intra-Mediator (c) Case 3: Confounder-Mediator

Figure 2: Visualization of p -value of independence test before and after conditioning on treatment T .

Table 1: Comparison of CiVAE with baselines on ATE estimation with latent post-treatment bias.

Dataset	MixedMediator		MixedCorrelator		Company	
Method	ATE.	Err.	ATE.	Err.	ATE.	Err.
LR0	0.975 ± 0.032	1.975	2.977 ± 0.032	1.977	0.131 ± 0.015	0.399
LR1	1.457 ± 0.167	2.457	3.400 ± 0.130	2.400	0.093 ± 0.071	0.361
TarNet	1.461 ± 0.172	2.461	3.414 ± 0.146	2.414	0.112 ± 0.085	0.380
CEVAE	1.550 ± 0.292	2.550	3.323 ± 0.167	2.323	0.106 ± 0.078	0.374
DR-CFR	1.239 ± 0.324	2.239	3.185 ± 0.319	2.185	0.094 ± 0.089	0.362
TEDVAE	1.042 ± 0.315	2.042	3.138 ± 0.281	2.138	0.097 ± 0.093	0.365
CiVAE	-0.822 ± 0.753	0.178	1.199 ± 0.765	0.199	-0.140 ± 0.137	-0.128
True ATE	-1.000 ± 0.000	0.000	1.000 ± 0.000	0.000	-0.268 ± 0.000	0.000

treatment and non-treatment groups. In addition, for `MixedMediator` and `Company` datasets, all methods except the proposed CiVAE fail to predict the negativity of the ATE.

Covariates disentanglement-based methods, i.e., DR-CFR and TEDVAE, achieve similar performance as CEVAE. The reason is that, these methods disentangle latent confounders C from latent instrumental variables I and latent adjusters A by utilizing their causal relations with T and Y , i.e., I is predictive only for T , A is predictive only for Y , whereas C is predictive for both T and Y . For example, TEDVAE includes three encoders to infer three sets of latent variables $\hat{I}, \hat{A}, \hat{C}$ from X and adds classification losses $p(T|\hat{I}, \hat{C})$ and $p(Y|T, \hat{C}, \hat{A})$ on the CEVAE loss. However, when latent post-treatment bias exists, since both latent confounders C and latent post-treatment variables M are correlated with both T and Y , \hat{C} inferred by TEDVAE still cannot disentangle C from M .

CiVAE achieves significantly better results compared to CEVAE and TEDVAE, which demonstrates its effectiveness in identifying and distinguishing latent confounders from post-treatment variables in proxies. However, we also notice that a downside of CiVAE is the comparatively large variance across ten dataset splits, as misidentifying latent mediators as confounders may result in severe performance degradation when the mediation effects are strong or the number of latent confounders is small.

4.3 DISENTANGLING OF LATENT CONFOUNDERS AND POST-TREATMENT VARIABLES

We show the p -value of the pairwise independence test of the true latent variables before and after conditioning on the assigned treatment T . From Fig. 2 we can find that the difference between the three cases discussed in Subsection 3.4 is significant. Here, we should note that the distinction of the intra-confounder case from other cases relies on the assumption that latent confounders are independent. If the latent confounders are correlated, we can first use causal discovery techniques such as the PC algorithm (Spirtes et al., 2000) to find direct parents of T , and use our algorithm as the refinement to determine the true confounders C from the misidentified post-treatment variables.

5 CONCLUSIONS

In this paper, we systematically investigated the latent post-treatment bias in causal inference from observational data. We first prove that unresolved latent post-treatment variables scrambled in the proxy of confounders can arbitrarily bias the ATE estimation. To address the bias, we proposed the Confounder-identifiable VAE (CiVAE), which, utilizing a mild assumption regarding the prior of latent factors, guarantees the identifiability of latent confounders up to bijective transformations. Finally, we show that controlling the latent confounders inferred by CiVAE can provide an unbiased estimation of the ATE. Experiments on both simulated and real-world datasets demonstrated that CiVAE has superior robustness to latent post-treatment bias compared with state-of-the-art methods.

REFERENCES

- Avidit Acharya, Matthew Blackwell, and Maya Sen. Explaining causal findings without bias: Detecting and assessing direct effects. *American Political Science Review*, 110(3):512–529, 2016.
- Animashree Anandkumar, Rong Ge, Daniel Hsu, Sham M Kakade, and Matus Telgarsky. Tensor decompositions for learning latent variable models. *Journal of Machine Learning Research*, 15: 2773–2832, 2014.
- David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- Lu Cheng, Ruocheng Guo, and Huan Liu. Causal mediation analysis with hidden confounders. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*, pp. 113–122, 2022.
- Thomas D Cook, Donald Thomas Campbell, and William Shadish. *Experimental and quasi-experimental designs for generalized causal inference*. Houghton Mifflin Boston, MA, 2002.
- Jessie K Edwards, Stephen R Cole, and Daniel Westreich. All your data are always missing: incorporating bias due to measurement error into the potential outcomes framework. *International Journal of Epidemiology*, 44(4):1452–1459, 2015.
- Felix Elwert and Christopher Winship. Endogenous selection bias: The problem of conditioning on a collider variable. *Annual Review of Sociology*, 40:31–53, 2014.
- Thomas A Glass, Steven N Goodman, Miguel A Hernán, and Jonathan M Samet. Causal inference in public health. *Annual Review of Public Health*, 34:61–75, 2013.
- Madelyn Glymour, Judea Pearl, and Nicholas P Jewell. *Causal inference in statistics: A primer*. John Wiley & Sons, 2016.
- Negar Hassanpour and Russell Greiner. Learning disentangled representations for counterfactual regression. In *International Conference on Learning Representations*, 2020.
- Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005.
- Kosuke Imai, Luke Keele, and Dustin Tingley. A general approach to causal mediation analysis. *Psychological Methods*, 15(4):309, 2010.
- Guido W Imbens and Donald B Rubin. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press, 2015.
- Fredrik Johansson, Uri Shalit, and David Sontag. Learning representations for counterfactual inference. In *International Conference on Machine Learning*, pp. 3020–3029, 2016.
- 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. PMLR, 2020.
- Gary King and Langche Zeng. The dangers of extreme counterfactuals. *Political Analysis*, 14(2): 131–159, 2006.
- Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- Manabu Kuroki and Judea Pearl. Measurement bias and effect restoration in causal inference. *Biometrika*, 101(2):423–437, 2014.
- Fan Li, Kari Lock Morgan, and Alan M Zaslavsky. Balancing covariates via propensity score weighting. *Journal of the American Statistical Association*, 113(521):390–400, 2018.
- Christos Louizos, Uri Shalit, Joris M Mooij, David Sontag, Richard Zemel, and Max Welling. Causal effect inference with deep latent-variable models. *Advances in Neural Information Processing Systems*, 30, 2017.

