# **CATE Estimation With Potential Outcome Imputation From Local Regression**

Ahmed Aloui \*1

Juncheng Dong \*1

Cat P. Le<sup>1</sup>

Vahid Tarokh<sup>1</sup>

<sup>1</sup> Department of Electrical and Computer Engineering, Duke University

## Abstract

One of the most significant challenges in Conditional Average Treatment Effect (CATE) estimation is the statistical discrepancy between distinct treatment groups. To address this issue, we propose a model-agnostic data augmentation method for CATE estimation. First, we derive regret bounds for general data augmentation methods suggesting that a small imputation error may be necessary for accurate CATE estimation. Inspired by this idea, we propose a contrastive learning approach that reliably imputes missing potential outcomes for a selected subset of individuals formed using a similarity measure. We augment the original dataset with these reliable imputations to reduce the discrepancy between different treatment groups while inducing minimal imputation error. The augmented dataset can subsequently be employed to train standard CATE estimation models. We provide both theoretical guarantees and extensive numerical studies demonstrating the effectiveness of our approach in improving the accuracy and robustness of numerous CATE estimation models.

# **1 INTRODUCTION**

One of the most significant challenges for Conditional Average Treatment Effect (CATE) estimation is the statistical disparity between distinct treatment groups [Goldsmith-Pinkham et al., 2022]. While Randomized Controlled Trials (RCTs) mitigate this issue [Rubin, 1974, Imbens and Rubin, 2015], they can be expensive, unethical, and sometimes unfeasible to conduct. Consequently, we are often constrained to relying on observational studies susceptible to the aforementioned issue due to various selection biases. To address this, we introduce a *model-agnostic data aug*-

*mentation method* based on potential outcome imputation, comprising two key steps. First, our approach identifies a subset of individuals whose counterfactual outcomes can be reliably imputed. Subsequently, it performs imputation for the counterfactual outcomes of these selected individuals, thereby augmenting the original dataset with the imputed values. Crucially, our method functions as a data preprocessing module that is independent of the chosen CATE estimation model.

Our key insight is that potential outcome imputation for individuals in the observational dataset can reduce the statistical discrepancy between treatment groups. In particular, our method is motivated by an observed trade-off between (i) the discrepancy across treatment groups and (ii) the imputation error of the missing counterfactual outcomes. To give a concrete example, consider the scenario with a binary treatment assignment. In this context, no individual can appear in both the control and treatment groups due to the inaccessibility of counterfactual outcomes [Holland, 1986]. Suppose that we randomly impute the missing counterfactual outcomes and subsequently integrate each individual along with their imputed outcomes back into the dataset. This data augmentation process ensures that individuals from both the control and treatment groups become identical, effectively eliminating all disparities. However, it becomes evident that any model trained on such a randomly augmented dataset would exhibit poor performance, primarily due to the substantial error introduced by the random imputation of the counterfactual outcomes. To tackle this challenge, we propose to identify a subset of individuals whose counterfactual outcomes can be reliably imputed and limit the imputation to this group. Such a risk-averse approach effectively control the imputation error so that the positive impact of disparity reduction will outweigh the negative impact of imputation error. We note that this is a general framework which can be instantiated with different choices of imputation methods. In Section 3, we formalize our insights by providing theoretical guarantees for this framework.

To instantiate this framework, this work considers a specific

<sup>\*</sup>Equal contribution.

realization — *imputation with local regression methods* such as local Gaussian Process (GP) using close neighbors of individuals to impute their missing counterfactual outcome. We choose to use the number of close neighbors as the criterion for imputation: our method only impute for individuals with sufficient number of neighbors. Specifically, we impute the counterfactual outcomes for these individuals by utilizing the factual outcomes of their proximate neighbors We explore two distinct methods for imputation: linear regression and GP. To comprehensively assess the efficacy of our method, we theoretically demonstrate that our approach asymptotically generates datasets whose probability densities converge to those of RCTs. In addition, we provide finite-sample generalization bounds for GP-based local regression in Section 4.2.

To further boost the performance of local regression methods, we propose to use contrastive learning to identify the close neighbors of individuals for potential outcome imputation. Through contrastive learning, we learn a representation space and a similarity measure, such that within this learned representation space, *close* individuals identified by the similarity measure exhibit *similar* potential outcomes. This smoothness property facilitates local regression methods as the outcomes are highly correlated locally. Consequently, *this ensures that the imputation can be achieved locally within the representation space* with simple models that require minimal tuning.

Our empirical results in Section 6 further demonstrate the efficacy of our method, showcasing consistent enhancements in the performance of state-of-the-art CATE estimation models, including TARNet, CFR-Wass, and CFR-MMD [Shalit et al., 2017], S-Learner and T-Learner integrated with neural networks, Bayesian Additive Regression Trees (BART) [Hill, 2011, Chipman et al., 2010, Hill et al., 2020] with X-Learner [Künzel et al., 2019], and Causal Forests (CF) [Athey and Imbens, 2016] with X-Learner.

**Related Works.** One of the fundamental tasks in causal inference is to estimate *Average Treatment Effects* (ATE) and *Conditional Average Treatment Effects* (CATE) [Neyman, 1923, Rubin, 2005]. Various methods have been proposed for ATE estimation, including Covariate Adjustment [Rubin, 1978], Propensity Scores [Rosenbaum and Rubin, 1983], Doubly Robust estimators [Funk et al., 2011], Inverse Probability Weighting [Hirano et al., 2003], and recently Reisznet [Chernozhukov et al., 2022]. While these methods are successful for ATE estimation, they are not directly applicable to CATE estimation.

On the other hand, recent advances in machine learning have led to new approaches for CATE estimation, such as decision trees [Athey and Imbens, 2016], Gaussian Processes [Alaa and Van Der Schaar, 2017], Multi-task deep learning ensemble [Jiang et al., 2023], Generative Modeling [Yoon et al., 2018], and representation learning with deep neural networks [Shalit et al., 2017, Johansson et al., 2016]. It is worth noting that alternative approaches for investigating causal relationships exist, such as do-calculus, proposed by Pearl [Pearl, 2009a,b]. Here, we adopt the Neyman-Rubin framework. At its core, the CATE estimation problem can be seen as a missing data problem [Rubin, 1974, Holland, 1986, Ding and Li, 2018] due to the unavailability of the counterfactual outcomes. In this context, we propose a new data augmentation approach for CATE estimation by imputing certain missing counterfactuals. Data augmentation, a well-established technique in machine learning, serves to enhance model performance and curb overfitting by artificially expanding the size of the training dataset [Van Dyk and Meng, 2001, Chawla et al., 2002, Han et al., 2005, Jiang et al., 2020, Chen et al., 2020a, Liu et al., 2020, Feng et al., 2021].

A crucial aspect of our methodology is the identification of similar individuals. There are various methods to achieve this goal, including propensity score matching [Rosenbaum and Rubin, 1983], and Mahalanobis distance matching [Imai et al., 2008]. Nonetheless, these methods pose significant challenges, particularly in scenarios with large sample sizes or high-dimensional data, where they suffer from the curse of dimensionality. Recently, Perfect Match [Schwab et al., 2018] is proposed to leverage importance sampling to generate replicas of individuals. It relies on propensity scores and other feature space metrics to balance the distribution between the treatment and control groups during the training process. In contrast, we utilize contrastive learning to construct a similarity metric within a representation space. Our method focuses on imputing missing counterfactual outcomes for a selected subset of individuals, without creating duplicates of the original data points. While the Perfect Match method is a universal CATE estimator, our method is a model-agnostic data augmentation method that serves as a data preprocessing step for other CATE estimation models.

# 2 THEORETICAL BACKGROUND

Let  $T \in \{0, 1\}$  be a binary treatment assignment,  $X \in \mathcal{X} \subset \mathbb{R}^d$  be the covariates (features), and  $Y \in \mathcal{Y} \subset \mathbb{R}$  be the factual (observed) outcome. For each  $j \in \{0, 1\}$ , we define  $Y_j$  as the *potential outcome* [Rubin, 1974], which represents the outcome that would have been observed if only the treatment T = j was administered. The random tuple (X, T, Y) jointly follows the *factual (observational) distribution* denoted by  $p_{\text{F}}(x, t, y)$ . Let  $D_{\text{F}} = \{(x_i, t_i, y_i)\}_{i=1}^n$  denote a dataset that consists of n observations independently sampled from  $p_{\text{F}}$  where n is the number of observations.

**Definition 2.1** (CATE). The Conditional Average Treatment Effect (CATE) is defined as:

$$\tau(x) = \mathbb{E}[Y_1 - Y_0 | X = x].$$
 (1)

CATE is identifiable under the assumptions of positivity, i.e.,

 $0 < p_{\rm F}(T = 1|X) < 1$ , and *conditional unconfoundedness*, i.e.,  $(Y_1, Y_0) \perp T|X$  [Robins, 1986, Imbens and Rubin, 2015]. Let  $\hat{\tau}(x) = h(x, 1) - h(x, 0)$  denote an estimator for CATE where h is a hypothesis  $h : \mathcal{X} \times \{0, 1\} \to \mathcal{Y}$ that estimates the underlying causal relationship f between (X, T) and Y.

**Definition 2.2** (PEHE). The Expected Precision in Estimating Heterogeneous Treatment Effect (PEHE) [Hill, 2011] is defined as:

$$\varepsilon_{\rm PEHE}(h) = \int_{\mathcal{X}} (\hat{\tau}(x) - \tau(x))^2 p_{\rm F}(x) dx \tag{2}$$

**Definition 2.3.** For a joint distribution p over (X, T, Y) and a hypothesis h the loss function is defined as:

$$\mathcal{L}_p(h) = \int (y - h(x, t))^2 p(x, t, y) \, dx \, dt \, dy,$$

*Remark* 2.4.  $\varepsilon_{\text{PEHE}}$  is widely-used as the performance metric for CATE estimation. However, directly estimating  $\varepsilon_{\text{PEHE}}$ from observational data  $D_{\text{F}}$  is a non-trivial task, as it requires knowledge of the counterfactual outcomes to compute the ground truth CATE values. This challenge underscores that models for CATE estimation need to be robust to overfitting the factual distribution. Our empirical results (in Section 6) indicate that our method mitigates the risk of overfitting for various CATE estimation models.

# **3 UNDERSTANDING DATA AUGMENTATION FOR CATE ESTIMATION**

In this section, we present a generalization bound for the performance of CATE estimation models *trained using the augmented dataset*. This theoretical result will motivate our proposed data augmentation algorithm. Given the observed factual dataset  $D_{\rm F}$  with *n* samples, a counterfactual data augmentation algorithm has two main components:

- Component I: identifying a subset  $\mathcal{R}_n \subset \mathcal{X} \times \{0, 1\}$ , where  $\mathcal{R}_n^t \subset \mathcal{X}$  for  $t \in \{0, 1\}$  is the projection for the treatment and control groups on which to perform data augmentation.
- Component II: imputing the missing potential outcomes for individuals in  $\mathcal{R}_n$  with an algorithm  $\tilde{f}_n$ :  $\mathcal{R}_n \to \mathcal{Y}$ .

The marginal distribution of (X,T) in the augmented dataset can be defined as follows:

$$p_{\text{AF}}(x,t) = \frac{1}{1+\beta}p_{\text{F}}(x,t) + \frac{\beta}{1+\beta}q(x,1-t),$$

where  $\frac{\beta}{1+\beta}\in[0,\frac{1}{2}]$  represents the ratio of the number of the select individuals for augmentation to the total number of

samples in the augmented dataset, and  $q = \frac{p_{\rm F}(x,1-t)}{\alpha} \mathbb{1}_{\mathcal{R}_n}$ , with  $\alpha$  as the normalizing constant, i.e.,  $\alpha = \int p_{\rm F}(x, 1-t) \mathbb{1}_{\mathcal{R}_n}(x, 1-t) dx dt$ . In other words, q is the factual distribution of the alternative treatment group with its probability mass normalized to the augmentation region  $\mathcal{R}_n$ . Intuitively, an effective data imputation method  $\tilde{f}_n(x,t)$ should approximate the true function f in the region  $\mathcal{R}_n$ , i.e.,  $\tilde{f}_n(x,t) \approx f(x,t)$  for  $x \in \mathcal{R}_n^t$ . Hence,  $p_{\rm AF}(y|x,t)$  can be defined as follows: it is equal to  $p_{\rm F}(y|x,t)$  when (x,t) is sampled from the factual distribution; for samples drawn from q(x, 1-t),  $p_{\rm AF}(y|x,t)$  is defined as a point mass function  $\delta(y = \tilde{f}_n(x,t))$ .

Let  $p_{\text{RCT}}(x, t, y)$  represent the distribution of (X, T, Y)when the observations are sampled from randomized controlled trials. To establish the generalization bound, we assume that there is a true potential outcome function f such that  $Y = f(X, T) + \eta$  with  $\eta$  verifying that  $\mathbb{E}[\eta] = 0$ .

