# Causal Inference

Registration Number:  2108737

**Abstract:**

Many key issues of causal inference, on the other hand, need assessing the causal consequences of certain treatments or exposures at the community level using data obtained from a sample of persons in that community. The need to assess the effectiveness of interventions when implemented on a large scale in institutions has sparked interest in methods for estimating the causal effects of community-based interventions. This is in contrast to individual subject effectiveness of these interventions, which can often be assessed in a conventional randomized controlled trial. Estimating treatment effects at the individual level (ITE) using observational data is a difficult and crucial field of causal machine learning that is being studied in a wide range of mission-critical applications. This work provides an information-theoretic technique to determining more trustworthy ITE representations. Use the principle of the Information Bottleneck (IB). This addresses the trade-off between brevity and predictive power in presentations (A., 2019).
The use of an expanded visual model for causal information bottlenecks increases the independence of learnt expressions and treatment kinds. From the standpoint of understanding ITE in the context of semi-supervised learning, we also provide an additional kind of regularization to ensure more trustworthy representation. Experiments indicate that our model obtains the most recent findings and incorporates the real dataset's uncertainty information for more trustworthy predicting performance (Hill, 2021).

**Introduction:**

One of the core issues of machine learning is estimating the therapeutic impact at the individual level (ITE), which infers the causal link between treatment and results from observational data, and it has applications in areas as diverse as medical politics and online advertising. This is a must for your application. There were only a few persons. The result of each treatment must be precisely anticipated for each individual data point in order to accurately deduce the causal link between treatment and outcome. However, in most (if not all) observational studies, each data point gets just one treatment option from a list of several, and the effects of alternative treatments are not included in the data (Louizos, 2018).

As a result, this ITE problem is linked to counterfactual questions like "Will a patient's illness be different if he or she is given a different prescription?" To forecast unobserved counterfactual outcomes, a basic method is to employ a model trained with just observed factual data. The implicit assumption in such reasoning is that the treatment and control groups' input distributions (for the two alternatives) must be identical. However, unless data is obtained in a rigorously randomized controlled experiment, it has been established that treatment decisions might be skewed, resulting in a covariate shift in the ITE problem (Rosenbaum, 2019).

Previous research has focused around representational learning techniques to tackle this challenge. The feature extractor is taught in this circuit, and any distance between processing and control becomes less in the learned representation space than in the covariance space. initial. Almost all of the previous work for ITE based on this representation learning technique can be viewed as a means to identify a balanced and maximum representation that can be employed in the project at the same time. forecast both real and hypothetical consequences. By boosting generalization, IB directors' models are empirically well
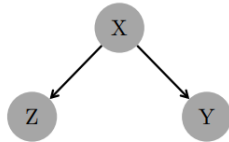
calibrated, resulting in high uncertainty in their forecasts. In a similar vein, it is reasonable to predict that the proposed CIB will automatically inherit the IB's benefits. In this regard, we empirically establish that the CIB outperforms the baseline when the models are permitted to state "I don't know" on situations other than the training data, and since then they haven't been sure. Since most tasks regard counterfactual inference to be irreversible and have substantial effects, this is one of the most critical criteria for imparting accuracy uncertainty to their predictions. serious outcomes (Rubin, 2019).

**Estimates for Individualized Treatment** (ITE) The set of potential variables for X and the range of different outcomes for Y are shown below. For the purpose of simplicity, we focus on a binary processing situation in which you pick one of the handling/control alternatives. We choose 0 to represent the controls and 1 to represent as treatment for symbolic reasons, therefore the treatment space is $T = \{0,1\}$. Let $x_i \in \chi$ be the covariate defining the features of patient I for example. $t_i = 1$ if the patient elect's therapy; otherwise, $t_i = 0$. After medical therapy is chosen, the patient's real progression factor (blood pressure, blood glucose) becomes $y_i^F \in y$ (Schaar, 2020).
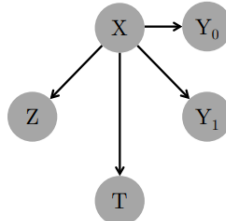
For n patients, an actual data set may be represented as $\{(x_i, t_i, y_i^F)\}_i^n = 1$. Antireal data may be written as $\{(x_i, 1 - t_i, y_i^{CF})\}_i^n = 1$, but the observed data does not have a $y_i^{CF}$I value. The conditional anticipated difference between the treatment result $y_i^{(1)}$ and the control outcome $y_i^{(0)}$ with the covariant $T(x_i) = \mathbb{E}[y_i^{(1)} - y_i^{(0)} \mid x_i]$, is the individual treatment effect (ITE) that we wished to estimate.