- Chaochao Lu, Yuhuai Wu, José Miguel Hernández-Lobato, and Bernhard Schölkopf. Invariant causal representation learning for out-of-distribution generalization. In *International Conference on Learning Representations*, 2021.
- Yulong Lu and Jianfeng Lu. A universal approximation theorem of deep neural networks for expressing probability distributions. In *Advances in Neural Information Processing Systems*, pp. 3094–3105, 2020.
- David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. Fairness through causal awareness: Learning causal latent-variable models for biased data. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pp. 349–358, 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.
- Ruth M Mickey and Sander Greenland. The impact of confounder selection criteria on effect estimation. *American Journal of Epidemiology*, 129(1):125–137, 1989.
- Jacob M Montgomery, Brendan Nyhan, and Michelle Torres. How conditioning on posttreatment variables can ruin your experiment and what to do about it. *American Journal of Political Science*, 62(3):760–775, 2018.
- Judea Pearl. On measurement bias in causal inference. *arXiv preprint arXiv:1203.3504*, 2012.
- Stuart J Pocock, Susan E Assmann, Laura E Enos, and Linda E Kasten. Subgroup analysis, covariate adjustment and baseline comparisons in clinical trial reporting: current practice and problems. *Statistics in Medicine*, 21(19):2917–2930, 2002.
- Mattia Proserpi, Yi Guo, Matt Sperrin, James S Koopman, Jae S Min, Xing He, Shannan Rich, Mo Wang, Iain E Buchan, and Jiang Bian. Causal inference and counterfactual prediction in machine learning for actionable healthcare. *Nature Machine Intelligence*, 2(7):369–375, 2020.
- Kenneth J Rothman, Sander Greenland, Timothy L Lash, et al. *Modern epidemiology*, volume 3. Wolters Kluwer Health/Lippincott Williams & Wilkins Philadelphia, 2008.
- 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, 2017.
- Claudia Shi, David Blei, and Victor Veitch. Adapting neural networks for the estimation of treatment effects. In *Advances in Neural Information Processing Systems*, 2019.
- Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, prediction, and search*. MIT press, 2000.
- 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.
- 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*, volume 31, 2018.
- Kun Zhang, Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Kernel-based conditional independence test and application in causal discovery. *arXiv preprint arXiv:1202.3775*, 2012.
- Weijia Zhang, Lin Liu, and Jiuyong Li. Treatment effect estimation with disentangled latent factors. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 10923–10930, 2021.

A PROOFS

A.1 PROOF OF LEMMA 2.1.

Proof. Let $\mathbf{Z} = g(\mathbf{X})$ and $\mathbf{z} = g(\mathbf{x})$. If g is injective and differentiable *a.e.*, and g^\dagger is the left-inverse, we have:

$$f_{Y|g(\mathbf{X})}(y|g(\mathbf{x})) = f_{Y|\mathbf{Z}}(y|\mathbf{z}) = \frac{f_{Y,\mathbf{Z}}(y,\mathbf{z})}{f_{\mathbf{Z}}(\mathbf{z})} = \frac{f_{Y,\mathbf{X}}(y,g^\dagger(\mathbf{z}))|\mathbf{J}_{g^\dagger}(\mathbf{z})|}{f_{\mathbf{X}}(g^\dagger(\mathbf{z}))|\mathbf{J}_{g^\dagger}(\mathbf{z})|} = \frac{f_{Y,\mathbf{X}}(y,\mathbf{x})}{f_{\mathbf{X}}(\mathbf{x})} = f_{Y|\mathbf{X}}(y|\mathbf{x}), \quad (13)$$

where f and $f_{\cdot|\cdot}$ represent the marginal and conditional density function, respectively, and $\mathbf{J}_{g^\dagger}(\mathbf{z})$ is the Jacobian matrix of function g^\dagger evaluated at \mathbf{z} . Based on Eq. (13), we have:

$$\mathbb{E}[Y|\mathbf{X}] = \int \mathbf{y} \cdot f_{Y|\mathbf{X}}(\mathbf{y}|\mathbf{x}) d\mathbf{y} = \int \mathbf{y} \cdot f_{Y|\mathbf{Z}}(\mathbf{y}|\mathbf{z}) d\mathbf{y} = \mathbb{E}[Y|\mathbf{Z} = \mathbf{z}] = \mathbb{E}[Y|g(\mathbf{X}) = g(\mathbf{x})]. \quad (14)$$

□

A.2 PROOF OF COROLLARY 2.3.