**Proposition 3.1** (Generalization Bound). Let *h* be a hypothesis, its  $\varepsilon_{PEHE}$  is upper bounded as follows:

$$\varepsilon_{PEHE}(h) \leq 4 \cdot \left(\underbrace{\mathcal{L}_{p_{AF}}(h)}_{(I)} + 2\underbrace{V(p_{RCT}(X,T), p_{AF}(X,T))}_{(II)} + \underbrace{\frac{\beta}{1+\beta} \cdot b_{\mathcal{A}}(n)}_{(III)}\right)$$
(3)

where  $V(g_1, g_2) = \frac{1}{2} \int_{\mathcal{S}} |g_1(s) - g_2(s)| ds$  is the total variation distance between two distributions, and,

$$b_{\mathcal{A}}(n) = \mathbb{E}_{X, T \sim q} \left[ \|f(X, T) - \tilde{f}_n(X, T)\|^2 \right]$$

**Interpretation.** We first note that term (I) in Proposition 3.1 is essentially the training loss of a hypothesis h on the augmented dataset. Term (II) characterizes the statistical similarity between the individuals' features in the augmented dataset and those generated from an RCT. As there is no statistical disparity across treatment groups when (X, T)follows  $p_{\text{RCT}}$ , the closer  $p_{\text{AF}}$  is to  $p_{\text{RCT}}$  the less is the statistical disparity in the augmented dataset. Meanwhile, (III) characterizes the accuracy of the data augmentation method. Hence, this theorem provides a rigorous illustration of the trade-off between the statistical disparity across treatment groups and the imputation error. It underscores that a data augmentation method with a low potential outcome imputation error can improve CATE estimation. Also note that as  $\frac{\beta}{1+\beta}$  (i.e., the ratio of imputed data points to all the data points) increases, (III) increases while (II) decreases. This captures another trade-off between the precision of data imputation and the discrepancy across treatment groups. It is also essential to highlight that if the local regression module can achieve more accurate estimation with more samples (e.g., local Gaussian Process)  $b_{\mathcal{A}}(n)$  will converge to 0, as proved in Section 4.2.

# 4 POTENTIAL OUTCOME IMPUTATION FROM LOCAL REGRESSION

While Section 3 proposes a general framework relying on potential outcome imputation, we consider a specific instantiation — imputation from local regression methods with simple function classes such as linear regression and GP. We opt for these relatively straightforward function classes motivated by the following three principles:

- *Local Approximation*: Complex functions can be locally estimated with simple functions, e.g., continuous functions and complex distributions can be approximated by a linear function [Rudin, 1953] and Gaussian distributions [Tjøstheim et al., 2021], respectively.
- \* *Sample Efficiency*: If a class of simple functions can estimate the true target function locally, then a class with less complexity will require fewer close neighbors for good approximations.
- † Practicality: A simpler function class requires less hyper-parameter tuning, which remains one of the most significant challenges in causal inference applications.

We refer to these approaches as <u>Potential Outcome via Local</u> Regression (POLO).

## 4.1 ALGORITHM

**Overview.** POLO have two components. The first component is a classifier g(x, x', t). For example, g can be a threshold function based on the Euclidean distance in  $\mathcal{X}$ :  $g(x, x', t) = \mathbb{1}\{||x - x'|| \le \epsilon_t\}$  where  $\epsilon_t \ge 0$  is a prespecified threshold. Recall that  $D_F = \{x_i, y_i, t_i\}_{i=1}^n$  is the factual dataset. For a given individual x, g identifies x's close neighbors  $D_x$ , that is, individuals in  $D_F$  who are likely to exhibit similar outcomes when subjected to the same treatment t. The second component is a local regressor  $\psi(x, D_x)$ , which imputes the counterfactual outcome for x after being fitted to its close neighbors  $D_x$ .

For  $t \in \{0, 1\}$ , we use  $D^t \subset D_F$  to denote the factual observations in treatment group t, i.e,  $D^t = \{(x_i, t_i, y_i) \in D_F | t_i = t\}$ . Note that  $D_F = D^0 \cup D^1$ . For a given individual x in  $D_F$  within treatment group t whose counterfactual outcome (i.e., potential outcome under treatment 1-t) needs to be imputed, the classifier g first selects from  $D^{1-t}$  a subset of individuals<sup>1</sup> who are close neighbors to x denoted by  $D_x$ . Specifically,

$$D_x = \{i \in [n] : t_i = 1 - t, g(x, x_i, 1 - t) = 1\}.$$
 (4)

Here,  $D_x$  are individuals in treatment group 1 - t who are likely to have similar potential outcomes to x under

treatment 1 - t. Subsequently, the non-parametric regressor  $\psi$  utilizes the factual outcomes in  $D_x$  to estimate the counterfactual outcome of x:  $\hat{y}_x = \psi(x, D_x)$ . Finally, the imputed outcome of x is incorporated into the dataset, i.e.,  $D_A^{1-t} = D^{1-t} \cup \{(x, 1-t, \hat{y}_x)\}$ . This process is repeated for every individuals in the factual dataset  $D_F$ . The augmented dataset  $D_A = D_A^0 \cup D_A^1$  will be used as the training dataset for CATE estimation models.

However, as discussed in Section 1 and shown by Propsition 3.1, the minimal error of the counterfactual imputation plays a crucial role in the success of data augmentation. *To ensure the reliability of these imputations, we only perform imputations for individuals who possess a sufficient number of close neighbors.* In the worst case, no individuals will meet these criteria for imputation, resulting in *no augmentation* of the dataset. To this end, POLO is a risk-averse method that ensures *performance at least does not degrade*.

**Gaussian Process.** While there are many choices for the non-parametric local regressor  $\psi$ , we focus on GP in this work and next elaborate how a local GP  $\psi(x, D_x)$  imputes the potential outcome. GP is a non-parametric method [Seeger, 2004] that offers robust solutions to regression problems. It is fully characterized by a mean function  $m : \mathcal{X} \to \mathbb{R}$  and a kernel  $K : \mathcal{X} \times \mathcal{X} \to \mathbb{R}^+$  and it is denoted as  $\mathcal{GP}(m, K)$ . A GP is a random process  $\phi(\mathcal{X})$  indexed by a set  $\mathcal{X}$  such that any finite collection of these random variables follows a multivariate Gaussian distribution. Consider a finite index set of n elements  $\mathbf{x}_n = \{x_i\}_{i=1}^n$ , then the n-dimensional random variable  $\phi(\mathbf{x}_n) = [\phi(x_1), \phi(x_2), \dots, \phi(x_n)]$  follows a Gaussian distribution:

$$\phi(\mathbf{x}_n) \sim \mathcal{N}\big(m(\mathbf{x}_n), K(\mathbf{x}_n, \mathbf{x}_n)\big) \tag{5}$$

where  $m(\mathbf{x}_n) = [m(x_1), \dots, m(x_n)]$  is the mean and the  $K(\mathbf{x}_n, \mathbf{x}_n)$  is a  $n \times n$  covariance matrix whose element on the *i*-th row and *j*-th column is defined as  $K(\mathbf{x}_n, \mathbf{x}_n)_{ij} = K(x_i, x_j)$ 

Based on the principle of Local Approximation, if an individual x in the factual dataset received treatment t, it is assumed that the potential outcome of the individual x and those of its close neighbors (i.e., the individuals within  $D_x$ ) under treatment 1 - t follow a GP. Note that by construction of  $D_x$ , the potential outcome of  $D_x$  under treatment 1 - t is the observed factual outcome. Thus, after constructing  $D_x$ , , the counterfactual outcome for x is imputed as

$$\widehat{y}_x^{1-t} = \psi(x, D_x) = \mathbb{E}[y^{1-t} | x, \{y_i\}_{i \in D_x}].$$
(6)

Under the assumption of GP,  $\hat{y}_x^{1-t}$  has a closed-form solution. Let  $\sigma(i)$  denote the *i*-th smallest index in  $D_x$  and K denote the kernel (covariance function) of GP. Then

$$\widehat{y}_x^{1-t} = \mathbf{K}_x^\top \mathbf{K}_{xx} \mathbf{y},\tag{7}$$

<sup>&</sup>lt;sup>1</sup>The terms "individual" and "indices of individuals" are used interchangeably.

where

$$\mathbf{K}_x = [K(x, x_{\sigma(1)}), \dots, K(x, x_{\sigma(|D_x|)})],$$
  
$$\mathbf{y} = [y_{\sigma(1)}, \dots, y_{\sigma(|D_x|)}],$$

and  $\mathbf{K}_{xx}$  is a  $|D_x| \times |D_x|$  matrix whose element on the *i*-th row and *j*-column is  $K(x_{\sigma(i)}, x_{\sigma(j)})$ . Finally, we append the tuple  $(x, 1-t, \hat{y}_x^{1-t})$  into the factual dataset to augment the training data.

#### 4.2 THEORETICAL ANALYSIS

Here we study the theoretical properties of the proposed method. Specifically, we present two main theoretical results regarding the efficacy of POLO: *(i)* Our first result characterizes the asymptotic behavior, demonstrating that the distribution of the augmented dataset converges towards the distribution of randomized controlled trials (RCTs); *(ii)* Our second result establishes finite-sample regret guarantees for POLO with GP, establishing that its imputation error can be effectively controlled.

**Notation.** We use  $\mathcal{O}$  to denote the standard big-O notation for asymptotic behaviors and  $\tilde{\mathcal{O}}$  to denote the big-O notation ignoring all the log terms.  $|| \cdot ||_2$  denote the Euclidean norm. For any two values  $a, b \in \mathbb{R}$ , we let  $a \vee b = \max(a, b)$  and  $a \wedge b = \min(a, b)$ .

Let  $n_1$  and  $n_0$  denote the number of individuals in the treatment and control groups, respectively. We define  $u = \mathbb{P}(T = 1)$  as the probability of an individual being in the treatment group, and let  $\psi = \frac{u}{1-u}$ . Moreover, let  $X^t \stackrel{d}{=} (X|T = t)$  and  $\gamma = \mathbb{P}(\rho(X^1, X^0) \ge \epsilon) \in (0, 1)$  where  $\rho(\cdot, \cdot)$  denotes the distance metric between features (e.g. the contrastive learning distance) of the treatment and control groups, and  $\epsilon$  is a pre-defined threshold. POLO defines augmentation regions for the control and the treatment groups denoted as  $\mathcal{R}_n^0$  and  $\mathcal{R}_n^1$ , respectively. For  $t \in \{0, 1\}$ , we have that,

$$\mathcal{R}_n^{1-t} = \{ x_j | j \in [n], t_j = 1 - t, \ \exists i_1 < \ldots < i_k \in [n], \\ t_{i_k} = t, \rho(x_{i_k}, x) \le \epsilon \}$$

where k is a positive constant denoting the number of neighbors. We remark that for any given individual x, the likelihood of encountering neighboring data points is sufficiently high as the number of data points grows, which facilitates reliable imputation of its counterfactual outcome. This concept is formally captured in the following Proposition.

**Proposition 4.1.** Let  $j \in \{0, 1\}$  and  $\alpha_{n_j} = \mathbb{P}(X^j \in \mathcal{R}_n^j)$ , be the probability of finding k close neighbors for X in the alternative group. Then

$$\alpha_{n_j} \ge 1 - n_j^k \gamma^{n_j} \sum_{i=0}^{k-1} \frac{1}{i!} \left(\frac{1-\gamma}{\gamma}\right)^i.$$

Hence,  $1 - \alpha_{n_j} = \mathcal{O}(n_j^k \gamma^{n_j}).$ 

This implies that with a sufficient number of samples, the probability of not encountering data points in close proximity to any given point x becomes very small as the exponential decay  $\gamma^{n_j}$  for  $\gamma < 1$  dominates. Hence, positivity ensures that within the big data regime, we will encounter densely populated regions, enabling us to approximate counterfactual distributions locally. This facilitates the application of our methods. Next, we prove that *our data augmentation method converges to an RCT*.

**Proposition 4.2** (Convergence to RCT). Let  $p_{AF}^1$  and  $p_{AF}^0$  be the distributions of the treatment and control groups, respectively, after data augmentation. The following upper bound holds:

$$V(p_{AF}^{1}, p_{AF}^{0}) \leq \frac{1 - \alpha_{n_{0}}}{1 + z^{-1}\alpha_{n_{1}}} + \frac{z\alpha_{n_{0}}\left(1 - \alpha_{n_{1}}\right)}{1 + \alpha_{n_{0}}z} + \frac{|1 - \alpha_{n_{0}}\alpha_{n_{1}}|}{\left(1 + z^{-1}\alpha_{n_{1}}\right)\left(1 + \alpha_{n_{0}}z\right)},$$
(8)

as  $n_1$  and  $n_0$  converge to infinity, we have that  $\alpha_{n_1}$  and  $\alpha_{n_0}$  converge to 1 with the rates proved in Proposition 4.1. Hence, the right-hand side of the bound converges to 0.

Now, we establish the finite-sample guarantees for the GP local regressor. From a functional perspective and by Mercer's decomposition Seeger [2004], a GP can be considered as a distribution on a function class  $\mathcal{F} \subset \{f : \mathcal{X} \to \mathbb{R}\}$ , and  $\mathcal{F}$  is fully specified by GP's kernel  $K : \mathcal{X} \times \mathcal{X} \to \mathbb{R}_0^+$ .

**Assumption 4.3.** The potential outcome functions belong to this function space  $\mathcal{F}$ , i.e,

$$\{f(X, T=t): \mathcal{X} \to \mathbb{R} \mid t \in \{0, 1\}\} \subset \mathcal{F}.$$

*Remark* 4.4. This assumption is not unreasonable because, by choosing a radial basis function (RBF) kernel, the function class  $\mathcal{F}$  is assumed to contain all continuous functions which commonly include the potential outcomes functions.

**Definition 4.5** (Lipschitz Constant for GP Kernel). Assume that  $K : \mathcal{X} \times \mathcal{X} \to \mathbb{R}^+$  is the kernel of a Gaussian Process (GP). Its Lipschitz constant  $L_K$  is defined as:

$$L_K(\mathcal{X}) = \sup_{x, x' \in \mathcal{X}} ||\nabla_x K(x, x')||_2.$$
(9)

*Remark* 4.6. For well-known kernels, such as RBF,  $L_K$  is known and finite if  $\mathcal{X}$  is a bounded space. Moreover,  $L_K(\mathcal{X})$  is an increasing function of the input space  $\mathcal{X}$ , i.e., if  $\mathcal{X} \subset \mathcal{X}', L_K(\mathcal{X}) \leq L_K(\mathcal{X}')$ .