We assume the following two criteria of RubinNeyman's causality throughout the study, much as some other studies on causal inference have done so to show counterfactual predictions (LaLonde, 2019).
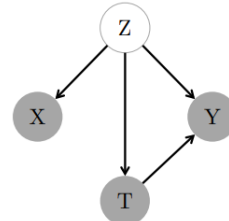
$$(\text{Overlap}) \ \mathbb{P}[t|x] \in (0,1), \forall x \in \mathcal{X}, \forall t \in \mathcal{T},$$



(a) A graphical representation of a traditional information bottleneck (b) Diagrammatic representation of a causal information bottleneck (c) CEVAE graphical model

$$(\text{Conditional ignorability}) \ Y \perp\!\!\!\perp T | X.$$

In order for the overlap hypothesis to work, all patients must have a larger than zero chance of getting all therapies. This assumption is required to guarantee that each patient has both treatment and control choices, allowing ITE to be accurately assessed. The assumption of conditional neglect ability states that

variables provide information regarding the overall impact of treatment decisions and outcomes. As a result, it's often called the case without a confounding hypothesis. For make the ITE estimate determinable, this no-confusion assumption is required (Vaughan, 2021).

**The principle of the information bottleneck (IB)** This information bottleneck idea was suggested for the first-time in. The following equation, which comprises two mutually informative terms I(,): one between input X and the Z representation and the other between the Z representation and the Y output, represents the trade-off between maximum expressiveness and maximum compression in information bottleneck theory:

$$\text{maximize}\, \mathcal{I}(Y; Z) - \beta\, \mathcal{I}(X; Z).$$

The traditional information bottleneck uses the graphical model in Figure 1(a) to calculate p(z|x) by sampling from the random encoder's input. The bottleneck principle of information is recognized to be difficult to implement in general because an exact assessment of mutual information is indissoluble. Recent work on variable approximation and adversarial learning, for example, has proposed efficient techniques for approximating mutual information and confirmed that information congestion is an efficient and robust regularization method for observed noise, resulting in better generalization performance and resilience against inversion attacks, among other things (Schaar, 2020).

**Bottleneck in Causal Information (CIB)**

It's worth noting that, even in the conditional ignoring scenario described above, the covariate X might include useless information for both Y and T in many real-world situations. Overfitting is an issue caused by noise in covariance. Because it is difficult to obtain significant volumes of training data in major application areas of causal inference, such as health care and political science, and because these are frequently crucial problems that cannot be solved. We need a method to eliminate the noise in the covariance that has been inverted. as well as a dependable depiction of the forecasts (Hill, 2021).

Finally, we offer an information bottleneck paradigm for causal inference in this section. Figure 1(b) adds variables for treatment and actual/counterfactual outcomes to the graphical model of Figure 1(a). Unlike graphical models that capture data creation with the concealed jammer Z (CEVAE in Fig. 1(c)), the CIB's graphical model increases the bottleneck directly. to find an appropriate representation of the discrimination for the causal inference problem establish (LaLonde, 2019).

- (Maximal Expressiveness) representation Z should have high mutual information with treatment outcome Y0, Y1.
- (Maximal Compressiveness) representation Z should have low mutual information with covariate X.