Proof. For $\mathbf{X} = \mathbf{x}$, let $[\mathbf{c}|\mathbf{m}] \doteq [f_C^\dagger(\mathbf{x})|f_M^\dagger(\mathbf{x})] \doteq f^\dagger(\mathbf{x}) = \mathbf{A}^\dagger(\mathbf{x} - \boldsymbol{\alpha}_X)$, where \mathbf{A}^\dagger is the left inverse of the full column-rank matrix \mathbf{A} in Eq. (2), we have:

$$\begin{aligned} CATE(\mathbf{x}) &= \mathbb{E}[Y|T = 1, \mathbf{C} = f_C^\dagger(\mathbf{x})] - \mathbb{E}[Y|T = 0, \mathbf{C} = f_C^\dagger(\mathbf{x})] \\ &= \mathbb{E}[Y|T = 1, \mathbf{C} = \mathbf{c}] - \mathbb{E}[Y|T = 0, \mathbf{C} = \mathbf{c}] \\ &= \mathbb{E}[\alpha_Y + \tau \cdot T + \sum \theta_j \cdot M_j + \sum \kappa_i \cdot C_i | T = 1, \mathbf{C} = \mathbf{c}] \\ &\quad - \mathbb{E}[\alpha_Y + \tau \cdot T + \sum \theta_j \cdot M_j + \sum \kappa_i \cdot C_i | T = 0, \mathbf{C} = \mathbf{c}] \\ &= \alpha_Y + \tau \cdot \mathbb{E}[T | T = 1, \mathbf{C} = \mathbf{c}] + \sum \theta_j \cdot \mathbb{E}[M_j | T = 1, \mathbf{C} = \mathbf{c}] + \sum \kappa_i \cdot \mathbb{E}[C_i | T = 1, \mathbf{C} = \mathbf{c}] \\ &\quad - \alpha_Y + \tau \cdot \mathbb{E}[T | T = 0, \mathbf{C} = \mathbf{c}] + \sum \theta_j \cdot \mathbb{E}[M_j | T = 0, \mathbf{C} = \mathbf{c}] + \sum \kappa_i \cdot \mathbb{E}[C_i | T = 0, \mathbf{C} = \mathbf{c}] \\ &= \tau \cdot (1 - 0) + \sum \theta_j \cdot (\gamma_j \cdot (1 - 0)) + \sum \kappa_i \cdot (c_i - c_i) \\ &= \tau + \sum \theta_j \cdot \gamma_j = \mathbb{E}[\tau + \sum \theta_j \cdot \gamma_j] = ATE, \end{aligned} \quad (15)$$

where the first equality is due to the definition of CATE in Eq. (2). In addition, the causal estimand and bias of a proxy-of-confounder-based causal inference model that controls the latent variable \mathbf{Z} inferred via $\hat{\mathbf{Z}} = \tilde{f}(\mathbf{X}) = \mathbf{B}^T \mathbf{X}$ (where \mathbf{B} is also a full column-rank matrix) can be formulated as:

$$\begin{aligned} DCEV(\mathbf{B}^T \mathbf{x}) &= \mathbb{E}[Y|T = 1, \hat{\mathbf{Z}} = \mathbf{B}^T \mathbf{x}] - \mathbb{E}[Y|T = 0, \hat{\mathbf{Z}} = \mathbf{B}^T \mathbf{x}] \\ &= \mathbb{E}[Y|T = 1, \hat{\mathbf{Z}} = \mathbf{B}^T \boldsymbol{\alpha}_X + \mathbf{B}^T \mathbf{A}[\mathbf{c}|\mathbf{m}]] - \mathbb{E}[Y|T = 0, \hat{\mathbf{Z}} = \mathbf{B}^T \boldsymbol{\alpha}_X + \mathbf{B}^T \mathbf{A}[\mathbf{c}|\mathbf{m}]] \\ &\stackrel{(a)}{=} \mathbb{E}[Y|T = 1, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] - \mathbb{E}[Y|T = 0, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] \\ &= \alpha_Y + \tau \cdot 1 + \sum \theta_j \cdot \mathbb{E}[M_j | T = 1, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] + \sum \kappa_i \cdot \mathbb{E}[C_i | T = 1, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] \\ &\quad - \alpha_Y + \tau \cdot 0 + \sum \theta_j \cdot \mathbb{E}[M_j | T = 0, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] + \sum \kappa_i \cdot \mathbb{E}[C_i | T = 0, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] \\ &= \tau \cdot (1 - 0) + \sum \theta_j \cdot (m_j - m_j) + \sum \kappa_i \cdot (c_i - c_i) \\ &= \tau = \mathbb{E}[\tau] = \mathbb{E}[DCEV(\mathbf{B}^T \mathbf{X})], \end{aligned} \quad (16)$$

where step (a) is due to the fact that, since both \mathbf{A} and \mathbf{B} are full column-rank matrices, $\mathbf{B}^T \mathbf{A}$ is an invertible matrix, and the mapping $\tilde{f} = \mathbf{B}^T \boldsymbol{\alpha}_X + \mathbf{B}^T \mathbf{A}$ is bijective. Therefore, we can invoke Lemma 2.1 and apply the left-inverse of \tilde{f} , i.e., $\tilde{f}^\dagger = (\mathbf{B}^T \mathbf{A})^{-1} - \mathbf{B}^T \boldsymbol{\alpha}_X$, to the condition of the expectation. The rest steps are based on the structural causal equations defined in Eq. (2). □

A.3 ANOTHER CASE OF LINEAR SCM WITH LATENT CORRELATORS

Corollary A.1. (*MixedCorrelator*) For another Linear Structural Causal Model defined as follows:

$$\begin{aligned}
T &\leftarrow \mathbb{1}(\alpha_T + \sum \beta_i \cdot C_i > a) \\
M_j &\leftarrow \alpha_M + \gamma_j \cdot T + \phi_j \cdot U_j \\
\mathbf{X} &\leftarrow \alpha_X + \mathbf{A}[\mathbf{M}|\mathbf{C}] \\
Y &\leftarrow \alpha_Y + \tau \cdot T + \sum \theta_j \cdot U_j + \sum \kappa_i \cdot C_i,
\end{aligned} \tag{17}$$

where the mixture function $f = \mathbf{A} \in \mathbb{R}^{K_X \times (K_C + K_M)}$ is a full column-rank matrix, the CATE, ATE, and the bias of proxy-of-confounder-based causal inference model that controls the latent variable $\hat{\mathbf{Z}}$ inferred via $\hat{\mathbf{Z}} = \hat{f}(\mathbf{X}) = \mathbf{B}^T \mathbf{X}$ can be formulated as follows:

$$\begin{aligned}
ATE &= CATE = \tau \\
DEV(\hat{\mathbf{Z}}) &= \mathbb{E}[DCEV(\hat{\mathbf{Z}})] = DCEV(\hat{\mathbf{Z}}) = \tau - \sum \frac{\theta_j \cdot \gamma_j}{\phi_j} \\
Bias &= ATE - DEV(\mathbf{B}^T \mathbf{X}) = \sum \frac{\theta_j \cdot \gamma_j}{\phi_j},
\end{aligned} \tag{18}$$

where $\mathbf{B} \in \mathbb{R}^{K_X \times (K_C + K_M)}$ is another full column-rank matrix. Since $\sum \frac{\theta_j \cdot \gamma_j}{\phi_j}$ is arbitrary, the estimator $DEV(\hat{\mathbf{Z}}) = \mathbb{E}[DCEV(\mathbf{B}^T \mathbf{X})]$ is arbitrarily biased for the estimation of ATE.