In this part, we assume that the data generation process is as follows,

$$Y = f(X, T) + \eta,$$

where  $\eta \sim \mathcal{N}(0, \sigma^2)$  and it is independent of (X, T). We also assume that  $\mathcal{X} \subset \mathbb{R}^d$  and the potential outcomes function f are bounded, and f is  $L_f$ -Lipschitz continuous.

Assume there is a dataset  $\{x_i, y_i\}_{i=1}^{\bar{n}_t}$  available with  $\bar{n}_t$  samples for the imputation of potential outcomes under treatment t. Let  $\sigma_{\bar{n}_t}(x) = K(x,x) - K(x,\mathbf{x}_{\bar{n}_t})(K(\mathbf{x}_{\bar{n}_t},\mathbf{x}_{\bar{n}_t}) + \sigma^2 \cdot I_{\bar{n}_t})^{-1}K(\mathbf{x}_{\bar{n}_t},x)$  be the posterior standard deviation of GP at x where

$$K(x, \mathbf{x}_{\bar{n}_t}) \in \mathbb{R}^{1 \times \bar{n}_t} = [K(x, x_1), \dots, K(x, x_{\bar{n}_t})],$$
  

$$K(\mathbf{x}_{\bar{n}_t}, x) \in \mathbb{R}^{\bar{n}_t \times 1} = [K(x, x_1), \dots, K(x, x_{\bar{n}_t})]^{\top},$$
  

$$K(\mathbf{x}_{\bar{n}_t}, \mathbf{x}_{\bar{n}_t}) \in \mathbb{R}^{\bar{n}_t \times \bar{n}_t}, K(\mathbf{x}_{\bar{n}_t}, \mathbf{x}_{\bar{n}_t})_{ij} = K(x_i, x_j).$$

Let  $\tilde{f}_{\bar{n}_t}(x,t)$  denote the GP-based imputation function given the dataset  $\{x_i, y_i\}_{i=1}^{\bar{n}_t} \subset D^t$ , i.e.,  $\tilde{f}_{\bar{n}_t}(x,t) = K(x, \mathbf{x}_{\bar{n}_t})(K(\mathbf{x}_{\bar{n}_t}, \mathbf{x}_{\bar{n}_t}) + \sigma^2 \cdot I_{\bar{n}_t})^{-1}\mathbf{y}_{\bar{n}_t}$  where  $\mathbf{y}_{\bar{n}_t} = [y_1, \ldots, y_{\bar{n}_t}]^\top$ . Note  $\tilde{f}_{\bar{n}_t}$  is a random function, varying with the observed dataset. The following result addresses its generalization error.

**Proposition 4.7.** For  $t \in \{0,1\}$ , let  $L_K^t = L_K(\mathcal{R}_n^{1-t})$ denote the Lipschitz constant of the kernel K in region  $\mathcal{R}_n^{1-t}$ and let  $U_K^t = \sup_{x,x' \in \mathcal{R}_n^{1-t}} K(x,x')$  denote the "width" of region  $\mathcal{R}_n^{1-t}$ . Then for  $t \in \{0,1\}$ , with probability at least  $1 - \delta$  where  $\delta \in (0,1)$ ,

$$\begin{split} \sup_{x \in \mathcal{R}_n^t} |f(x,t) - \tilde{f}_{\bar{n}_t}(x,t)| \\ &\leq \left(\sqrt{\frac{C_K^t}{\bar{n}_t}} + \sqrt{\sup_{x \in \mathcal{R}_n^{1-t}} \sigma_{\bar{n}_t}(x)}\right) \sqrt{d \log\left(\frac{1 + \bar{n}_t^2 r_t}{\delta}\right)} \\ &+ \mathcal{O}(1/\bar{n}_t), \end{split}$$

where  $C_K^t = 4L_K^t + 2U_K^t/\sigma^2$  is only related to the kernel K and unrelated to the number of sample  $\bar{n}_t$ ;  $r_t = \max_{x,x' \in \mathcal{R}_n^{1-t}} ||x - x'||$  is the radius of the augmentation region.

(10)

*Remark* 4.8. To control the error, observe that  $\sigma_{\bar{n}_t}(x)$  is a decreasing function of  $\bar{n}_t$  while  $\sup_{x \in \mathcal{R}_n^t} \sigma_{\bar{n}_t}(x)$  is an increasing function of the size of the augmentation region  $\mathcal{R}_n^t$ . Therefore, the data augmentation region must be chosen carefully such that it can be controlled and diminishes asymptotically to zero.

Proposition 4.7 is a sufficient condition for controlling term (III) in Proposition 3.1 due to the fact that

$$\mathbb{E}_{X,T\sim q}\left[\left\|f(X,T) - \tilde{f}_n(X,T)\right\|\right]$$
  
$$\leq \sup_{t\in\{0,1\}} \sup_{X\in\mathcal{R}_n^{1-t}} |f(X,t) - \tilde{f}_n(X,t)|,$$

and the following result:

**Proposition 4.9.** With probability at least  $1 - \delta$  where

 $\delta \in (0, 1),$ 

$$\sup_{t \in \{0,1\}} \sup_{x \in \mathcal{R}_n^{1-t}} |f(x,t) - \tilde{f}_{\bar{n}_t}(x,t)| \\
\leq \sqrt{d} \tilde{\mathcal{O}} \left( \sqrt{\frac{C_K^0 \vee C_K^1}{\bar{n}_0 \wedge \bar{n}_1}} + \sqrt{\sup_{x \in \mathcal{R}_n^1} \sigma_{\bar{n}_0}(x) \vee \sup_{x \in \mathcal{R}_n^0} \sigma_{\bar{n}_1}(x)} \right) \\
+ \mathcal{O}(1/(\bar{n}_0 \wedge \bar{n}_1)), \tag{11}$$

with all the constants defined in Proposition 4.7.

*Remark* 4.10. As proved in Proposition 4.1, for any number of required neighbors  $\bar{n}_t$ , the probability of a fixed x not having more than  $\bar{n}_t$  neighbors *decreases approximately exponentially* to 0. This implies that *the imputation error with local GP can be effectively controlled*. As the righthand side of Equation (10) converges to 0 as  $n \to +\infty$ , this demonstrates that asymptotically POLO with local GP will lead to unbiased learning of CATE.

*Remark* 4.11. POLO carefully selects the subset of individuals for counterfactual outcome imputation so that

- By only selecting individuals with a sufficient amount of close neighbors, R<sup>1−t</sup><sub>n</sub> is reduced. σ<sub>n̄t</sub>(x) is also decreased as the posterior of GP has less variance with more close neighbors. Hence, sup<sub>x∈R<sup>1−t</sup><sub>n</sub></sub> σ<sub>n̄t</sub>(x) is significantly reduced, leading to reduced error.
- Smaller  $\mathcal{R}_n^{1-t}$  decrease both  $L_K^t$  and  $U_K^t$ , further decreasing the error.

*Remark* 4.12. The effect of the complexity of the true causal function f is captured in  $C_K^t$  and  $\sigma_{\bar{n}_t}(x)$ : a simpler f implies smoother kernel thus smaller  $C_K^t$  and faster decrease of  $\sigma_{\bar{n}_t}(x)$ .

# 5 COCOA: CONTRASTIVE COUNTERFACTUAL AUGMENTATION

To further boost the performance of local regression methods, we propose to employ contrastive learning for an enhanced classifier g to select neighbors. Our motivation is that the success of local regression methods depend on the strength of local correlation. In other words, neighbors selected by the classifier g should exhibit similar outcomes when subjected to the same treatment. *This property, however, may not hold in the original feature space under an arbitrarily selected metric* g. To this end, we propose to use contrastive learning to learn a classifier  $g_{\theta}$  with parameters  $\theta$  and a latent representation space which verify this desired property.

Contrastive (representation) learning methods [Wu et al., 2018, Bojanowski and Joulin, 2017, Dosovitskiy et al., 2014, Caron et al., 2020, He et al., 2020, Chen et al., 2020b, Trinh et al., 2019, Misra and Maaten, 2020, Tian et al., 2020] are based on the principle that similar individuals should



Figure 1: t-SNE visualization of IHDP features and potential outcome  $Y_0$  in the ambient space (left) and the latent space (right) learned by contrastive learning. Groups are defined by dividing the potential outcome  $Y_0$  values into five equal intervals from smallest to largest, with each individual labeled based on the value of its potential outcome.

be associated with closely related representations within an embedding space. This is achieved by training models to perform an auxiliary task: predicting whether two individuals are similar or dissimilar. In the context of CATE estimation, we consider two individuals as similar individuals if they show similar outcomes under the same treatment. Figure 1 illustrates this: *with contrastive learning, the features of the individuals with similar potential outcomes are more clustered in the representation space*, demonstrating the smoothness property that enables reliable local imputation.

**Module Training.** The degree of similarity between outcomes is measured using a particular metric in the potential outcome space  $\mathcal{Y}$ . In our case, we employ the Euclidean norm in  $\mathbb{R}^1$  for this purpose. With this perspective, given the factual (original) dataset  $D_{\mathsf{F}} = \{(x_i, t_i, y_i)\}_{i=1}^n$ , we construct a *positive dataset*  $D_{\epsilon}^+$  that includes pairs of similar individuals. Specifically, we define  $D_{\epsilon}^+ = \{(x_i, x_j, t_i) : i, j \in [n], i \neq j, t_i = t_j, ||y_i - y_j|| \leq \epsilon\}$  where  $\epsilon$  is user-defined sensitivity parameter specifying the desired level of precision. We also create a *negative dataset*  $D^- = \{(x_i, x_j, t_i) : i, j \in [n], i \neq j, t_i = t_j, ||y_i - y_j|| > \epsilon\}$  containing pairs of individuals deemed dissimilar. Let  $\ell : \{0, 1\} \times \{0, 1\} \rightarrow \mathbb{R}$  be any loss function for classification task . We learn a parametric classifier (neural network)  $g_{\theta} : \mathcal{X} \times \mathcal{X} \rightarrow \{0, 1\}$  with parameter  $\theta$  by optimizing the following objective function:

$$\min_{\theta} \sum_{(x,x',t)\in D_{\epsilon}^{+}} \ell(g_{\theta}(x,x',t),1) + \sum_{(x,x',t)\in D_{\epsilon}^{-}} \ell(g_{\theta}(x,x',t),0)$$

**Neighbor Identification.** For a given individual x in  $D_F$  within treatment group t, we utilize trained  $g_{\theta}$  to identify its close neighbors  $D_x$  for counterfactual imputation following Equation (4) with  $g_{\theta}$  as the classifier. We term this proposed approach with a learned classifier as <u>COntrastive</u> <u>COunterfactual Augmentation</u> (COCOA). After selecting the neighbors, the potential outcome imputation of COCOA is the same as POLO. See Algorithm 1 in the Appendix for the Pseudocode of COCOA.

# 6 EMPIRICAL STUDIES

While the theoretical results in Sections 3 and 4.2 provide large-sample guarantees, here we empirically demonstrate that COCOA works for practical scenarios where the number of samples is only moderate. In particular, we observe that COCOA consistently improves the CATE estimation performance across state-of-the-art CATE models. More importantly, we observe that *COCOA prevents CATE models from overfitting to the factual data* during training. We believe this property is particularly important in the setting of CATE estimation because the true performance of models cannot be validated in practice, making robustness to overfitting an especially desirable property.

**Evaluation Setup.** We test our proposed methods on various benchmark datasets: the IHDP dataset [Ramey et al., 1992, Hill, 2011], the News dataset [Johansson et al., 2016, Newman et al., 2008], and the Twins dataset [Louizos et al., 2017]. Additionally, we apply our methods to two synthetic datasets: one with linear functions for potential outcomes and the other with non-linear functions, we include these results in Appendix E.1. A detailed description of these datasets is provided in Appendix D. To estimate the variance of our method, we randomly divide each of these datasets into a train (70%) dataset and a test (30%) dataset with varying seeds. Moreover, we demonstrate the efficacy of our methods across a variety of CATE estimation models.

**Performance Improvements.** Table 1 summarizes the experimental results verifying COCOA's effect on *consistently improving* the performance of various CATE estimation models. We observe significant improvements for certain models over specific benchmarks (e.g., Twins with CFR-Wass, IHDP with CD), lead to new state-of-the-art performance. Moreover, even in cases where the improvement is marginal, we note substantial enhancements in models' robustness to overfitting the factual distribution, as described in the following paragraph.

Robustness Improvements. In the context of CATE esti-



Figure 2: Effects of COCOA on *preventing overfitting*. From left to right: IHDP with TARNet, CFR-Wass, and T-learner. X-axis has the training epochs; Y-axis shows the performance measure (not accessible in practice). The performance of the models trained without data augmentation decreases as the epoch number increases beyond the optimal stopping epoch (blue curves), overfitting to the factual distribution. In contrast, *the error of the models trained with the augmented dataset barely increase* (red curves), demonstrating the effect of COCOA on preventing overfitting.

Table 1:  $\sqrt{\varepsilon_{\text{PEHE}}}$  across models, with COCOA augmentation (w/ aug.) and without augmentation (w/o aug.) on Twins, News, and IHDP datasets. Lower  $\sqrt{\varepsilon_{\text{PEHE}}}$  corresponds to better performance.