With a configurable Lagrange multiplier, these concepts may be stated as the following optimization problem:

$$\text{maximize}\, \mathcal{I}(Y_0; Z) + \mathcal{I}(Y_1; Z) - \beta\mathcal{I}(X; Z).$$

The following variation approximation may be used to determine the first term of goal (1), the maximum expression of Z given the reference result Y0:

$$\mathcal{I}(Y_0; Z) = \int\int p(y_0, z) \log \frac{p(y_0, z)}{p(y_0)p(z)} dy_0 dz = \int\int p(y_0, z) \log \frac{p(y_0|z)}{p(y_0)} dy_0 dz$$

$$\geq \int\int p(y_0, z) \log \frac{q_\theta(y_0|z)}{p(y_0)} dy_0 dz \tag{2}$$

$$= \int\int p(y_0, z) \log q_\theta(y_0|z) dy_0 dz + \mathcal{H}(y_0) \tag{3}$$

where qθ(y|z) is a variable approximation of the conditional distribution p(y|z) and H(y0) is the entropy of the random variable y0. Inequality (2) comes from the non-negative property of the Kullback Leibler (KL) divergence (Rosenbaum, 2019):

$$D_{KL}(p(y_0|z)\|q_\theta(y_0|z)) \geq 0 \implies \int p(y_0|z) \log p(y_0|z) dy_0 \geq \int p(y_0|z) \log q_\theta(y_0|z) dy_0, \forall z.$$

(3) may now be determined using the random encoder p(z|x) as follows:

$$\int\int p(y_0, z) \log q_\theta(y_0|z) dy_0 dz = \int\int p(y_0|x)p_\phi(z|x)p(x) \log q_\theta(y_0|z) dy_0 dz dx. \tag{4}$$

This follows from the fact that $p(y_0, z) = \int p(y_0, z, x) dx = \int p(y_0, z|x)p(x) dx = \int p(y_0|x)p_\phi(z|x)p(x) dx$ under the graph structure Figure 1(b).

Because of the nature of observation data, extra issues occur in it (4); we right edge P (Y0 | X) On the complete Population P (x), but we only know Y0 under conditions Population P (x | T = 0). To answer this question, we may use the following formula: (4) under overlapping hypotheses p (t = 0 | x) (0, 1):

$$\int\int p(y_0|x)p_\phi(z|x)p(x) \log q_\theta(y_0|z) dy_0 dz dx$$

$$= \int\int p(y_0|x)p_\phi(z|x) \frac{p(x)}{p(x|t=0)} p(x|t=0) \log q_\theta(y_0|z) dy_0 dz dx$$

$$= \int\int p(y_0|x)p_\phi(z|x) \frac{p(t=0)}{p(t=0|x)} p(x|t=0) \log q_\theta(y_0|z) dy_0 dz dx.$$

We now present the score classifier sv, which uses data to assess the likelihood of treatment and define the loss:

$$\mathcal{L}_0(\phi, \theta) := \int\int p(y_0|x)p_\phi(z|x) \frac{p(t=0)}{s_\nu(t=0|x)} p(x|t=0) \log q_\theta(y_0|z) dy_0 dz dx.$$

In a similar fashion, we may obtain the variational bound for the mutual information between Z and Y1. The maximal compressiveness, on the other hand, may be determined as a continuous extension of variation approximation for conventional IB [3]: I(X;Z)
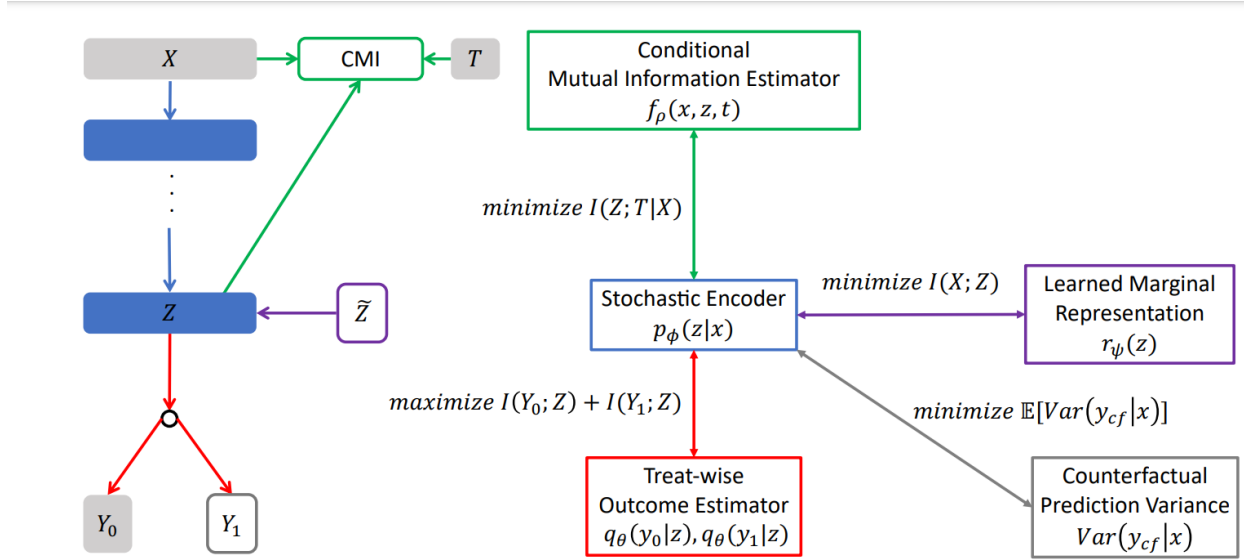
Figure 2: Overall architecture of Causal Information Bottleneck (CIB).