Proof. The proof of the CATE and ATE is trivial. The causal estimand and the bias of a proxy-of-confounder-based causal inference model that controls the latent variables $\hat{\mathbf{Z}}$ inferred via $\hat{\mathbf{Z}} = \hat{f}(\mathbf{X}) = \mathbf{B}^T \mathbf{X}$ (where \mathbf{B} is also a full column-rank matrix) can be formulated as follows:

$$\begin{aligned}
DCEV(\mathbf{B}^T \mathbf{x}) &= \mathbb{E}[Y|T=1, \hat{\mathbf{Z}} = \mathbf{B}^T \mathbf{x}] - \mathbb{E}[Y|T=0, \hat{\mathbf{Z}} = \mathbf{B}^T \mathbf{x}] \\
&= \mathbb{E}[Y|T=1, \hat{\mathbf{Z}} = \alpha_X + \mathbf{B}^T \mathbf{A}[\mathbf{c}|\mathbf{m}]] - \mathbb{E}[Y|T=0, \hat{\mathbf{Z}} = \alpha_X + \mathbf{B}^T \mathbf{A}[\mathbf{c}|\mathbf{m}]] \\
&\stackrel{(a)}{=} \mathbb{E}[Y|T=1, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] - \mathbb{E}[Y|T=0, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] \\
&= \alpha_Y + \tau \cdot 1 + \sum \theta_j \cdot \mathbb{E}[U_j|T=1, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] + \sum \kappa_i \cdot \mathbb{E}[C_i|T=1, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] \\
&\quad - \alpha_Y + \tau \cdot 0 + \sum \theta_j \cdot \mathbb{E}[U_j|T=0, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] + \sum \kappa_i \cdot \mathbb{E}[C_i|T=0, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] \\
&= \tau \cdot (1 - 0) + \sum \theta_j \cdot (\phi_j^{-1} \cdot (m_j - \alpha_M - \gamma_j) - \phi_j^{-1} \cdot (m_j - \alpha_M)) + \sum \kappa_i \cdot (c_i - c_i) \\
&= \tau - \sum \frac{\theta_j \cdot \gamma_j}{\phi_j} = \mathbb{E} \left[\tau - \sum \frac{\theta_j \cdot \gamma_j}{\phi_j} \right] = \mathbb{E}[DCEV(\mathbf{B}^T \mathbf{X})],
\end{aligned} \tag{19}$$

□

where step (a) and the rest of the proof follow the same logic as the proof in Section 2.3.

A.4 PROOF OF THEOREM 3.1

The strict definitions of the exponential family, strong exponential (which is assumed for the factorized part of the conditional prior), and identifiability follow (Khemakhem et al., 2020; Lu et al., 2021), and can be referred to in Appendix E, F of (Lu et al., 2021), which we omit to avoid redundancy. The proof of Theorem 3.1 is largely based on the NF-iVAE paper (Lu et al., 2021), where most of the details can be found, with the new assumption introduced in CiVAE that each $\mathcal{S}_{f,i}$ has at least one invertible dimension incorporated to ensure that each dimension of the inferred latent variables is a bijective transformation of the corresponding true latent variable.

A.4.1 PART I

Step I. In this step, we transform the equality of noisy conditional marginal distribution of \mathbf{X} given Y, T of two models with parameter $\theta, \hat{\theta} \in \Theta$ into the equality of noise-free distributions.

$$\begin{aligned}
& p_{\theta}(\mathbf{X} | Y, T) = p_{\hat{\theta}}(\mathbf{X} | Y, T) \\
& \implies \int_{\mathcal{Z}} p_f(\mathbf{X} | \mathbf{Z}) p_{S, \lambda}(\mathbf{Z} | Y, T) d\mathbf{Z} = \int_{\mathcal{Z}} p_{\tilde{f}}(\mathbf{X} | \mathbf{Z}) p_{\tilde{S}, \tilde{\lambda}}(\mathbf{Z} | Y, T) d\mathbf{Z} \\
& \implies \int_{\mathcal{Z}} p_{\varepsilon}(\mathbf{X} - f(\mathbf{Z})) p_{S, \lambda}(\mathbf{Z} | Y, T) d\mathbf{Z} = \int_{\mathcal{Z}} p_{\varepsilon}(\mathbf{X} - \tilde{f}(\mathbf{Z})) p_{\tilde{S}, \tilde{\lambda}}(\mathbf{Z} | Y, T) d\mathbf{Z} \\
& \stackrel{(a)}{\implies} \int_{\mathcal{X}} p_{\varepsilon}(\mathbf{X} - \bar{\mathbf{X}}) p_{S, \lambda}(f^{\dagger}(\bar{\mathbf{X}}) | Y, T) \text{vol}(\mathbf{J}_{f^{\dagger}}(\bar{\mathbf{X}})) d\bar{\mathbf{X}} = \\
& \quad \int_{\mathcal{X}} p_{\varepsilon}(\mathbf{X} - \bar{\mathbf{X}}) p_{\tilde{S}, \tilde{\lambda}}(\tilde{f}^{\dagger}(\bar{\mathbf{X}}) | Y, T) \text{vol}(\mathbf{J}_{\tilde{f}^{\dagger}}(\bar{\mathbf{X}})) d\bar{\mathbf{X}} \\
& \stackrel{(b)}{\implies} \int_{\mathbb{R}^d} p_{\varepsilon}(\mathbf{X} - \bar{\mathbf{X}}) \tilde{p}_{f, S, \lambda, Y, T}(\bar{\mathbf{X}}) d\bar{\mathbf{X}} = \int_{\mathbb{R}^d} p_{\varepsilon}(\mathbf{X} - \bar{\mathbf{X}}) \tilde{p}_{\tilde{f}, \tilde{S}, \tilde{\lambda}, \tilde{Y}, \tilde{T}}(\bar{\mathbf{X}}) d\bar{\mathbf{X}} \\
& \implies (\tilde{p}_{f, S, \lambda, Y, T} * p_{\varepsilon})(\mathbf{X}) = (\tilde{p}_{\tilde{f}, \tilde{S}, \tilde{\lambda}, \tilde{Y}, \tilde{T}} * p_{\varepsilon})(\mathbf{X}) \\
& \stackrel{(c)}{\implies} F[\tilde{p}_{f, S, \lambda, Y, T}](\boldsymbol{\omega}) \varphi_{\varepsilon}(\boldsymbol{\omega}) = F[\tilde{p}_{\tilde{f}, \tilde{S}, \tilde{\lambda}, \tilde{Y}, \tilde{T}}](\boldsymbol{\omega}) \varphi_{\varepsilon}(\boldsymbol{\omega}) \\
& \stackrel{(d)}{\implies} F[\tilde{p}_{f, S, \lambda, Y, T}](\boldsymbol{\omega}) = F[\tilde{p}_{\tilde{f}, \tilde{S}, \tilde{\lambda}, \tilde{Y}, \tilde{T}}](\boldsymbol{\omega}) \\
& \implies \tilde{p}_{f, S, \lambda, Y, T}(\mathbf{X}) = \tilde{p}_{\tilde{f}, \tilde{S}, \tilde{\lambda}, \tilde{Y}, \tilde{T}}(\mathbf{X}).
\end{aligned} \tag{20}$$