	Twins		News		IHDP	
Model	w/o aug.	w/ aug.	w/o aug.	w/ aug.	w/o aug.	w/ aug.
TARNet	$0.59 \pm .29$	$0.57 \pm .32$	$5.34 \pm .34$	$5.31 {\pm}.17$	$0.92 \pm .01$	$0.87 {\pm}.01$
CFR-Wass	$0.50 \pm .13$	$0.14 \pm .10$	$3.51 \pm .08$	$3.47 {\pm} .09$	$0.85 \pm .01$	$0.83 {\pm .01}$
CFR-MMD	$0.19 \pm .09$	$0.18 \pm .12$	$5.05 \pm .12$	$4.92 {\pm} .10$	$0.87 \pm .01$	$0.85 {\pm .01}$
T-Learner	$0.11 \pm .03$	$0.10 \pm .03$	$4.79 \pm .17$	$4.73 \pm .18$	$2.03 \pm .08$	$1.69 {\pm} .03$
S-Learner	$0.90 \pm .02$	$0.81 {\pm}.06$	$3.83 \pm .06$	$3.80 {\pm} .06$	$1.85 \pm .12$	$0.86 {\pm .01}$
BART	$0.57 \pm .08$	$0.56 \pm .08$	$3.61 \pm .02$	$3.55 {\pm}.00$	$0.67 \pm .00$	$0.67 {\pm}.00$
CF	$0.57 \pm .08$	$0.51 \pm .11$	$3.58 \pm .01$	$3.56{\scriptstyle \pm .01}$	$0.72 \pm .01$	$0.63 {\pm}.01$

mation, it is essential to notice the absence of a validation dataset due to the unavailability of the counterfactual outcomes. This poses a challenge in preventing the models from overfitting to the factual distribution. Our proposed data augmentation technique effectively addresses this challenge, as illustrated in Figure 2, resulting in a significant enhancement of the overall effectiveness of various CATE estimation models. Notably, counterfactual balancing frameworks [Johansson et al., 2016, Shalit et al., 2017] significantly benefit from COCOA. This improvement can be attributed to the fact that data augmentation in dense regions helps narrow the discrepancy between the distributions of the control and the treatment groups. We include more results in Appendix E.7.

Ablation Studies. We conducted ablation studies to assess the impact of the embedding ball size (R) and the number of neighbors (k) on the performance of CATE estimation models trained on the IHDP dataset. Detailed results are in Appendix E.6. These experiments illustrate the trade-off between the quality of imputation and the discrepancy of the treatment groups. *COCOA is robust to the choice of these hyperparameters*, with a wide range of values leading to performance improvements. Table 2 compares our contrastive learning method to propensity scores and Euclidean distance as similarity measures. *The significantly improved*  Table 2:  $\sqrt{\varepsilon_{\text{PEHE}}}$  across different similarity measures: Contrastive Learning (CL), propensity scores (PS), and Euclidean distance (ED), using CFR-Wass across IHDP, News, and Twins datasets.

	ED	PS	CL
IHDP	$3.32{\pm}1.13$	$3.94{\pm}0.21$	$0.83 \pm 0.01$
News	$4.98 \pm 0.10$	$4.82 \pm 0.11$	$3.47 \pm 0.09$
Twins	$0.23 \pm 0.10$	$0.48 {\pm} 0.09$	$0.14{\scriptstyle \pm 0.10}$

*performance of contrastive learning over other close neighbor classifiers proves its efficacy.* Moreover, Appendix E.4 includes ATE estimation results, and Appendix E.5 covers ablations on GP and local linear regression kernels.

# 7 CONCLUSION

We present a data augmentation method for CATE estimation based on potential outcome imputation and local regression. We propose a generalization bound motivating our approach. We provide both asymptotic and finite sample guarantees to support the proposed method. Notably, we enhance both the performance and robustness of various CATE estimation models across various datasets.

#### Acknowledgments

Ahmed Aloui, Juncheng Dong, and Vahid Tarokh were supported in part by the National Science Foundation (NSF) under the National AI Institute for Edge Computing Leveraging Next Generation Wireless Networks Grant # 2112562.



#### References

- Ahmed M Alaa and Mihaela Van Der Schaar. Bayesian inference of individualized treatment effects using multitask gaussian processes. *Advances in neural information processing systems*, 30, 2017.
- Susan Athey and Guido Imbens. Recursive partitioning for heterogeneous causal effects. *Proceedings of the National Academy of Sciences*, 113(27):7353–7360, 2016.
- Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1):151–175, 2010.
- Piotr Bojanowski and Armand Joulin. Unsupervised learning by predicting noise. In *International Conference on Machine Learning*, pages 517–526. PMLR, 2017.
- Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems*, 33: 9912–9924, 2020.
- Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic minority oversampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.
- Pengguang Chen, Shu Liu, Hengshuang Zhao, and Jiaya Jia. Gridmask data augmentation. *arXiv preprint arXiv:2001.04086*, 2020a.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020b.

- Victor Chernozhukov, Whitney Newey, Victor M Quintas-Martinez, and Vasilis Syrgkanis. Riesznet and forestriesz: Automatic debiased machine learning with neural nets and random forests. In *International Conference on Machine Learning*, pages 3901–3914. PMLR, 2022.
- Hugh A Chipman, Edward I George, and Robert E McCulloch. Bart: Bayesian additive regression trees. *The Annals of Applied Statistics*, 2010.
- Peng Ding and Fan Li. Causal inference. *Statistical Science*, 33(2):214–237, 2018.
- Alexey Dosovitskiy, Jost Tobias Springenberg, Martin Riedmiller, and Thomas Brox. Discriminative unsupervised feature learning with convolutional neural networks. *Advances in neural information processing systems*, 27, 2014.
- Steven Y Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. A survey of data augmentation approaches for nlp. arXiv preprint arXiv:2105.03075, 2021.
- Michele Jonsson Funk, Daniel Westreich, Chris Wiesen, Til Stürmer, M Alan Brookhart, and Marie Davidian. Doubly robust estimation of causal effects. *American journal of epidemiology*, 173(7):761–767, 2011.
- Paul Goldsmith-Pinkham, Karen Jiang, Zirui Song, and Jacob Wallace. Measuring changes in disparity gaps: An application to health insurance. In AEA Papers and Proceedings, volume 112, pages 356–360. American Economic Association 2014 Broadway, Suite 305, Nashville, TN 37203, 2022.
- Hui Han, Wen-Yuan Wang, and Bing-Huan Mao. Borderline-smote: a new over-sampling method in imbalanced data sets learning. In Advances in Intelligent Computing: International Conference on Intelligent Computing, ICIC 2005, Hefei, China, August 23-26, 2005, Proceedings, Part I 1, pages 878–887. Springer, 2005.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020.
- Jennifer Hill, Antonio Linero, and Jared Murray. Bayesian additive regression trees: A review and look forward. *Annual Review of Statistics and Its Application*, 7:251–278, 2020.
- Jennifer L Hill. Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, 20(1):217–240, 2011.

- Keisuke Hirano, Guido W Imbens, and Geert Ridder. Efficient estimation of average treatment effects using the estimated propensity score. *Econometrica*, 71(4):1161–1189, 2003.
- Paul W Holland. Statistics and causal inference. Journal of the American statistical Association, 81(396):945–960, 1986.
- Kosuke Imai, Gary King, and Elizabeth A Stuart. Misunderstandings between experimentalists and observationalists about causal inference. *Journal of the royal statistical society: series A (statistics in society)*, 171(2):481–502, 2008.
- Guido W Imbens and Donald B Rubin. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press, 2015.
- Wei Jiang, Kai Zhang, Nan Wang, and Miao Yu. Meshcut data augmentation for deep learning in computer vision. *Plos one*, 15(12):e0243613, 2020.
- Ziyang Jiang, Zhuoran Hou, Yiling Liu, Yiman Ren, Keyu Li, and David Carlson. Estimating causal effects using a multi-task deep ensemble. *arXiv preprint arXiv:2301.11351*, 2023.
- Fredrik Johansson, Uri Shalit, and David Sontag. Learning representations for counterfactual inference. In *International conference on machine learning*, pages 3020–3029. PMLR, 2016.
- Sören R Künzel, Jasjeet S Sekhon, Peter J Bickel, and Bin Yu. Metalearners for estimating heterogeneous treatment effects using machine learning. *Proceedings of the national academy of sciences*, 116(10):4156–4165, 2019.
- Armin Lederer, Jonas Umlauft, and Sandra Hirche. Uniform error bounds for gaussian process regression with application to safe control. *Advances in Neural Information Processing Systems*, 32, 2019.
- Pei Liu, Xuemin Wang, Chao Xiang, and Weiye Meng. A survey of text data augmentation. In 2020 International Conference on Computer Communication and Network Security (CCNS), pages 191–195. IEEE, 2020.
- 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.
- Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6707–6717, 2020.

- Matthew L Newman, Carla J Groom, Lori D Handelman, and James W Pennebaker. Gender differences in language use: An analysis of 14,000 text samples. *Discourse processes*, 45(3):211–236, 2008.
- Jersey Neyman. Sur les applications de la théorie des probabilités aux experiences agricoles: Essai des principes. *Roczniki Nauk Rolniczych*, 10(1):1–51, 1923.
- Judea Pearl. *Causality: Models, Reasoning and Inference.* Cambridge University Press, USA, 2nd edition, 2009a.
- Judea Pearl. Causal inference in statistics: An overview. *Statistics Surveys*, 3(none):96 146, 2009b.
- Craig T Ramey, Donna M Bryant, Barbara H Wasik, Joseph J Sparling, Kaye H Fendt, and Lisa M La Vange. Infant health and development program for low birth weight, premature infants: Program elements, family participation, and child intelligence. *Pediatrics*, 89(3):454– 465, 1992.
- James Robins. A new approach to causal inference in mortality studies with a sustained exposure period—application to control of the healthy worker survivor effect. *Mathematical modelling*, 7(9-12):1393–1512, 1986.
- Paul R Rosenbaum and Donald B Rubin. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55, 1983.
- Donald B Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of educational Psychology*, 66(5):688, 1974.
- Donald B Rubin. Bayesian inference for causal effects: The role of randomization. *The Annals of statistics*, pages 34–58, 1978.
- Donald B Rubin. Causal inference using potential outcomes: Design, modeling, decisions. *Journal of the American Statistical Association*, 100(469):322–331, 2005.
- Walter Rudin. *Principles of Mathematical Analysis.* McGraw-Hill Book Co., New York, 1st edition, 1953.
- Patrick Schwab, Lorenz Linhardt, and Walter Karlen. Perfect match: A simple method for learning representations for counterfactual inference with neural networks. *arXiv preprint arXiv:1810.00656*, 2018.
- Matthias Seeger. Gaussian processes for machine learning. *International journal of neural systems*, 14(02):69–106, 2004.
- Uri Shalit, Fredrik D Johansson, and David Sontag. Estimating individual treatment effect: generalization bounds and algorithms. In *International Conference on Machine Learning*, pages 3076–3085. PMLR, 2017.

- Niranjan Srinivas, Andreas Krause, Sham M. Kakade, and Matthias W. Seeger. Information-theoretic regret bounds for gaussian process optimization in the bandit setting. *IEEE Transactions on Information Theory*, 58(5):3250– 3265, 2012. doi: 10.1109/TIT.2011.2182033.
- Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What makes for good views for contrastive learning? *Advances in neural information processing systems*, 33:6827–6839, 2020.
- Dag Tjøstheim, Håkon Otneim, and Bård Støve. *Statistical Modeling Using Local Gaussian Approximation*. Academic Press, 2021.
- Trieu H Trinh, Minh-Thang Luong, and Quoc V Le. Selfie: Self-supervised pretraining for image embedding. *arXiv preprint arXiv:1906.02940*, 2019.
- AW van der Vaart and Jon A Wellner. Empirical processes. In Weak Convergence and Empirical Processes: With Applications to Statistics, pages 127–384. Springer, 2023.
- David A Van Dyk and Xiao-Li Meng. The art of data augmentation. *Journal of Computational and Graphical Statistics*, 10(1):1–50, 2001.
- Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3733–3742, 2018.
- Jinsung Yoon, James Jordon, and Mihaela Van Der Schaar. Ganite: Estimation of individualized treatment effects using generative adversarial nets. In *International Conference on Learning Representations*, 2018.

# CATE Estimation With Potential Outcome Imputation From Local Regression (Supplementary Material)

Ahmed Aloui <sup>†1</sup>

Juncheng Dong \*1

Cat P. Le<sup>1</sup>

Vahid Tarokh<sup>1</sup>

<sup>1</sup> Department of Electrical and Computer Engineering, Duke University

# A SCHEMATIC ILLUSTRATION OF COCOA

Figure 3a provides an overview of the COCOA framework, where similarity learning and local imputations are leveraged to augment observational data, reducing statistical discrepancies while minimizing imputation error. This augmented dataset is then used to train CATE estimation models, improving accuracy and robustness.

Additionally, Figure 3b illustrates the trade-off between statistical discrepancy and imputation error as the augmentation level varies. This observation motivates the COCOA approach by highlighting the balance required between data alignment and the reliability of imputed counterfactuals.





(a) Similarity learning is used to select a subset of individuals, followed by reliable local imputations to generate their counterfactuals. These imputations augment the original dataset, reducing the statistical discrepancy between treatment groups while minimizing imputation error. The augmented data is then used to train off-the-shelf CATE estimation models, improving their accuracy and robustness.

(b) Trade-off between statistical discrepancy and imputation error across different augmentation levels (0 to 1). A full description of the synthetic toy dataset and implementation details can be found in Appendix E.2.

Figure 3: (a) Overview of the proposed model-agnostic data augmentation method for CATE estimation, and (b) the observed trade-off that motivated the proposed method.