*Figure 1 : Causal Information Bottleneck's overall architecture (CIB).*

$$\int \int p(x,z) \log p_\phi(z|x)dxdz - \int p(z) \log r_\psi(z)dz =: \mathcal{L}_C(\phi, \psi)$$ where $r_\psi(z)$ is a variational approximation of $p(z)$.

However, the graph model in Figure 1(b) implies that the covariance X is given a conditional independence presentation of Z with treatment T. This assumption cannot be met using the random encoder p(z|x) if there are no further limitations. As a result, we offer a new approach known as mutually informative instruction decoupling, which will be discussed in the next section. Furthermore, semi-supervised learning may be viewed as ITE estimation, in which we were asked to estimate unlabeled counterfactual outcomes. We use predictive variance regularization, a prominent semi-supervised learning approach, to counterfactual prediction based on this interpretation (Rubin, 2019).

**Conditioning for Counterfactual Predictive Variance (CPVR)** Variance reduction is added to the predicted distribution of unlabeled data as a regular term in predictive variance conditioning, a popular semi-supervised learning approach. This technique contains the notion that an unlabeled sample in the representation space is close to a labeled sample in the conventional Bayesian framework. In this study, we incorporate an inductive bias to ensure that the stochastic encoder's true back predictions are consistent with each other.

$$\mathcal{L}_V(\phi, \theta) := - \mathbb{E}_{p(x)}\left[\mathrm{Var}_{q_\theta(y|x)}[y^{CF}|x]\right].$$

In that both give approximate goals for counterfactual data, regular predictive variance is analogous to the closest neighbor approach widely utilized in genuine backward prediction. Predictive variance regularization, on the other hand, employs learnt network predictions, whereas closest-based approaches send real data to unlabeled counterfactuals. When there is adequate reference data, the

closest neighbor approach can be a decent approximation, but if there is insufficient data to compute the neighbors, correcting for the prediction variance can offer superior inductive bias (Schaar, 2020).

Overall, CIB's ultimate goal is characterized as (with graphical description in Figure 2)

$$\underset{\theta,\phi,\psi}{\text{maximize}} \; \mathcal{L}(\theta, \phi, \psi) := \mathcal{L}_0(\phi, \theta) + \mathcal{L}_1(\phi, \theta) - \beta \cdot \mathcal{L}_C(\phi, \psi) + \lambda_M \cdot \mathcal{L}_M(\phi, \rho) + \lambda_V \cdot \mathcal{L}_V(\phi)$$

**Generative models for causal inference**

Causative Inference Generative Model For causal inference, the CIB's random encoder p(z|x) may be thought of as a composite model that samples representations that meet the information congestion principle. fruit. The supervised learning variant of VAE is the bottleneck on variable information. Although the network architecture and loss function are similar, they are interpreted differently: the VIB with the random encoder approximates the marginal representation p(z) and the conditional distribution p(y | z), whereas the VAE with the sampler and pre-generator approximates the posterior distribution. p(z|x). This discrepancy is due to the z provisions. Our model inherits this trait and has a comparable relationship to CEVAE since it is a natural extension of VIB for causal inference.