Step (a) is based on the rule of change-of-variable, where $\text{vol}(\mathbf{A}) = \sqrt{\det(\mathbf{A}^T \mathbf{A})}$. In step (b), we define $\tilde{p}_{f, S, \lambda, Y, T}(\mathbf{X}) \triangleq p_{S, \lambda}(f^{\dagger}(\mathbf{X}) | Y, T) \text{vol}(\mathbf{J}_{f^{\dagger}}(\mathbf{X})) \mathbb{I}_{\mathcal{X}}(\mathbf{X})$. In step (c), we use $F[\cdot]$ to denote the Fourier transform. In step (d), we drop $\varphi_{\varepsilon}(\boldsymbol{\omega})$ as it is non-zero *a.e.* (see Assumption 3).

Step II. In this step, we transform the equality of the noise-free distributions into the relationship of the sufficient statistics \mathbf{S} and $\tilde{\mathbf{S}}$. By taking logarithm of both sides of Eq. (20), we have:

$$\begin{aligned}
& \log \text{vol}(\mathbf{J}_{f^{\dagger}}(\mathbf{X})) + \log \mathcal{Q}(f^{\dagger}(\mathbf{X})) - \log \mathcal{C}(Y, T) + \langle \mathbf{S}(f^{\dagger}(\mathbf{X})), \boldsymbol{\lambda}(Y, T) \rangle \\
& = \log \text{vol}(\mathbf{J}_{\tilde{f}^{\dagger}}(\mathbf{X})) + \log \tilde{\mathcal{Q}}(\tilde{f}^{\dagger}(\mathbf{X})) - \log \tilde{\mathcal{C}}(Y, T) + \langle \tilde{\mathbf{S}}(\tilde{f}^{\dagger}(\mathbf{X})), \tilde{\boldsymbol{\lambda}}(Y, T) \rangle.
\end{aligned} \tag{21}$$

Let $(Y, T)_0, \dots, (Y, T)_k$ be the $k+1$ distinct points defined in Assumption 3 - (iv). We obtain $k+1$ equations by evaluating the Eq. (21) at these points, where the first equation is subtracted from the remaining ones, which leads to the following equation system:

$$\begin{aligned}
& \langle \mathbf{S}(f^{\dagger}(\mathbf{X})), \boldsymbol{\lambda}((Y, T)_l) - \boldsymbol{\lambda}((Y, T)_0) \rangle + \log \frac{\mathcal{C}((Y, T)_0)}{\mathcal{C}((Y, T)_l)} \\
& = \langle \tilde{\mathbf{S}}(\tilde{f}^{\dagger}(\mathbf{X})), \tilde{\boldsymbol{\lambda}}((Y, T)_l) - \tilde{\boldsymbol{\lambda}}((Y, T)_0) \rangle + \log \frac{\tilde{\mathcal{C}}((Y, T)_0)}{\tilde{\mathcal{C}}((Y, T)_l)}, \quad l = 1, \dots, k.
\end{aligned} \tag{22}$$

Let \mathbf{L} be the invertible matrix defined in Assumption 3 - (iv) and $\tilde{\mathbf{L}}$ be the counterpart for $\tilde{\boldsymbol{\lambda}}$, if we summarize all terms irrelevant to \mathbf{X} into a constant \mathbf{b} , we have:

$$\begin{aligned}
& \mathbf{L}^T \mathbf{S}(f^{\dagger}(\mathbf{X})) = \tilde{\mathbf{L}}^T \tilde{\mathbf{S}}(\tilde{f}^{\dagger}(\mathbf{X})) + \mathbf{b} \\
& \implies \mathbf{S}(f^{\dagger}(\mathbf{X})) = \mathbf{A} \tilde{\mathbf{S}}(\tilde{f}^{\dagger}(\mathbf{X})) + \mathbf{c},
\end{aligned} \tag{23}$$

where $\mathbf{A} = \mathbf{L}^{-T} \tilde{\mathbf{L}} \in \mathbb{R}^{k \times k}$, and $\mathbf{c} = \mathbf{L}^{-T} \mathbf{b} \in \mathbb{R}^k$.