# **B PROOFS OF THE THEORETICAL RESULTS**

In this section, we include the proofs for the theoretical results presented in the main text.

\*Equal contribution.

<sup>†</sup>Equal contribution.

#### **B.1 PROOF OF PROPOSITION 3.1**

To prove the generalization bound, we first define a notion of consistency for data augmentation. And, we demonstrate a lemma proving that the proposed consistency is equivalent to emulating RCTs.

**Definition B.1** (Consistency of Factual Distribution). A factual distribution  $p_F$  is consistent if for every hypothesis  $h: \mathcal{X} \times \{0,1\} \to \mathcal{Y}, \mathcal{L}_{\mathsf{F}}(h) = \mathcal{L}_{\mathsf{CF}}(h).$ 

**Definition B.2** (Consistency of Data Augmentation). A data augmentation method is said to be consistent if the augmented data follows a factual distribution that is consistent.

**Lemma B.3** (Consistency is Equivalent Randomized Controlled Trials). Suppose we have a factual distribution  $p_F$  and its corresponding counterfactual distribution  $p_{CF}$  such that for every hypothesis  $h : \mathcal{X} \times \{0, 1\} \to \mathcal{Y}, \mathcal{L}_F(h) = \mathcal{L}_{CF}(h)$ . This implies that the data must originate from a randomized controlled trial, i.e.,  $p_F(X|T=1) = p_F(X|T=0)$ .

Proof of Proposition B.3. Suppose that for every hypothesis  $h : \mathcal{X} \times \{0,1\} \to \mathcal{Y}, \mathcal{L}_{F}(h) = \mathcal{L}_{CF}(h)$ . By definition,

$$\mathcal{L}_{\rm F}(h) = \int (y - h(x,t))^2 p_{\rm F}(x,t,y) \, dx \, dt \, dy$$

and

$$\mathcal{L}_{\rm CF}(h) = \int (y - h(x, t))^2 p_{\rm CF}(x, t, y) \, dx \, dt \, dy$$

We can write this as

$$\mathbb{E}_{p_{\mathsf{F}}}\left[\left(Y - h(X, T)^{2}\right)\right] = \mathbb{E}_{p_{\mathsf{CF}}}\left[\left(Y - h(X, T)^{2}\right)\right]$$

Since this holds for every function h, consider two Borel sets A and B in  $\mathcal{X} \times \mathcal{T} \times \mathcal{Y}$ , and we let  $h_1(X,T) = \mathbb{E}[Y|X,T] - \mathbb{1}_A$ and  $h_2(X,T) = \mathbb{E}[Y|X,T] - \mathbb{1}_B$ . Hence we have that,

$$\mathbb{E}_{p_{\mathrm{F}}}\left[\left(Y - h_{1}(X,T)\right)^{2}\right] = \mathbb{E}_{p_{\mathrm{F}}}\left[\left(Y - \mathbb{E}\left[Y|X,T\right] + \mathbb{1}_{A}\right)^{2}\right]$$
$$= \mathbb{E}_{p_{\mathrm{F}}}\left[\left(Y - \mathbb{E}\left[Y|X,T\right]\right)^{2}\right] + \mathbb{E}_{p_{\mathrm{F}}}\left[\mathbb{1}_{A}\right] + 2\mathbb{E}_{p_{\mathrm{F}}}\left[\mathbb{1}_{A}\left(Y - \mathbb{E}\left[Y|X,T\right]\right)\right]$$

And we have that,  $\mathbb{E}_{p_{\mathrm{F}}}\left[\mathbbm{1}_{A}\left(Y - \mathbb{E}\left[Y|X,T\right]\right)\right] = 0$  since by definition of the conditional expectation we have that  $\mathbb{E}[Y\mathbbm{1}_{A}] = \mathbb{E}\left[\mathbb{E}\left[Y|X,T\right]\mathbbm{1}_{A}\right]$ . We denote by  $MSE(p_{\mathrm{F}}) = \mathbb{E}_{p_{\mathrm{F}}}\left[\left(Y - \mathbb{E}\left[Y|X,T\right]\right)^{2}\right]$ . Therefore we have that

$$\mathbb{E}_{p_{\mathrm{F}}}\left[\left(Y - h_{1}(X, T)\right)^{2}\right] = MSE(p_{\mathrm{F}}) + \mathbb{E}_{p_{\mathrm{F}}}\left[\mathbbm{1}_{A}\right]$$

Using the same argument for  $p_{\rm CF}$  we have the following result:

$$\mathbb{E}_{p_{\mathsf{CF}}}\left[\left(Y - h_1(X, T)\right)^2\right] = MSE(p_{\mathsf{CF}}) + \mathbb{E}_{p_{\mathsf{CF}}}\left[\mathbbm{1}_A\right]$$

Similarly, we have the following for  $h_2$ :

$$\mathbb{E}_{p_{\mathrm{F}}}\left[\left(Y - h_2(X, T)\right)^2\right] = MSE(p_{\mathrm{F}}) + \mathbb{E}_{p_{\mathrm{F}}}\left[\mathbbm{1}_B\right]$$
$$\mathbb{E}_{p_{\mathrm{CF}}}\left[\left(Y - h_2(X, T)\right)^2\right] = MSE(p_{\mathrm{CF}}) + \mathbb{E}_{p_{\mathrm{CF}}}\left[\mathbbm{1}_B\right]$$

Therefore we have

$$MSE(p_{\rm F}) - MSE(p_{\rm CF}) = \mathbb{E}_{p_{\rm F}}\left[\mathbb{1}_A\right] - \mathbb{E}_{p_{\rm CF}}\left[\mathbb{1}_A\right]$$

and

$$MSE(p_{\rm F}) - MSE(p_{\rm CF}) = \mathbb{E}_{p_{\rm F}}\left[\mathbb{1}_B\right] - \mathbb{E}_{p_{\rm CF}}\left[\mathbb{1}_B\right]$$

Therefore

$$\mathbb{E}_{p_{\mathrm{F}}}\left[\mathbbm{1}_{A}\right] - \mathbb{E}_{p_{\mathrm{CF}}}\left[\mathbbm{1}_{A}\right] = \mathbb{E}_{p_{\mathrm{F}}}\left[\mathbbm{1}_{B}\right] - \mathbb{E}_{p_{\mathrm{CF}}}\left[\mathbbm{1}_{B}\right]$$

Hence it follows,

$$\mathbb{E}_{p_{\mathrm{F}}}\left[\mathbbm{1}_{A\cap B}\right] = \mathbb{E}_{p_{\mathrm{CF}}}\left[\mathbbm{1}_{A\cap B}\right]$$

And as this holds for every Borel measurable set A and B, therefore we have that  $p_{\rm F} = p_{\rm CF}$ .

Denote by  $u = p_{\rm F}(T=1)$  we have  $p_{\rm F}(X) = up_{\rm F}(X|T=1) + (1-u)p_{\rm F}(X|T=0)$ . Similarly we have that  $p_{\rm CF}(X) = (1-u)p_{\rm CF}(X|T=1) + up_{\rm CF}(X|T=0)$ . Therefore, since  $p_{\rm F} = p_{\rm CF}$ ,

$$up_{\rm F}(X|T=1) + (1-u)p_{\rm F}(X|T=0) = (1-u)p_{\rm CF}(X|T=1) + up_{\rm CF}(X|T=0)$$
$$= (1-u)p_{\rm F}(X|T=1) + up_{\rm F}(X|T=0)$$

Hence

$$(2u-1) p_{\rm F}(X|T=1) = (2u-1) p_{\rm F}(X|T=0)$$

Therefore we conclude the result that,

$$p_{\rm F}(X|T=1) = p_{\rm F}(X|T=0)$$

This concludes the proof.

For completeness, we also include this result.

**Lemma B.4** (Consistency of Randomized Controlled Trials). *The factual distribution of any randomized controlled trial* =verifying  $p_F(T = 1) = p_F(T = 0)$  is consistent, i.e., if  $p_F(X|T = 1) = p_F(X|T = 0)$  and  $p_F(T = 1) = p_F(T = 0)$ , then for all  $h : \mathcal{X} \times \{0, 1\} \to \mathcal{Y}$ ,

$$\mathcal{L}_{\scriptscriptstyle F}(h) = \mathcal{L}_{\scriptscriptstyle CF}(h)$$

*Proof.* Let  $u = p_F(T = 1) = \frac{1}{2}$ ,  $p_F(T = 1) = p_{CF}(T = 0)$ 

$$\begin{split} \mathcal{L}_{\rm F}(h) &= \int (y - h(x,t))^2 p_{\rm F}(x,t,y) \, dx, \, dt \, dy \\ &= u \int (y - h(x,1))^2 p_{\rm F}(x,y|T=1) \, dx \, dy + (1-u) \int (y - h(x,0))^2 p_{\rm F}(x,y|T=0) \, dx \, dy \\ &= u \int (y - h(x,1))^2 p_{\rm F}(x,y|T=0) \, dx \, dy + (1-u) \int (y - h(x,0))^2 p_{\rm F}(x,y|T=1) \, dx \, dy \\ &= u \int (y - h(x,1))^2 p_{\rm CF}(x,y|T=1) \, dx \, dy + (1-u) \int (y - h(x,0))^2 p_{\rm CF}(x,y|T=0) \, dx \, dy \\ &= \int (y - h(x,t))^2 p_{\rm CF}(x,t,y) \, dx \, dy \\ &= \mathcal{L}_{\rm CF}(h) \end{split}$$

To prove Proposition 3.1 we also include a new definition for an "ideal" factual distribution. Subsequently, we will prove its consistency. The ideal factual distribution is defined as follows:

$$p_{\rm IF} = \frac{1}{2} p_{\rm F} + \frac{1}{2} p_{\rm CF}.$$
 (12)

In other words, to sample a dataset from  $p_{IF}$ , we sample from the factual distribution  $p_F$  half of the time and from the counterfactual distribution  $p_{CF}$  in the other half of the times. Let  $p_{ICF}$  denote the counterfactual distribution corresponding to  $p_{IF}$ . We next show that  $p_{IF}$  is consistent (thus called ideal distribution).

**Lemma B.5** (Consistency of  $p_{\text{IF}}$ .). The error of the ideal factual distribution equals the error of its corresponding counterfactual distribution, i.e., for every hypothesis  $h : \mathcal{X} \times \{0, 1\} \to \mathcal{Y}$ , we have that  $\mathcal{L}_{IF}(h) = \mathcal{L}_{ICF}(h)$ .

*Proof.* We observe that 
$$p_{ICF} = \frac{1}{2}p_{CF} + \frac{1}{2}p_{F}$$
. Therefore,  $p_{ICF} = p_{IF}$  and the result follows.

Intuitively, this result is saying that the ideal counterfactual augmentation gives us a factual distribution that perfectly balances the factual and counterfactual worlds. It follows from Lemma B.3 that achieving this property guarantees that the dataset is identically distributed to the one generated from a Randomized Controlled Trial. However, it is impossible to sample from  $p_{\rm CF}$ .

Also, we cite this Theorem that we will use in our proof:

**Theorem B.6** (Theorem 1 in Ben-David et al. [2010]). Let f be the true function for a learning task such that  $f(x) = \mathbb{E}[Y|X = x]$  where X has a density p and let another true function  $g(x) = \mathbb{E}[Y|X = x]$  modeling another learning task, where X has a density q. Let h by a hypothesis function estimating the true function f, therefore we have

$$\mathbb{E}_{X \sim q(x)}[\|g(X) - h(X)\|^2] \le \mathbb{E}_{X \sim p(x)}[\|f(X) - h(X)\|^2] + 2V(p(x), p(x)) + \mathbb{E}_{X \sim p(x)}[\|f(X) - g(X)\|^2]$$

We can now prove Proposition 3.1.

*Proof.* We have  $f : \mathcal{X} \times \{0, 1\} \to \mathcal{Y}$  to be the function underlying the true causal relationship between (X, T) and Y.It follows from Theorem B.6 that:

$$\mathcal{L}_{\mathrm{IF}}(h) \leq \mathcal{L}_{\mathrm{AF}}(h) + 2V(p_{\mathrm{IF}}, p_{\mathrm{AF}}) + \mathbb{E}_{x, t \sim p_{\mathrm{AF}}}[\|f(x, t) - \tilde{f}(x, t)\|^2]$$

where  $\mathcal{L}_{IF}$  is the factual loss with respect to the ideal density and  $\mathcal{L}_{AF}$  is the factual loss with respect to the density of the augmented data.

By decomposition of the  $\varepsilon_{\text{PEHE}}$  we have that,

$$\begin{split} \varepsilon_{\text{PEHE}}(h) &= \int_{\mathcal{X}} \left( h(x,1) - h(x,0) - f(x,1) + f(x,0) \right)^2 p_{\text{IF}}(x) dx \\ &= \int_{\mathcal{X}} \left( h(x,1) - h(x,0) - f(x,1) + f(x,0) \right)^2 p_{\text{IF}}(x|T=1) p(T=1) dx dt \\ &+ \int_{\mathcal{X}} \left( h(x,1) - h(x,0) - f(x,1) + f(x,0) \right)^2 p_{\text{IF}}(x|T=0) p(T=0) dx dt \\ &\leq 2 \cdot \mathcal{L}_{\text{IF}}(h) + 2 \cdot \mathcal{L}_{\text{ICF}}(h) \end{split}$$

Therefore, it follows from Lemma B.5 that,

$$\varepsilon_{\text{PEHE}}(h) \leq 4 \cdot \left( \mathcal{L}_{\text{AF}}(h) + 2V(p_{\text{RCT}}(x,t), p_{\text{AF}}(x,t)) + \mathbb{E}_{x,t \sim p_{\text{AF}}}[\|f(x,t) - \tilde{f}_n(x,t)\|^2] \right)$$

And since we have that,

$$\mathbb{E}_{x,t \sim p_{\text{AF}}}[\|f(x,t) - \tilde{f}_{n}(x,t)\|^{2}]) = \\ (\frac{1}{1+\beta}) \cdot \mathbb{E}_{x,t \sim p_{F}}[||f(x,t) - \tilde{f}_{n}(x,t)||] + \frac{\beta}{1+\beta} \mathbb{E}_{x,t \sim q}[||f(x,t) - \tilde{f}_{n}(x,t)||]$$

And by observing that the first term  $\mathbb{E}_{x,t\sim p_{\rm F}}[\|f(x,t) - \tilde{f}_n(x,t)\|^2] = 0$ , since the algorithm keeps the samples from the factual distribution to be the same.