Parallel work recently exploited the information bottleneck concept to infer cause and effect. CEIB's graph structure presupposes that the representation Z has a direct impact on the outcome Y. The CIB, on the other hand, naturally extends the standard information bottleneck principle and inherits its benefits. The architectural decisions and desired functions are affected by the differences between the CIB and the CEIB (Shalit, 2019).

Table 1: Comparisons of counterfactual errors: 10 repeat/in-sample case

| Dataset | IHDP($\sqrt{\epsilon_{PEHE}}$) | Jobs($\mathcal{R}_{pol}$) | Twins(AUC) |
|---|---|---|---|
| TARNET | $0.729 \pm 0.088$ | $0.228 \pm 0.004$ | $0.849 \pm 0.002$ |
| CFR-M | $0.663 \pm 0.068$ | $0.213 \pm 0.006$ | $0.852 \pm 0.001$ |
| CFR-W | $0.649 \pm 0.089$ | $0.225 \pm 0.004$ | $0.850 \pm 0.002$ |
| CEVAE | $(2.7 \pm 0.1)$ | $0.15 \pm 0.0$ | not reported |
| SITE | $0.604 \pm 0.093$ | $0.224 \pm 0.004$ | $0.862 \pm 0.002$ |
| CIB | $0.663 \pm 0.193$ | $0.256 \pm 0.006$ | $0.870 \pm 0.002$ |

## Table 2: Comparisons of counterfactual errors: 10 repeat/out-sample case

| Dataset | IHDP($\sqrt{\epsilon_{PEHE}}$) | Jobs($\mathcal{R}_{pol}$) | Twins(AUC) |
|---------|--------------------------------|---------------------------|------------|
| TARNET  | $1.342 \pm 0.597$ | $0.234 \pm 0.012$ | $0.840 \pm 0.006$ |
| CFR-M   | $1.202 \pm 0.550$ | $0.231 \pm 0.009$ | $0.840 \pm 0.006$ |
| CFR-W   | $1.152 \pm 0.527$ | $0.225 \pm 0.010$ | $0.842 \pm 0.005$ |
| CEVAE   | $(2.6 \pm 0.1)$ | $0.26 \pm 0.0$ | not reported |
| SITE    | $\mathbf{0.656 \pm 0.108}$ | $\mathbf{0.219 \pm 0.009}$ | $\mathbf{0.853 \pm 0.006}$ |
| CIB     | $\mathbf{0.613 \pm 0.118}$ | $\mathbf{0.211 \pm 0.017}$ | $\mathbf{0.861 \pm 0.005}$ |

*Bold font indicates that the mean belongs to 95% confidence interval of the best performing model on each dataset.

**Baseline** Our technique is compared to traditional regression, closest neighbor, and representation learning. Ordinary least squares with the treatment as a feature (OLS/LR 1), and OLS with independent regressors for each treatment are examples of traditional regression methods (OLS2). HilbertSchmidt Independent Neighbor Match based on Closest Criteria (HSICNNM), Trend Score Match with Logistic Regression (PSM), and nearest neighbor approaches are all based on nearest neighbors (kNN). Balanced Neural Networks (BNN), Treatment Agnostic Representation Networks (TARNETs), Anti-Reality Regression with MMD/Wasserstein Metric (CFRM/W), Common Adversarial Networks for ITE Inference (GANITE), a variant autoencoder with causal effect (CEVAE), and locally similar preserved individual treatment effect are examples of representation learning methods (SITE) (Vaughan, 2021).

The MMD/Wasserstein index (SITEM/W) is also used to report SITE findings. Tables 1 and 2 findings, omitting our model, are shown. Because their results are based on somewhat different parameters: 1000 runs for the IHDP dataset and 100 runs for the Employment (GANITE) dataset, we have emphasized the CEVAE and GANITE outcomes (Rubin, 2019).