Step III. Ideally, to prove the element-wise bijective identifiability of the latent variables \mathbf{Z} , the transformation of the sufficient statistics \mathbf{S} derived in Eq. (23) should be bijective. We claim that if the conditional prior $p_{S, \lambda}(\mathbf{Z} | Y, T)$ is strongly exponential and \mathbf{L} is invertible, $\tilde{\mathbf{L}}$ and \mathbf{A} must also be invertible. The proof is omitted, and can be referred to in Appendix H.1.1 of (Lu et al., 2021).

A.4.2 PART II

In this part, we prove that, if Assumptions 1, 2 and 3 hold, we can identify the factorized part of the sufficient statistics $\mathbf{S}(\mathbf{Z})$, i.e., $\mathbf{S}_f(\mathbf{Z})$, up to permutation and element-wise transformation. Specifically, if we use \mathbf{v} to denote the composite map $\tilde{f}^\dagger \circ f : \mathcal{Z} \rightarrow \mathcal{Z}$, Eq. (23) can be rewritten into:

$$\mathbf{S}(\mathbf{Z}) = \mathbf{A}\tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z})) + \mathbf{c}. \quad (24)$$

We aim to prove that \mathbf{A} in Eq. (24) is a block permutation matrix.

Step I. We start by showing that \mathbf{v} is a component-wise function. If we differentiate both sides of Eq. (24) with respect to Z_s and Z_t , where $s \neq t$, we have:

$$\begin{aligned} \frac{\partial \mathbf{S}(\mathbf{Z})}{\partial Z_s} &= \mathbf{A} \sum_{i=1}^{K_Z} \frac{\partial \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z})} \cdot \frac{\partial v_i(\mathbf{Z})}{\partial Z_s} \\ \frac{\partial^2 \mathbf{S}(\mathbf{Z})}{\partial Z_s \partial Z_t} &= \mathbf{A} \sum_{i=1}^{K_Z} \sum_{j=1}^{K_Z} \frac{\partial^2 \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z}) \partial v_j(\mathbf{Z})} \cdot \frac{\partial v_j(\mathbf{Z})}{\partial Z_t} \cdot \frac{\partial v_i(\mathbf{Z})}{\partial Z_s} + \mathbf{A} \sum_{i=1}^{K_Z} \frac{\partial \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z})} \cdot \frac{\partial^2 v_i(\mathbf{Z})}{\partial Z_s \partial Z_t}. \end{aligned} \quad (25)$$

Note that for the factorized part of the sufficient statistics \mathbf{S} , i.e., \mathbf{S}_f , all *cross-derivatives* are zero, and for the non-factorized part of \mathbf{S} , i.e., \mathbf{S}_{nf} , which is a neural network with ReLU activation (i.e., linear *a.e.*), all *second-order derivatives* are zero. Therefore, the *second order cross-derivatives* on the LHS. of Eq. (25) are zero, which leads to the following equality:

$$\mathbf{0} = \mathbf{A} \sum_{i=1}^{K_Z} \frac{\partial^2 \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z})^2} \cdot \frac{\partial v_i(\mathbf{Z})}{\partial Z_t} \cdot \frac{\partial v_i(\mathbf{Z})}{\partial Z_s} + \mathbf{A} \sum_{i=1}^{K_Z} \frac{\partial \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z})} \cdot \frac{\partial^2 v_i(\mathbf{Z})}{\partial Z_s \partial Z_t}. \quad (26)$$

Eq. (26) can be written into the matrix-vector product form as follows:

$$\mathbf{0} = \mathbf{A}\tilde{\mathbf{S}}''(\mathbf{Z})\mathbf{v}'_{s,t}(\mathbf{Z}) + \mathbf{A}\tilde{\mathbf{S}}'(\mathbf{Z})\mathbf{v}''_{s,t}(\mathbf{Z}), \quad (27)$$

where

$$\begin{aligned} \tilde{\mathbf{S}}''(\mathbf{Z}) &= \left[\frac{\partial^2 \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_1(\mathbf{Z})^2}, \dots, \frac{\partial^2 \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_{K_Z}(\mathbf{Z})^2} \right] \in \mathbb{R}^{k \times K_Z}, \\ \mathbf{v}'_{s,t}(\mathbf{Z}) &= \left[\frac{\partial v_1(\mathbf{Z})}{\partial Z_t} \cdot \frac{\partial v_1(\mathbf{Z})}{\partial Z_s}, \dots, \frac{\partial v_{K_Z}(\mathbf{Z})}{\partial Z_t} \cdot \frac{\partial v_{K_Z}(\mathbf{Z})}{\partial Z_s} \right]^T \in \mathbb{R}^{K_Z}, \end{aligned}$$

and

$$\begin{aligned} \tilde{\mathbf{S}}'(\mathbf{Z}) &= \left[\frac{\partial \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_1(\mathbf{Z})}, \dots, \frac{\partial \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_{K_Z}(\mathbf{Z})} \right] \in \mathbb{R}^{k \times K_Z}, \\ \mathbf{v}''_{s,t}(\mathbf{Z}) &= \left[\frac{\partial^2 v_1(\mathbf{Z})}{\partial Z_s \partial Z_t}, \dots, \frac{\partial^2 v_{K_Z}(\mathbf{Z})}{\partial Z_s \partial Z_t} \right]^T \in \mathbb{R}^{K_Z}. \end{aligned}$$

If we denote the concatenation as $\tilde{\mathbf{S}}'''(\mathbf{Z}) = \left[\tilde{\mathbf{S}}''(\mathbf{Z}), \tilde{\mathbf{S}}'(\mathbf{Z}) \right] \in \mathbb{R}^{k \times 2K_Z}$ and $\mathbf{v}''_{s,t}(\mathbf{Z}) = \left[\mathbf{v}'_{s,t}(\mathbf{Z})^T, \mathbf{v}''_{s,t}(\mathbf{Z})^T \right]^T \in \mathbb{R}^{2K_Z}$, we have:

$$\mathbf{0} = \mathbf{A}\tilde{\mathbf{S}}'''(\mathbf{Z})\mathbf{v}''_{s,t}(\mathbf{Z}). \quad (28)$$