## **B.2** PROOF OF PROPOSITION 4.1 AND PROPOSITION 4.2

Proof of Proposition 4.1. We have that,

$$\mathbb{P}(X^{0} \in \mathcal{R}_{n}^{0}) = \sum_{i=k}^{n_{1}} {\binom{n_{1}}{i}} (1-\gamma)^{i} \gamma^{n_{1}-i}$$

$$= 1 - \sum_{i=0}^{k-1} {\binom{n_{1}}{i}} (1-\gamma)^{i} \gamma^{n_{1}-i}$$

$$= 1 - \sum_{i=1}^{k-1} \frac{n_{1}!}{(n_{1}-i)!i!} (1-\gamma)^{i} \gamma^{n_{1}-i}$$

$$\geq 1 - \frac{n_{1}!}{(n_{1}-k+1)!} \gamma^{n_{1}} \sum_{i=0}^{k-1} \frac{1}{i!} \left(\frac{1-\gamma}{\gamma}\right)^{i}$$

$$\geq 1 - n_{1}^{k} \gamma^{n_{1}} \sum_{i=0}^{k-1} \frac{1}{i!} \left(\frac{1-\gamma}{\gamma}\right)^{i}$$

Similarly, we have,

$$\mathbb{P}(X^{1} \in \mathcal{R}_{n}^{1}) = \sum_{i=k}^{n_{0}} \binom{n_{0}}{i} (1-\gamma)^{i} \gamma^{n_{0}-i}$$
$$= 1 - \sum_{i=1}^{k-1} \frac{n_{0}!}{(n_{0}-i)!i!} (1-\gamma)^{i} \gamma^{n_{0}-i}$$
$$\ge 1 - n_{0}^{k} \gamma^{n_{0}} \sum_{i=0}^{k-1} \frac{1}{i!} \left(\frac{1-\gamma}{\gamma}\right)^{i}$$

Therefore we have,

$$1 - \alpha_{n_0} = \mathcal{O}(n_0^k \gamma^{n_0}),$$
  

$$1 - \alpha_{n_1} = \mathcal{O}(n_1^k \gamma^{n_1}),$$

Proof of Proposition 4.2. We start by defining the probability densities of the control and treatment groups resulting from the augmentation process as,  $p_{\rm AF}^1 = \frac{1}{1 + \beta_{n_1}} p^1 + \frac{\beta_{n_1}}{1 + \beta_{n_1}} \frac{p^0 \mathbb{1}_{\mathcal{R}_0}}{\alpha_{n_1}}$ 

and,

$$p_{\rm AF}^{0} = \frac{1}{1 + \beta_{n_0}} p^0 + \frac{\beta_{n_0}}{1 + \beta_{n_0}} \frac{p^1 \mathbb{1}_{\mathcal{R}_1}}{\alpha_{n_0}}$$

with,

and,

$$\beta_{n_0} = \alpha_{n_0} \left( \frac{u}{1-u} \right)$$

 $\beta_{n_1} = \alpha_{n_1} \left( \frac{1-u}{u} \right)$ 

$$\begin{split} V(p_{\rm AF}^1, p_{\rm AF}^0) &= \frac{1}{2} \int |p_{\rm AF}^1 - p_{\rm AF}^0| \\ &= \frac{1}{2} \int \left| \frac{1}{1 + \beta_{n_1}} p^1 + \frac{\beta_{n_1}}{1 + \beta_{n_1}} \frac{p^0 \mathbbm{1}_{\mathcal{R}_0}}{\alpha_{n_1}} - \frac{1}{1 + \beta_{n_0}} p^0 - \frac{\beta_{n_0}}{1 + \beta_{n_0}} \frac{p^1 \mathbbm{1}_{\mathcal{R}_1}}{\alpha_{n_0}} \right| \\ &\leq \frac{1}{2} \int \left| \frac{1}{1 + \beta_{n_1}} p^1 - \frac{\beta_{n_0}}{1 + \beta_{n_0}} \frac{p^1 \mathbbm{1}_{\mathcal{R}_1}}{\alpha_{n_0}} \right| + \frac{1}{2} \int \left| \frac{\beta_{n_1}}{1 + \beta_{n_1}} \frac{p^0 \mathbbm{1}_{\mathcal{R}_0}}{\alpha_{n_1}} - \frac{1}{1 + \beta_{n_1}} p^0 \right| \\ &\leq \frac{1}{1 + \beta_{n_1}} V(p^1, \frac{p^1 \mathbbm{1}_{\mathcal{R}_1}}{\alpha_{n_0}}) + \frac{\beta_{n_0}}{1 + \beta_{n_0}} V(p^0, \frac{p^0 \mathbbm{1}_{\mathcal{R}_0}}{\alpha_{n_1}}) + \left| \frac{1}{1 + \beta_{n_1}} - \frac{\beta_{n_0}}{1 + \beta_{n_0}} \right| \end{split}$$

We have that,

$$V(p^{1}, \frac{p^{1} \mathbb{1}_{\mathcal{R}_{1}}}{\alpha_{n_{0}}}) = \frac{1}{2} \left( \int_{\mathcal{R}_{1}} |p^{1} - \frac{p^{1}}{\alpha_{n_{0}}}| + \int_{\mathcal{R}_{1}^{c}} p^{1} \right)$$
$$= \frac{1}{2} \left( \int_{\mathcal{R}_{1}} p^{1} |1 - \frac{1}{\alpha_{n_{0}}}| + (1 - \alpha_{n_{0}}) \right)$$
$$= \frac{1}{2} \left( \frac{|\alpha_{n_{0}} - 1|}{\alpha_{n_{0}}} \int_{\mathcal{R}_{1}} p^{1} + (1 - \alpha_{n_{0}}) \right)$$
$$= \frac{1}{2} \left( \frac{|\alpha_{n_{0}} - 1|}{\alpha_{n_{0}}} \alpha_{n_{0}} + (1 - \alpha_{n_{0}}) \right)$$
$$= (1 - \alpha_{n_{0}})$$

Similarly,

$$V(p^{0}, \frac{p^{0} \mathbb{1}_{\mathcal{R}_{0}}}{\alpha_{n_{1}}}) = (1 - \alpha_{n_{1}})$$

Substituting this into the bound and letting  $z = \frac{u}{1-u}$  we have that,

$$\begin{split} V(p_{\rm AF}^1, p_{\rm AF}^0) &\leq \frac{1 - \alpha_{n_0}}{1 + \beta_{n_1}} + \frac{\beta_{n_0}(1 - \alpha_{n_1})}{1 + \beta_{n_0}} \alpha_{n_1}) + \left| \frac{1}{1 + \beta_{n_1}} - \frac{\beta_{n_0}}{1 + \beta_{n_0}} \right| \\ &= \frac{1 - \alpha_{n_0}}{1 + \psi^{-1}\alpha_{n_1}} + \frac{\psi \alpha_{n_0} \left(1 - \alpha_{n_1}\right)}{1 + \alpha_{n_0} \psi} + \frac{\left| 1 - \alpha_{n_1} \alpha_{n_0} \right|}{\left(1 + \psi^{-1} \alpha_{n_1}\right) \left(1 + \alpha_{n_0} \psi\right)} \end{split}$$

## **B.3 PROOF OF PROPOSITION 4.7 AND PROPOSITION 4.9**

*Proof of Proposition 4.7.* The proof for t = 0 and t = 1 is symmetric, thus fix  $t \in \{0, 1\}$ . For notational simplicity, we use z in the proof to denote  $\bar{n}_t$ , and let

$$A = (K(\mathbf{x}_z, \mathbf{x}_z) + \sigma^2 \cdot I_z)^{-1} \in \mathbb{R}^{z \times z}.$$

and

$$U_K^t = \max_{x, x' \in \mathcal{R}_n^t} K(x, x').$$

Consider  $\tau > 0$ . A set S is a  $\tau$ -cover for  $\mathcal{R}_n^{1-t}$  if  $\forall x \in \mathcal{R}_n^{1-t}, \exists x' \in S$  such that  $||x' - x|| \leq \tau$ . Let  $\mathcal{C}(\tau, \mathcal{R}_n^{1-t})$  be the covering number of  $\mathcal{R}_n^{1-t}$  with radius  $\tau$ :

$$\mathcal{C}(\tau, \mathcal{R}_n^{1-t}) = \inf\{|S| : S \text{ is } \tau \text{-cover of } \mathcal{R}_n^{1-t}\}.$$

Since  $\mathcal{R}_n^{1-t} \subset \mathbb{R}^d$ , we have Vaart and Wellner [2023]

$$\mathcal{C}(\tau, \mathcal{R}_n^{1-t}) \leq \left(1 + \frac{r}{\tau}\right)^d,$$

where  $r = \max_{x,x' \in \mathcal{R}_n^{1-t}} ||x - x'||$ . Consider a minimum  $\tau$ -cover  $\mathcal{C}_{\tau}$  for  $\mathcal{R}_n^{1-t}$  with (by definition of covering number)  $\mathcal{C}(\tau, \mathcal{R})$  elements. We have that Srinivas et al. [2012], with probability at least  $1 - \mathcal{C}(\tau, \mathcal{R}) \exp(-\xi(\tau)/2)$ ,

$$\sup_{x \in \mathcal{C}_{\tau}} |f(x,t) - \tilde{f}_n(x,t)| \le \sqrt{\xi(\tau)} \sup_{x \in \mathcal{C}_{\tau}} \sigma_n(x).$$

Choosing  $\xi(\tau) = 2 \log(\mathcal{C}(\tau, \mathcal{R})/\delta)$ , we have with probability  $1 - \delta$ ,

$$\sup_{x \in \mathcal{C}_{\tau}} |f(x,t) - \tilde{f}_n(x,t)| \le \sqrt{\xi(\tau)} \sup_{x \in \mathcal{C}_{\tau}} \sigma_n(x).$$

Moreover, by definition of  $C_{\tau}$ ,  $\max_{x \in \mathcal{R}_n^t} \min_{x' \in \mathcal{C}_{\tau}} ||x - x'|| \leq \tau$ . Because f(x, t) is  $L_f$ -Lipschitz continuous, we have for all  $x \in \mathcal{R}_n^{1-t}$ 

$$\min_{x'\in\mathcal{C}_{\tau}}|f(x,t)-f(x',t)|\leq\tau L_f.$$

With the fact that Lederer et al. [2019]  $\tilde{f}_z(x,t)$  and  $\sigma_z(x)$  is Lipschitz continuous with respective Lipschitz constant

$$C_1 = L_K \sqrt{z} ||A\mathbf{y}_n||,\tag{13}$$

$$C_2(\tau) = \sqrt{2\tau L_K (1 + z \cdot ||A||_F \cdot U_K^t)},$$
(14)

we have with probability at least  $1 - \delta$  that

$$\sup_{x \in \mathcal{R}_n^{1-t}} |\tilde{f}_z(x,t) - f(x,t)| \le \sqrt{\xi(\tau)} \sup_{x \in \mathcal{R}_n^{1-t}} \sigma_z(x) + C_2(\tau)\sqrt{\xi(\tau)} + (C_1 + L_f)\tau$$

To continue, we will proceed to upper bound  $C_1$ :

$$C_1 = L_K \sqrt{z} ||A\mathbf{y}_z|| \le L_K \sqrt{z} ||A||_F ||\mathbf{y}_z|| \le L_K \sqrt{z} \frac{||\mathbf{y}_z||}{\sigma^2}$$

due to the fact that  $||A||_F \leq 1/\sigma^2$ . Assume that  $f(x,t) \leq F \leq +\infty$ , by the assumption of the data generation process  $y = f(x,t) + \epsilon$ ,  $\epsilon \sim \mathcal{N}(0,\sigma^2)$ , and triangular inequality of norm,

$$||\mathbf{y}_z|| \le ||f(\mathbf{x}_z, \mathbf{t}_z)|| + ||\gamma_z||$$
(15)

$$\leq \sqrt{z}F + ||\gamma_z||,\tag{16}$$

where  $\gamma_z$  is a multi-variate Gaussian random variable in  $\mathbb{R}^z$  with mean 0 and covariance matrix  $\sigma^2 \cdot I_z$ . Hence  $||\gamma_z||/\sigma^2$  is a Chi-squared random variable with degrees of freedom equal to z. Then we have with probability at least  $1 - \delta/2$ ,

$$C_1 \le L_K (zF + 2z\sqrt{\eta_z \sigma^2})/\sigma^2,$$

where  $\eta_z = \log(\pi^2 z^2/\delta)$ . On the other hand,  $C_2$  can be upper bounded as

$$C_2(\tau) \le \sqrt{2\tau L_K (1 + z \cdot U_K^t / \sigma^2)}.$$

Hence, by choosing  $\tau = 1/z^2$ , we have

$$(C_1 + L_f)\tau \in \mathcal{O}(1/z),$$

and with probability at least  $1 - \delta$ , we have

$$\sup_{X \in \mathcal{R}} |f(X,t) - \tilde{f}_n(X,t)| \le \sqrt{\frac{4L_K + 2U_K/\sigma^2}{z}d\log(1+z^2r)} + \sqrt{2d\log(1+z^2r)} \sup_{x \in \mathcal{R}_n^{1-t}} \sigma_n(x) + \mathcal{O}(1/z)$$

After reorganizing terms, the proof is complete.