Implementation As in TARNET, we employ a common encoder for data processing/control and independent regressors. Our stochastic encoder predicts the mean and standard deviation of specific conditional Gaussian covariates using the Gaussian re-algebra method provided in VIB. The conditional mutual information estimator and the encoder both employ the same number of hidden classes. For the regressors, we employ a single-layer network. The CIB's experimental performance was not improved by the score classifier; therefore, it was omitted. The appendix contains more information regarding the deployment (Vaughan, 2021).
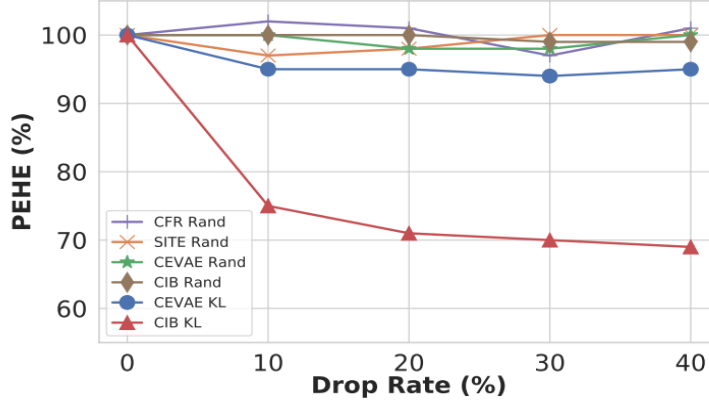
**Results** Tables 1 and 2 represent the means and standard errors of 10 out of three perceptions/separations datasets. CIB obtained the best results for the content error on the Twins dataset and the above sample error IHDP and Twins datasets. On the Employment dataset, CIB shows comparable performance to other top models. The tables here show comparisons with key baselines only, and the full comparisons are provided in the appendix due to space constraints.

As a consequence, the results reveal that representation-based learning approaches outperform traditional regression and closest neighbor-based learning methods. This is likewise the case with our CIB technique. Other representation equalization approaches, such as BNN, CFR, and SITE, take group equilibrium representations into account, but they don't define them to equalize (LaLonde, 2019).

Table 3: Results of CIB without MIGDR or CPVR

| Dataset | IHDP($\sqrt{\epsilon_{PEHE}}$) | | Jobs($\mathcal{R}_{pol}$) | | Twins(AUC) | |
| | In-sample | Out-sample | In-sample | Out-sample | In-sample | Out-sample |
| --- | --- | --- | --- | --- | --- | --- |
| CIB | $0.663 \pm 0.193$ | $0.613 \pm 0.118$ | $0.256 \pm 0.006$ | $0.211 \pm 0.017$ | $0.870 \pm 0.002$ | $0.861 \pm 0.005$ |
| w/o MIGDR | $0.664 \pm 0.191$ | $0.614 \pm 0.118$ | $0.246 \pm 0.004$ | $0.230 \pm 0.013$ | $0.864 \pm 0.001$ | $0.858 \pm 0.006$ |
| w/o CPVR | $0.686 \pm 0.186$ | $0.649 \pm 0.122$ | $0.245 \pm 0.013$ | $0.230 \pm 0.013$ | $0.865 \pm 0.001$ | $0.858 \pm 0.006$ |
| No regularizer | $0.686 \pm 0.186$ | $0.649 \pm 0.122$ | $0.245 \pm 0.004$ | $0.230 \pm 0.013$ | $0.865 \pm 0.001$ | $0.858 \pm 0.006$ |

The CIB's MIGDR examines selecting a balanced representative for each case, whereas representation is chosen for a representative group. The CIB does not suffer from the computation of the Sinkhorn Divergence in CFRW or the optimal kernel selection in CFRM since it adds an extra network for the equilibrium representation. Furthermore, unlike SITE, CIB does not require particular instances to analyze local similarities (Rosenbaum, 2019).



Figure 3: Results on removing top $k\%$ "I don't know" samples

Resection study We have proposed two additional regulations above, MIGDR and CPVR in addition to the basic CIB framework. Here, we confirm that our complement regulation is critical to the experimental success of our model. For this, we compare our CIB with CIB without MIGDR, CIB without CPVR and CIB without additional regulation across three data sets. Table 3 summarizes the results. CPVR regularization is necessary for the results of the IHDP dataset, while MIGDR and CPVR are important for other data sets (Hill, 2021).