Finally, if we denote the rows of $\tilde{\mathbf{S}}'''(\mathbf{Z})$ that correspond to the factorized part of \mathbf{S} by $\tilde{\mathbf{S}}'''_f(\mathbf{Z})$, according to Lemma 5 of the iVAE paper (Khemakhem et al., 2020) and the assumption that $k \geq 2K_Z$, we have that the rank of $\tilde{\mathbf{S}}'''_f(\mathbf{Z})$ is $2K_Z$. Since $k \geq 2K_Z$, the rank of $\tilde{\mathbf{S}}'''(\mathbf{Z})$ is also $2K_Z$. Since the rank of \mathbf{A} is k , the rank of $\mathbf{A}\tilde{\mathbf{S}}'''(\mathbf{Z})$ is $2K_Z$, which implies that $\mathbf{v}''_{s,t}(\mathbf{Z}) \in \mathbb{R}^{2K_Z}$ is a zero vector. Therefore, we have $\mathbf{v}'_{s,t}(\mathbf{Z}) = \mathbf{0}, \forall s \neq t$, and we have demonstrated that \mathbf{v} is a component-wise function.

Step II. Based on **Step I**, we demonstrate that \mathbf{A} is a block permutation matrix. Without loss of generality, we assume that the permutation in \mathbf{v} is Identity, where $\mathbf{v}(\mathbf{Z}) = [v_1(Z_1), \dots, v_{K_Z}(Z_{K_Z})]^T$ and each v_i is a nonlinear univariate scalar function. Since f and \tilde{f} are injective, \mathbf{v} is bijective and

$\mathbf{v}^{-1}(\mathbf{Z}) = [v_1^{-1}(Z_1), \dots, v_{K_Z}^{-1}(Z_{K_Z})]^T$. If we denote $\bar{\mathbf{S}}(\mathbf{v}(\mathbf{Z})) = \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z})) + \mathbf{A}^{-1}\mathbf{c}$, Eq. (24) can be reformulated as $\mathbf{S}(\mathbf{Z}) = \mathbf{A}\bar{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))$. We then apply \mathbf{v}^{-1} to \mathbf{Z} on both sides, which gives

$$\mathbf{S}(\mathbf{v}^{-1}(\mathbf{Z})) = \mathbf{A}\bar{\mathbf{S}}(\mathbf{Z}). \quad (29)$$

Let t be the index of an entry in \mathbf{S} that corresponds to the factorized part \mathbf{S}_f . For all $s \neq t$, we have:

$$0 = \frac{\partial \mathbf{S}(\mathbf{v}^{-1}(\mathbf{Z}))_t}{\partial Z_s} = \sum_{j=1}^k a_{tj} \frac{\partial \bar{\mathbf{S}}(\mathbf{Z})_j}{\partial Z_s}. \quad (30)$$

Since the entries of $\tilde{\mathbf{S}}$ are linearly independent, a_{tj} is zero for any j such that $\frac{\partial \bar{\mathbf{S}}(\mathbf{Z})_j}{\partial Z_s} \neq 0$. This includes the entries S_j that correspond to **1**) the factorized part that does not depend on Z_t ; and **2**) the non-factorized part \mathbf{S}_{nf} . Therefore, when t is the index of an entry in the sufficient statistics \mathbf{S} that corresponds to factor i in the factorized part \mathbf{S}_f , i.e., $\mathbf{S}_{f,i}$, the only non-zero a_{tj} are the ones that map between $\mathbf{S}_{f,i}(Z_i)$ and $\bar{\mathbf{S}}_{f,i}(v_i(Z_i))$. Therefore, we can construct an invertible submatrix \mathbf{A}'_i with all non-zero elements a_{tj} for all t that corresponds to factor i , such that

$$\mathbf{S}_{f,i}(Z_i) = \mathbf{A}'_i \bar{\mathbf{S}}_{f,i}(v_i(Z_i)) = \mathbf{A}'_i \tilde{\mathbf{S}}_{f,i}(v_i(Z_i)) + \mathbf{c}_i, \quad i = 1, \dots, K_Z, \quad (31)$$

where \mathbf{c}_i denotes the corresponding elements of \mathbf{c} . Eq. (31) means that for each $i = 1, \dots, K_Z$, the matrix block \mathbf{A}'_i of \mathbf{A} affinely transforms the i -specific sufficient statistics vector $\mathbf{S}_{f,i}(Z_i)$ into $\tilde{\mathbf{S}}_{f,i}(v_i(Z_i))$. In addition, there is also an additional block \mathbf{A}' that affinely transforms $\mathbf{S}_{nf}(\mathbf{Z})$ in into $\mathbf{S}_{nf}(\mathbf{v}(\mathbf{Z}))$. This completes the proof that \mathbf{A} is a block permutation matrix.

A.4.3 PART III

Let $\tilde{Z}_i = v_i(Z_i) = \tilde{f}^\dagger(\mathbf{X})_i$ be the i th inferred latent variable. Assume again that the permutation in \mathbf{v} is Identity. In this part, we prove that if Assumption 2 holds, each inferred latent variable \tilde{Z}_i is the bijective transformation of the true latent variable. The proof is as follows.

Proof. Plugging \tilde{Z}_i into Eq. (31), we have:

$$\mathbf{S}_{f,i}(Z_i) = \mathbf{A}'_i \bar{\mathbf{S}}_{f,i}(\tilde{Z}_i). \quad (32)$$

According to Assumption 2, there exists one dimension of $\mathbf{S}_{f,i}$, i.e., j , such that $S_{f,ij}$ is bijective. This implies that $\mathbf{S}_{f,i}$ is injective, and therefore it has a left-inverse $\mathbf{S}_{f,i}^\dagger$. we apply $\mathbf{S}_{f,i}^\dagger$ to both sides of Eq. (32), which gives:

$$Z_i = \mathbf{S}_{f,i}^\dagger \mathbf{A}'_i \bar{\mathbf{S}}_{f,i}(\tilde{Z}_i). \quad (33)$$