*Proof of Proposition 4.9.* With Proposition 4.7, we have probability at least  $(1 - \delta)^2$  that for both t = 0 and t = 1

$$\sup_{x \in \mathcal{R}_n^{1-t}} |f(x,t) - \tilde{f}_{\bar{n}_t}(x,t)| \le \left(\sqrt{\frac{C_K^t}{\bar{n}_t}} + \sqrt{\sup_{x \in \mathcal{R}_n^{1-t}} \sigma_{\bar{n}_t}(x)}\right) \sqrt{d \log\left(\frac{1 + \bar{n}_t^2 r_t}{\delta}\right)} + \mathcal{O}(1/\bar{n}_t),$$

This implies that

$$\begin{split} \sup_{t \in \{0,1\}} \sup_{x \in \mathcal{R}_{n}^{1-t}} |f(x,t) - \tilde{f}_{\bar{n}_{t}}(x,t)| &\leq \sup_{t \in \{0,1\}} \left\{ \left( \sqrt{\frac{C_{K}^{t}}{\bar{n}_{t}}} + \sqrt{\sup_{x \in \mathcal{R}_{n}^{1-t}} \sigma_{\bar{n}_{t}}(x)} \right) \sqrt{d \log \left(\frac{1 + \bar{n}_{t}^{2} r_{t}}{\delta}\right)} + \mathcal{O}(1/\bar{n}_{t}) \right\} \\ &\leq \sup_{t \in \{0,1\}} \left\{ \left( \sqrt{\frac{C_{K}^{t}}{\bar{n}_{t}}} + \sqrt{\sup_{x \in \mathcal{R}_{n}^{1-t}} \sigma_{\bar{n}_{t}}(x)} \right) \sqrt{d \log \left(\frac{1 + \bar{n}_{t}^{2} r_{t}}{\delta}\right)} \right\} + \mathcal{O}(1/\bar{n}_{0} \wedge \bar{n}_{1}) \\ &\leq \sup_{t \in \{0,1\}} \left\{ \left( \sqrt{\frac{C_{K}^{0} \vee C_{K}^{1}}{\bar{n}_{0} \wedge \bar{n}_{1}}} + \sqrt{\sup_{x \in \mathcal{R}_{n}^{1-t}} \sigma_{\bar{n}_{t}}(x)} \right) \sqrt{d \log \left(\frac{1 + \bar{n}_{t}^{2} r_{t}}{\delta}\right)} \right\} + \mathcal{O}(1/\bar{n}_{0} \wedge \bar{n}_{1}) \\ &\leq \sqrt{d} \left( \sqrt{\frac{C_{K}^{0} \vee C_{K}^{1}}{\bar{n}_{0} \wedge \bar{n}_{1}}} + \sup_{t \in \{0,1\}} \sqrt{\sup_{x \in \mathcal{R}_{n}^{1-t}} \sigma_{\bar{n}_{t}}(x)} \right) \sqrt{\log \left(\frac{1 + (\bar{n}_{0} \vee \bar{n}_{1})^{2} r_{t}}{\delta}\right)} + \mathcal{O}(1/\bar{n}_{0} \wedge \bar{n}_{1}) \end{split}$$

By change of variable  $(1 - \delta)^2 = 1 - \delta'$ , we have with probability  $1 - \delta'$  for  $\delta' \in (0, 1)$ ,

$$\sup_{t \in \{0,1\}} \sup_{x \in \mathcal{R}_{n}^{1-t}} |f(x,t) - f_{\bar{n}_{t}}(x,t)| \\
\leq \sqrt{d} \left( \sqrt{\frac{C_{K}^{0} \vee C_{K}^{1}}{\bar{n}_{0} \wedge \bar{n}_{1}}} + \sup_{t \in \{0,1\}} \sqrt{\sup_{x \in \mathcal{R}_{n}^{1-t}} \sigma_{\bar{n}_{t}}(x)} \right) \sqrt{\log \left(\frac{1 + (\bar{n}_{0} \vee \bar{n}_{1})^{2} r_{t}}{\sqrt{1 - \sqrt{1 - \delta'}}}\right)} + \mathcal{O}(1/\bar{n}_{0} \wedge \bar{n}_{1}) \\
= \sqrt{d} \tilde{\mathcal{O}} \left( \sqrt{\frac{C_{K}^{0} \vee C_{K}^{1}}{\bar{n}_{0} \wedge \bar{n}_{1}}} + \sqrt{\sup_{x \in \mathcal{R}_{n}^{1}} \sigma_{\bar{n}_{0}}(x) \vee \sup_{x \in \mathcal{R}_{n}^{0}} \sigma_{\bar{n}_{1}}(x)} \right) + \mathcal{O}(1/\bar{n}_{0} \wedge \bar{n}_{1})$$

# 

## C PESUDOCODE

In this section, we include the pseudocode for COCOA 1.

Algorithm 1 Contrastive Counterfactual Augmentation

**Input:** Factual dataset  $D_{\rm F} = \{(x_i, t_i, y_i)\}_{i=1}^n$ ; sensitivity parameter  $\epsilon$ ; threshold k **Output:** Augmented factual dataset  $D_{\rm AF}$  as training dataset for CATE estimation models Initialize  $D_{\rm A} = \emptyset$ Construct datasets  $D_{\epsilon}^+$  and  $D_{\epsilon}^-$  from  $D_{\rm F}$ Optimize a parametric model  $g_{\theta}$  with contrastive learning and  $(D_{\epsilon}^+, D_{\epsilon}^-)$ for i = 1 to n do Determine  $N_i = \{(x_j, y_j) | j \in [n], t_j = 1 - t_i, g_{\theta}(x_i, x_j) = 1\}$ if  $|N_i| \ge k$  then Estimate  $\hat{y}_i$  with  $\psi(x_i, N_i)$ Add  $(x_i, 1 - t_i, \hat{y}_i)$  to  $D_A$ end if end for Set  $D_{\rm AF} = D_{\rm A} \cup D_{\rm F}$ 

# **D** DATASET DESCRIPTIONS

**IHDP** The IHDP dataset is a semi-synthetic dataset that was introduced based on real covariates available from the Infant Health and Development Program (IHDP) to study the effect of development programs on children. The features (covariates) in this dataset come from a Randomized Control Trial. The potential outcomes were simulated following Setting B in Hill [2011]. The IHDP dataset consists of 747 individuals (139 in the treatment group and 608 in the control group), each with 25 features. The potential outcomes are generated as follows:

$$Y_0 \sim \mathcal{N}(\exp(\beta^T (X+W)), 1)$$

and

$$Y_1 \sim \mathcal{N}(\beta^T (X + W) - \omega, 1)$$

where W has the same dimension as X with all entries equal 0.5 and  $\omega = 4$ . The regression coefficient  $\beta$  is a vector of length 25 where each element is randomly sampled from a categorical distribution with the support (0, 0.1, 0.2, 0.3, 0.4) and the respective probability masses  $\mu = (0.6, 0.1, 0.1, 0.1, 0.1)$ .

**News** The News Dataset is a semi-synthetic dataset designed to assess the causal effects of various news topics on reader responses. It was first introduced in Johansson et al. [2016]. The documents were sampled from news items from the NY Times corpus (downloaded from UCI Newman et al. [2008]). The covariates available for CATE estimation are the raw word counts for the 100 most probable words in each topic. The treatment  $t \in \{0, 1\}$  denotes the viewing device. t = 0

means with computer and t = 1 means with mobile. A topic model is trained on a comprehensive collection of documents to generate  $z(x) \in \mathbb{R}^k$  that represents the topic distribution of a given news item x [Johansson et al., 2016].

Let the treatment effects be represented by  $z_{c_1}$  (for t = 1) and  $z_{c_0}$  (for t = 0)  $z_{c_1}$  is defined as the topic distribution of a randomly selected document while  $z_{c_0}$  is the average topic representation across all documents. The reader's opinion of news item x on device t is influenced by the similarity between z(x) and  $z_{c_t}$ , expressed as:

$$y(x,t) = C \cdot \left( z(x)^T z_{c_0} + t \cdot z(x)^T z_{c_1} \right) + \epsilon$$

where C = 50 is a scaling factor and  $\epsilon \sim \mathcal{N}(0, 1)$ . The assignment of a news item x to a device  $t \in \{0, 1\}$  is biased towards the preferred device for that item, modeled using the softmax function:

$$p(t=1|x) = \frac{e^{\kappa \cdot z(x)^T z_{c_1}}}{e^{\kappa \cdot z(x)^T z_{c_0}} + e^{\kappa \cdot z(x)^T z_{c_1}}}$$

Here,  $\kappa$  determines the strength of the bias and it is assigned to be 10.

**Twins** The Twins dataset Louizos et al. [2017] is based on the collected birthday data of twins born in the United States from 1989 to 1991. It is assumed that twins share significant parts of their features. Consider the scenario where one of the twins was born heavier than the other as the treatment assignment. The outcome is whether the baby died in infancy (i.e., the outcome is mortality). Here, the twins are divided into two groups: the treatment and the control groups. The treatment group consists of heavier babies from the twins. On the other hand, the control group consists of lighter babies from the twins. The potential outcomes,  $Y_0$  and  $Y_1$ , are generated through:

$$Y_0 \sim \mathcal{N}(\exp(\beta^T X), 0.2)$$

and

$$Y_1 \sim \mathcal{N}(\alpha^T X, 0.2)$$

Where  $\beta$  and  $\alpha$  are sampled from a high dimensional standard normal distribution.

**Linear dataset** We synthetically generate a dataset with N = 1500 samples and d = 10 features. The feature vectors  $X = (x_1, x_2, \dots, x_d)^T \in \mathbb{R}^d$  are drawn from a standard normal distribution. The treatment assignment  $t \in \{0, 1\}$  is biased, with the probability of treatment being

$$p(t = 1|x) = \frac{1}{1 + \exp(-(x_1 + x_2))}$$

We generate potential outcomes using two linear functions with coefficients  $\beta_0 = (0.5, \dots, 0.5) \in \mathbb{R}^d$  and  $\beta_1 = (0.3, \dots, 0.3) \in \mathbb{R}^d$  as follows:

$$Y_0 = \beta_0 X + \mathcal{N}(0, 0.01)$$
  
$$Y_1 = \beta_1 X + \mathcal{N}(0, 0.01)$$

**Non-Linear dataset** We construct a synthetic dataset consisting of N = 1500 instances with d = 10 features. The feature vectors, denoted by  $X = (x_1, x_2, \dots, x_d)^T \in \mathbb{R}^d$ , are sampled from a standard normal distribution. The treatment assignment  $t \in \{0, 1\}$  is biased, with the probability of treatment being

$$p(t = 1|x) = \frac{1}{1 + \exp(-(x_1 + x_2))}$$

We generate potential outcomes using two linear functions with coefficients  $\beta_0 = (0.5, \dots, 0.5) \in \mathbb{R}^d$  and  $\beta_1 = (0.3, \dots, 0.3) \in \mathbb{R}^d$  as follows:

$$Y_0 = \exp(\beta_0 X) + \mathcal{N}(0, 0.01)$$
$$Y_1 = \exp((\beta_1 X) + \mathcal{N}(0, 0.01)$$

Table 3:  $\sqrt{\varepsilon_{\text{PEHE}}}$  across various CATE estimation models with and without COCOA augmentation on Linear and Non-Linear synthetic datasets. Lower  $\sqrt{\varepsilon_{\text{PEHE}}}$  corresponds to better performance.

	Lin	ear	Non-linear		
Model	w/o aug.	w/ aug.	w/o aug.	w/ aug.	
TARNet	$0.93 {\pm} .09$	$0.81 {\pm}.02$	$7.41 \pm .23$	$6.64 \pm .11$	
CFR-Wass	$0.87 {\pm} .05$	$0.74 {\pm} .05$	$7.32 {\pm}.21$	$6.22 \pm .07$	
CFR-MMD	$0.91 {\pm} .04$	$0.78 {\pm}.06$	$7.35 {\pm}.19$	$6.28 \pm .10$	
T-Learner	$0.90 {\pm}.01$	$0.89 {\pm .01}$	$7.68 {\pm}.12$	$7.51 {\pm}.07$	
S-Learner	$0.64 \pm .01$	$0.63 {\pm}.01$	$7.22 {\pm}.01$	$6.92 {\pm}.01$	
BART	$0.65 {\pm}.00$	$0.30 {\pm}.00$	$5.49 {\pm}.00$	$4.50 {\pm}.00$	
CF	$0.63 {\pm}.00$	$0.27 {\pm}.00$	$5.46 \pm .00$	$4.46 \pm .00$	

# **E** ADDITIONAL EMPIRICAL RESULTS

In this section, we present additional results for the completeness of the empirical study for COCOA. Specifically, we (i) add the results for the synthetic datasets, (ii) provide details for the toy example used to generate Figure 3b, (ii) present more visualizations illustrating the effect of contrastive learning, (iv) study the performance of our proposed method on ATE estimation, (v) conduct ablation studies on the local regression module, (vi) present additional results to demonstrate robustness against overfitting, and (vii) perform ablation studies on different parameters for the contrastive learning module.