**Conclusion:**

By widening the information bottleneck, we've established a new CIB methodology for estimating ITE. We've offered two new elements in addition to the information bottleneck framework to investigate a more trustworthy representation. We demonstrate that the CIB technique outperforms contemporary models in detecting the extent to which a given covariate x is a sample out of distribution (OOD) for the entire set p. (x). This CIB characteristic is critical for crucial causal inference applications.

**The whole exam results**

The mean is inside the 95 percent confidence interval of the highest performing model per data set when the text is bold.

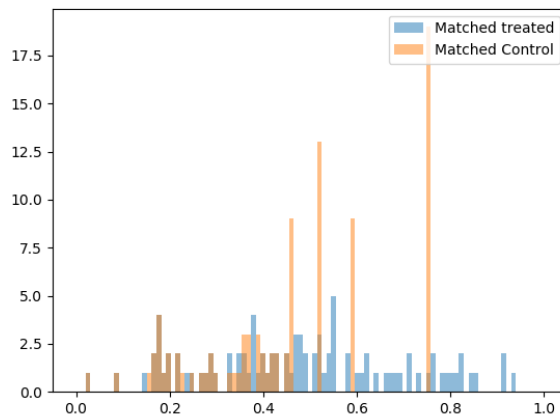Table 4: Mean/Std Err for counterfactual error 10 repeat/in-sample

| Dataset | IHDP($\sqrt{\epsilon_{PEHE}}$) | Jobs($\mathcal{R}_{pol}$) | Twins(AUC) |
|---|---|---|---|
| OSL/LR1 | $10.761 \pm 4.350$ | $0.310 \pm 0.017$ | $0.660 \pm 0.005$ |
| OSL/LR2 | $10.280 \pm 3.794$ | $0.228 \pm 0.012$ | $0.660 \pm 0.004$ |
| HSIC-NNM | $2.439 \pm 0.445$ | $0.291 \pm 0.019$ | $0.762 \pm 0.011$ |
| PSM | $7.188 \pm 2.679$ | $0.292 \pm 0.019$ | $0.500 \pm 0.003$ |
| $k$-NN | $4.432 \pm 2.345$ | $0.230 \pm 0.016$ | $0.609 \pm 0.010$ |
| BNN | $3.827 \pm 2.044$ | $0.232 \pm 0.008$ | $0.690 \pm 0.008$ |
| TARNET | $\mathbf{0.729 \pm 0.088}$ | $0.228 \pm 0.004$ | $0.849 \pm 0.002$ |
| CFR-M | $\mathbf{0.663 \pm 0.068}$ | $0.213 \pm 0.006$ | $0.852 \pm 0.001$ |
| CFR-W | $\mathbf{0.649 \pm 0.089}$ | $0.225 \pm 0.004$ | $0.850 \pm 0.002$ |
| GANITE | $(1.9 \pm 0.4)$ | $(0.13 \pm 0.01)$ | not reported |
| CEVAE | $(2.7 \pm 0.1)$ | $\mathbf{0.15 \pm 0.0}$ | not reported |
| SITE-M | $1.162 \pm 0.118$ | $0.194 \pm 0.015$ | $0.710 \pm 0.003$ |
| SITE-W | $0.993 \pm 0.112$ | $0.190 \pm 0.015$ | $0.849 \pm 0.003$ |
| SITE | $\mathbf{0.604 \pm 0.093}$ | $0.224 \pm 0.004$ | $0.862 \pm 0.002$ |
| CIB | $\mathbf{0.663 \pm 0.193}$ | $0.256 \pm 0.006$ | $\mathbf{0.870 \pm 0.002}$ |