Since \mathbf{A}'_i is a block of an invertible block permutation matrix, \mathbf{A}_i is also an invertible matrix, and therefore \mathbf{A}'_i is a bijective mapping. In addition, since $\tilde{\mathbf{S}}_{f,i}$ is injective, $\bar{\mathbf{S}}_{f,i}$ is also injective, and therefore the composite map $\mathbf{S}_{f,i}^\dagger \mathbf{A}'_i \bar{\mathbf{S}}_{f,i} : \mathbb{R} \rightarrow \mathbb{R}$ that applies on \tilde{Z}_i is a bijective. This completes the proof that each inferred latent variable \tilde{Z}_i is the bijective transformation of the true latent variable in the case of no noise, where $\mathbf{Z} = f^\dagger(\mathbf{X})$ are the true latent variables. If noise ε exists, the posterior distribution of the latent variables can be identified up to an analogous bijective indeterminacy. \square

A.4.4 CONSISTENCY

Proof. If the family of the variational posterior $q_\phi(\mathbf{Z}|\mathbf{X}, Y, T)$ contains the true posterior $p_\theta(\mathbf{Z}|\mathbf{X}, Y, T)$, then by optimizing the loss of Eq. (10) (with the KL term replaced by the score matching loss defined in Eq. (11)) over its parameter ϕ , the score matching term will eventually vanish. Therefore, the ELBO term in Eq. (10) will be equal to the log-likelihood. Under this circumstance, CiVAE inherits all the properties of maximum likelihood estimation (MLE). Since the identifiability of CiVAE is guaranteed up to permutation and component-wise bijective transformation of the latent variables, the consistency property of MLE means that the model will converge to the true parameter θ^* up to such mild indeterminacy of the latent variables in the limit of infinite data. \square

A.5 PROOF OF THEOREM 3.2

Proof. Let \mathbf{C} be the true latent confounders and $\hat{\mathbf{C}}$ be the transformed confounders, where the transformation function g is bijective and differentiable *a.e.* Let g^{-1} denote its inverse. The ATE estimator that controls transformed confounders $\hat{\mathbf{C}}$ can be formulated as:

$$DEV(\hat{\mathbf{C}}) = \mathbb{E}_{p(\hat{\mathbf{C}})}[\mathbb{E}[Y|T = 1, \hat{\mathbf{C}} = \hat{\mathbf{c}}] - \mathbb{E}[Y|T = 0, \hat{\mathbf{C}} = \hat{\mathbf{c}}]]. \quad (34)$$

Specifically, for the continuous case where density functions exist, for each term, we have:

$$\mathbb{E}_{p(\hat{\mathbf{C}})}[\mathbb{E}[Y|T = t, \hat{\mathbf{C}} = \hat{\mathbf{c}}]] = \int f_{\hat{\mathbf{C}}}(\hat{\mathbf{c}}) \int y \cdot f_{Y|T, \hat{\mathbf{C}}}(y|t, \hat{\mathbf{c}}) dy d\hat{\mathbf{c}}. \quad (35)$$

For the marginal density $f_{\hat{\mathbf{C}}}(\hat{\mathbf{c}})$, the following equality holds:

$$f_{\hat{\mathbf{C}}}(\hat{\mathbf{c}}) = f_{\mathbf{C}}(g^{-1}(\hat{\mathbf{c}})) |J_{g^{-1}}(\hat{\mathbf{c}})| = f_{\mathbf{C}}(\mathbf{c}) |J_{g^{-1}}(\hat{\mathbf{c}})|. \quad (36)$$

As for the conditional density $f_{Y|T, \hat{\mathbf{C}}}(y|t, \hat{\mathbf{c}})$, since g is bijective, according to Eq. (13), we have:

$$f_{Y|T, \hat{\mathbf{C}}}(y|t, \hat{\mathbf{c}}) = f_{Y|T, \mathbf{C}}(y|t, \mathbf{c}). \quad (37)$$

Combining Eqs. (36) and (37), and given that $d\hat{\mathbf{c}} = |J_g(\mathbf{c})| d\mathbf{c}$, we have:

$$\begin{aligned} (35) &= \int f_{\mathbf{C}}(\mathbf{c}) |\mathbf{J}_{g^{-1}}(\hat{\mathbf{c}})| \int y \cdot f_{Y|T, \mathbf{C}}(y|t, \mathbf{c}) dy |\mathbf{J}_g(\mathbf{c})| d\mathbf{c} \\ &= |\mathbf{J}_{g^{-1}}(\hat{\mathbf{c}})| \cdot |\mathbf{J}_g(\mathbf{c})| \int f_{\mathbf{C}}(\mathbf{c}) \int y \cdot f_{Y|T, \mathbf{C}}(y|t, \mathbf{c}) dy d\mathbf{c} \\ &\stackrel{(a)}{=} \int f_{\mathbf{C}}(\mathbf{c}) \int y \cdot f_{Y|T, \mathbf{C}}(y|t, \mathbf{c}) dy d\mathbf{c} \\ &= \mathbb{E}_{p(\mathbf{C})}[\mathbb{E}[Y|T = t, \mathbf{C} = \mathbf{c}]], \end{aligned} \quad (38)$$

where the term $|\mathbf{J}_{g^{-1}}(\hat{\mathbf{c}})| \cdot |\mathbf{J}_g(\mathbf{c})|$ vanishes in step (a) as the two factors have the product of one. Therefore, if we plug Eq. (38) into Eq. (34), it leads to the following equality:

$$\begin{aligned} DEV(\hat{\mathbf{C}}) &= \mathbb{E}_{p(\hat{\mathbf{C}})}[\mathbb{E}[Y|T = 1, \hat{\mathbf{C}} = \hat{\mathbf{c}}] - \mathbb{E}[Y|T = 0, \hat{\mathbf{C}} = \hat{\mathbf{c}}]] \\ &= \mathbb{E}_{p(\mathbf{C})}[\mathbb{E}[Y|T = 1, \mathbf{C} = \mathbf{c}] - \mathbb{E}[Y|T = 0, \mathbf{C} = \mathbf{c}]] \\ &= DEV(\mathbf{C}) = ATE, \end{aligned} \quad (39)$$

where the last step is due to Eq. (2) in Definition 2, which completes our proof that controlling bijectively transformed confounders provides an unbiased estimation of ATE. \square