#### E.1 RESULTS FOR SYNTHETIC DATA

In this section, we present the  $\sqrt{\varepsilon_{\text{PEHE}}}$  results for various CATE estimation models on synthetic datasets, both linear and nonlinear. Table 3 summarizes the performance of each model with COCOA augmentation (w/ aug.) and without augmentation (w/o aug.). Lower  $\sqrt{\varepsilon_{\text{PEHE}}}$  indicates better performance. The results demonstrate that COCOA augmentation consistently improves the performance across different models and datasets.

#### E.2 TRADE-OFF TOY EXAMPLE

In this section, we synthetically generate a dataset for a binary treatment scenario with 1000 samples per treatment group and d = 4 features. We sample a vector of coefficients,

$$\beta \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$$

where  $\mathbf{0} \in \mathbb{R}^d$  is the zero vector and  $\mathbf{I}_d$  is the  $d \times d$  identity matrix.

Next, we generate feature vectors  $X \in \mathbb{R}^d$  for the two treatment groups:

$$X_0 \sim \mathcal{N}(-1, 0.5\mathbf{I}_d)$$

and,

$$X_1 \sim \mathcal{N}(\mathbf{1}, 0.5\mathbf{I}_d)$$

where  $-\mathbf{1} \in \mathbb{R}^d$  and  $\mathbf{1} \in \mathbb{R}^d$  are vectors with all elements equal to -1 and 1, respectively, and  $\mathbf{I}_d$  is the  $d \times d$  identity matrix. The potential outcomes are generated as follows:

$$Y_0 = (\beta^T X_0)^3 + \mathcal{N}(0, 0.1)$$

and

$$Y_1 = (\beta^T X_1)^2 + \mathcal{N}(0, 0.1)$$

We implement a function to augment the datasets using a nearest-neighbor approach with a specified radius (radius is set to 8). The augmentation involves imputing potential outcomes for individuals from the opposite treatment group if they have at least three close neighbors within the specified radius. We then perform linear regression to impute the outcomes. We include further empirical results in Figure 4.



Figure 4: Trade-off between imputation error and statistical disparity. The first plot displays the percentage of augmentation as a function of the radius. The second and third plots show the Total Variation (TV) distance and imputation error, respectively, for different radius values.



Figure 5: Comparison between euclidean distance and latent distance lerned by contrastive learning for the IHDP dataset (treatment group). The first heatmap illustrates the outcome distances. The second heatmap shows the feature distances, reflecting differences between feature vectors. The third heatmap presents the embedding distances, demonstrating how the learned embeddings capture the same similarities as the potential outcome.

#### E.3 CONTRASTIVE LEARNING MOTIVATION

In this section, we provide more motivation for the use of contrastive learning to learn a representation space in which we identify similar individuals instead of using traditional methods (e.g., euclidean distance the ambient space). Figures 5 and 6 illustrate this. We also include an ablation on the effect of the embedding dimension for contrastive learning on the learned representation for the IHDP dataset as illustrated in Figure 7.

#### E.4 ATE ESTIMATION PERFORMANCE

In this section, we provide additional empirical results when applying our methods to ATE estimation. The Average Treatment Effect (ATE) is defined as:

$$\tau_{\text{ATE}} = \mathbb{E}[Y_1 - Y_0].$$

The error of ATE estimation is defined as:

$$\varepsilon_{\text{ATE}} = \left| \hat{\tau}_{\text{ATE}} - \tau_{\text{ATE}} \right|,\tag{17}$$

Our results are summarized in Tables 4, 5, and 6. We observe that our methods, while not tailored for ATE estimation, still bring some benefits for a subset of the estimation models.

#### E.5 LOCAL REGRESSION MODULE

In this section, we compare the performance of using Gaussian Processes (GP)with different kernels vs. local linear regression. We next define the local linear regression module and present the empirical results in Table 7.

**Local Linear Regression.** For a fixed individual x who received treatment t, and has a selected neighbors  $D_{x,t}$ . Under the assumption that we can locally approximate the true function with a linear function. Suppose  $X_D$  is the matrix of the



Figure 6: Comparison between euclidean distance and latent distance lerned by contrastive learning for the IHDP dataset (control group). The first heatmap illustrates the outcome distances. The second heatmap shows the feature distances, reflecting differences between feature vectors. The third heatmap presents the embedding distances, demonstrating how the learned embeddings capture the same similarities as the potential outcome.



Figure 7: t-SNE visualizations of the IHDP dataset control group embeddings for different embedding dimensions. The figure illustrates t-SNE plots for the control group with embedding dimensions of 2, 4, 6, and 8. The points are colored based on outcome groups, created by dividing the outcomes into four quantiles. Each subplot shows how the embeddings distribute in a 2D space, capturing the relationship between the learned embeddings and outcome groups. Outcome groups represent different quantile ranges of potential outcomes: Group 0 (yellow) includes the lowest quantile, Group 1 (cyan) includes the second lowest, Group 2 (blue) includes the second highest, and Group 3 (magenta) includes the highest quantile.

Table 4:  $\varepsilon_{ATE}$  across various CATE estimation models, with COCOA augmentation (w/ aug.) and without augmentation (w/o aug.) in Twins, Linear, and Non-Linear datasets. Lower  $\varepsilon_{ATE}$  corresponds to the better performance.

	TWINS		LINEAR		Non-linear	
MODEL	W/O AUG.	W/ AUG.	W/O AUG.	W/ AUG.	W/O AUG.	W/ AUG.
TARNET	$0.33 \pm .19$	$0.41 \pm .29$	$0.10 \pm .02$	$0.04 {\pm}.02$	$0.23 \pm .13$	$0.04 {\pm}.02$
CFR-WASS	$0.47 \pm .16$	$0.14 \pm .09$	$0.13 {\pm}.04$	$0.06 {\pm}.01$	$0.19 {\pm}.09$	$0.03 {\pm}.01$
CFR-MMD	$0.19 {\pm} .09$	$0.18 \pm .12$	$0.12 \pm .05$	$0.05 {\pm}.03$	$0.25 \pm .15$	$0.04 {\pm}.01$
T-LEARNER	$0.02 \pm .02$	$0.05 {\pm}.03$	$0.01 {\pm}.01$	$0.01 {\pm}.01$	$0.05{\pm}0.02$	$0.05 {\pm}.01$
S-LEARNER	$0.89 {\pm}.03$	$0.79 {\pm}.07$	$0.03 {\pm}.01$	$0.05 {\pm}.01$	$0.45 \pm .05$	$0.27 {\pm}.02$
BART	$0.28 \pm .08$	$0.21 \pm .10$	$0.37 {\pm}.00$	$0.07 {\pm}.01$	$0.80 {\pm}.00$	$0.26 {\pm}.00$
CF	$0.28 \pm .06$	$0.14 \pm .15$	$0.39 {\pm}.00$	$0.06 {\pm}.01$	$0.77 {\pm}.00$	$0.32 {\pm}.00$

Table 5:  $\varepsilon_{ATE}$  across various CATE estimation models, with COCOA augmentation (w/ aug.), without augmentation (w/o aug.), and with Perfect Match augmentation in News and IHDP datasets. Lower  $\varepsilon_{ATE}$  corresponds to the better performance.

	NEWS		IHI	)P
MODEL	W/O AUG.	W/ AUG.	W/O AUG.	W/ AUG.
TARNET	$0.97 {\pm}.45$	$0.96 \pm .38$	$0.12 \pm .05$	$0.07 \pm .03$
CFR-WASS	$1.00 \pm .29$	$0.75 \pm .22$	$0.10 \pm .03$	$0.05 {\pm}.02$
CFR-MMD	$0.89 \pm .38$	$0.71 \pm .22$	$0.16 {\pm}.04$	$0.09 {\pm} .04$
T-LEARNER (NN)	$0.49 \pm .26$	$0.76 \pm .20$	$0.27 {\pm}.06$	$0.07 {\pm} .03$
S-LEARNER (NN)	$0.40 {\pm}.06$	$0.49 \pm .27$	$1.72 \pm .21$	$0.40 \pm .02$
BART	$0.77 \pm .13$	$0.60 {\pm} .00$	$0.02 {\pm}.01$	$0.02 {\pm}.01$
CAUSAL FORESTS	$0.72 {\pm .01}$	$0.60 {\pm}.00$	$0.11 \pm .01$	$0.03 {\pm}.02$
PERFECT MATCH	$2.00{\pm}1.01$		$0.24 \pm .20$	

observed feature values in  $D_{x,t}$  augmented with a column of ones for the intercept, and  $Y_D$  is the column vector of observed factual outcomes. The local linear regression coefficients,  $\hat{\beta}$ , are computed as:

$$\hat{\beta} = (X_D^T X_D)^{-1} X_D^T Y_D$$

Then we impute the value of x as  $\hat{y} = [1, x]^T \hat{\beta}$ .

## E.6 ABLATION FOR CONTRASTIVE LEARNING PARAMETERS

In this section, we provide a comprehensive set of ablation studies for the effect of the hyper-parameters of the contrastive learning module.

Ablation on K and R. We provide extra ablation studies on the IHDP dataset and the Non-linear dataset to study the effect of (*i*) the number of neighbors (K) and (*ii*) the embedding radius (R) on both  $\varepsilon_{PEHE}$  and  $\varepsilon_{ATE}$ . We observe a consistently enhanced performance across different CATE estimation models. See results in figures 10 and 11. We also provide ablation studies on the sensitivity of the proposed Contrastive Learning module to the parameter  $\epsilon$ , which is used to create the training points for the contrastive learning module by creating positive and a negative dataset.

Ablation on the sensitivity parameter  $\epsilon$  We provide ablation on the sensitivity parameter  $\epsilon$ , a similarity classifier for the potential outcomes. The results for the  $\varepsilon_{PEHE}$  as a function of  $\epsilon$  are presented in Figure 8. It can be observed that the error of CATE estimation models is consistent for a wide range of  $\epsilon$ , demonstrating the robustness of COCOA to the choice of hyper-parameters.

Table 6:  $\varepsilon_{ATE}$  across different similarity measures: Contrastive Learning (CL), propensity scores (PS), and Euclidean distance (ED), using CFR-Wass across IHDP, News, and Twins datasets.

MEASURE OF SIMILARITY	ED	PS	CL
IHDP	$3.12 \pm 1.33$	$3.85 \pm .22$	$0.05 \pm .02$
NEWS	$0.68 \pm .20$	$0.54 \pm .25$	$0.75 \pm .22$
TWINS	$0.13 \pm .15$	$0.46 {\pm}.09$	$0.14 \pm .09$

Table 7: Comparison of  $\varepsilon_{\text{PEHE}}$  and  $\varepsilon_{\text{ATE}}$  across different local regression modules: Gaussian Process (GP) with various kernels (DotProduct, RBF, and Matern) and Linear Regression. The first three rows present  $\sqrt{\varepsilon_{\text{PEHE}}}$ , while the subsequent three rows display  $\varepsilon_{\text{ATE}}$ .

LR	<b>GP</b> (DOTPRODUCT)	GP (RBF)	GP (MATERN)	LINEAR REGRESSION
IHDP	<b>0.63</b> ±.01	$0.63 \pm .00$	$0.65 \pm .02$	$0.75 \pm .01$
NEWS	$3.56 {\pm}.01$	$3.55 {\pm}.04$	$3.44 \pm .05$	$3.53 {\pm}.08$
TWINS	$0.51 \pm .11$	$0.51 \pm .02$	$0.54 \pm .04$	$0.68 \pm .08$
IHDP	$0.02 \pm .01$	<b>0.01</b> ±.00	$0.03 \pm .01$	$0.09 {\pm}.01$
NEWS	$0.60 \pm .00$	$0.24 \pm .12$	$0.05 \pm .03$	$0.21 {\pm}.10$
TWINS	<b>0.21</b> ±.10	$0.24 {\pm}.04$	$0.29 {\pm}.04$	$0.38 \pm .10$

## E.7 OVERFITTING TO THE FACTUAL DISTRIBUTION

In this section, we provide more empirical results on the robustness against overfitting to the factual distribution for the Linear and Non-Linear synthetic datasets, as presented in Figure 9.



Figure 8:  $\varepsilon_{\text{PEHE}}$  as a function of the similarity sensitivity parameter  $\epsilon$ . The figure on the left presents results for the IHDP dataset, while the one on the right is for the News dataset. Performances of two different models (CFR-Wass and Causal Forests) are plotted for both datasets.



Figure 9: Comparison of training progression for TARNet, CFR-Wass, and T-learner models on linear and non-linear datasets. Top row: Models trained on the linear dataset, showcasing TARNet, CFR-Wass, and T-learner, respectively. Bottom row: The same models trained on the non-linear dataset. This visualization demonstrates the effects of COCOA on preventing overfitting across different data complexities and the performance of three CATE estimation models trained with various levels of data augmentation.



Figure 10: Ablation studies on the impact of the size of the  $\epsilon$ -Ball (R) and the number of neighbors (K) on the performance. The first row from left to right: IHDP with TARNet, BART, S-Learner, and Causal Forests. The second row: IHDP with Causal Forests, T-Learner, BART, and TARNet. These studies illustrate the trade-off between minimizing the discrepancy between the distributions—achieved by reducing K and increasing R—and the quality of the imputed data points, which is achieved by decreasing R and increasing K.



Figure 11: Ablation studies on the Non-linear dataset. Top row from left to right: Causal Forests (PEHE), BART (PEHE), TARNet (PEHE). Bottom row from left to right: Causal Forests (ATE), BART (ATE), TARNet (ATE). Each pair of images represents the performance of the respective models evaluated in terms of Precision in Estimation of Heterogeneous Effect (PEHE) and the error in Average Treatment Effect (ATE) estimation on a non-linear dataset.