Table 5: Mean/Std Err for counterfactual error 10 repeat/out-sample

| Dataset | IHDP($\sqrt{\epsilon_{PEHE}}$) | Jobs($\mathcal{R}_{pol}$) | Twins(AUC) |
|---|---|---|---|
| OSL/LR1 | $7.354 \pm 2.914$ | $0.279 \pm 0.067$ | $0.500 \pm 0.028$ |
| OSL/LR2 | $5.245 \pm 0.986$ | $0.733 \pm 0.103$ | $0.500 \pm 0.016$ |
| HSIC-NNM | $2.401 \pm 0.367$ | $0.311 \pm 0.069$ | $0.501 \pm 0.017$ |
| PSM | $7.290 \pm 3.389$ | $0.307 \pm 0.053$ | $0.506 \pm 0.011$ |
| $k$-NN | $4.303 \pm 2.077$ | $0.262 \pm 0.038$ | $0.492 \pm 0.012$ |
| BNN | $4.874 \pm 2.850$ | $\mathbf{0.240 \pm 0.012}$ | $0.676 \pm 0.008$ |
| TARNET | $1.342 \pm 0.597$ | $\mathbf{0.234 \pm 0.012}$ | $0.840 \pm 0.006$ |
| CFR-M | $1.202 \pm 0.550$ | $\mathbf{0.231 \pm 0.009}$ | $0.840 \pm 0.006$ |
| CFR-W | $1.152 \pm 0.527$ | $\mathbf{0.225 \pm 0.010}$ | $0.842 \pm 0.005$ |
| GANITE | $(2.4 \pm 0.4)$ | $(0.14 \pm 0.01)$ | not reported |
| CEVAE | $(2.6 \pm 0.1)$ | $0.26 \pm 0.0$ | not reported |
| SITE-M | $1.242 \pm 0.153$ | $\mathbf{0.218 \pm 0.010}$ | $0.705 \pm 0.006$ |
| SITE-W | $1.459 \pm 0.481$ | $\mathbf{0.232 \pm 0.011}$ | $0.762 \pm 0.007$ |
| SITE | $\mathbf{0.656 \pm 0.108}$ | $\mathbf{0.219 \pm 0.009}$ | $\mathbf{0.853 \pm 0.006}$ |
| CIB | $\mathbf{0.613 \pm 0.118}$ | $\mathbf{0.211 \pm 0.017}$ | $\mathbf{0.861 \pm 0.005}$ |

## DATASET

We use 3 Realworld data sets to evaluate the existing methodological methodologies for the proposed CIB assessment. An arbitrary RealWorld RealWorld test experience led to Covaries of a Health Program and Infant Development Program (IHDP). IHDP data sets are purposely tied to the selection deviation and utilized as reference points for ITE study by omitting some RCT data processing. There are 747 covariates in 25 dimensions of data. With the recommended 63/27/10 ratio, we separated shipping, authentication, and testing into 10 groups. The employment data collection is a composite observation based on observational research and the National Support Work Program (Shalit, 2019).

# References

A., F., 2019. *Uncertainty in the variational information bottleneck..* 3rd ed. London: arXiv preprint arXiv:1807.00906..

Hill, J. L., 2021. *Bayesian nonparametric modeling for causal inference..* 3rd ed. Bejing: Journal of Computational and Graphical Statistics.

LaLonde, R., 2019. *Evaluating the econometric evaluations of training programs with experimental data..* 2nd ed. US: American Economic Review, 76(4):604–20..

Louizos, C. S. U. M. J. M., 2018. *Causal effect inference with deep latent-variable models..* 4th ed. UK: In Advances in Neural Information Processing Systems.

Murphy, K., 2021. *Deep variational information bottleneck.* 4th ed. London: In ICLR..

Rosenbaum, P. R., 2019. *The central role of the propensity score in observational studies for causal effects..* 3rd ed. UK: Biometrika.

Rubin, D. B., 2019. *Causal inference using potential outcomes..* 2nd ed. US: Journal of the American Statistical Association.

Schaar, v. d., 2020. *Bayesian inference of individualized treatment effects using multi-task gaussian processes..* 3rd ed. London: In Advances in Neural Information Processing Systems.

Shalit, U. J., 2019. *Estimating individual treatment effect: generalization bounds and algorithms..* 3rd ed. US: In Proceedings of the 34th International Conference on Machine Learning,.

Vaughan, J. W., 2021. *A theory of learning from different domains..* 5th ed. UK: Machine learning, 79(1-2):151–175..