000 001 002

003

004

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

# REVISITING COVARIATE AND HYPOTHESIS ROLES IN ITE ESTIMATION: A NEW APPROACH USING LAPLA-CIAN REGULARIZATION

Anonymous authors

Paper under double-blind review

#### Abstract

The recent surge in data availability across many fields, such as medicine, social science, and marketing, has brought to the forefront the problem of estimating Individual Treatment Effect (ITE) from observational data to effectively tailor treatment to personalized characteristics. ITE estimation is known to be a challenging task because we can only observe the outcome with or without treatment, but never both. Moreover, observational datasets exhibit selection bias induced by the treatment assignment policy. In this paper, we present a new approach consisting of two novel aspects. First, we depart from conventional approaches that minimize the covariate shift. Instead, we incorporate it as a crucial element in ITE estimation, recognizing that it stems from highly predictive features that exhibit significant imbalance in observational data. Second, unlike existing methods, our approach utilizes hypothesis functions to directly estimate outcomes under covariate shift, enhancing reliability across observed and unobserved outcomes. To support this approach theoretically, we derive a new upper bound of the expected ITE loss and show that it explicitly depends on the discrepancy between the hypothesis functions, which are absent from the objectives of existing methods. Based on this new approach, we present LITE: Laplacian Individual Treatment Effect, a novel method that leverages Laplacian-regularized representation and incorporates both the covariate shift and the hypothesis functions for ITE estimation, effectively bridging observed and unobserved outcomes. We demonstrate LITE on illustrative simulations and two leading benchmarks, where we show superior results compared to state-of-the-art methods.

032 033 034

#### 1 INTRODUCTION

Individual Treatment Effect (ITE) has come to the forefront of precision medicine (Prosperi et al., 2020; Glass et al., 2013), targeted marketing (Lemmens & Gupta, 2020), personalized education (Beemer et al., 2018), and various other fields (Wang et al., 2016) that require individual-level predictions. ITE specifically aims to quantify the unique outcome of an action (also referred to as a treatment or intervention) for each individual based on their specific characteristics, which is a departure from traditional methods that focus on average treatment effect (Abadie & Imbens, 2016).

The focus on individual outcomes is critical, especially since individuals are unlikely to conform to average behavioral patterns (Schork, 2015). Customizing treatments based on unique characteristics is therefore essential for achieving both effective and efficient interventions. In this approach, an individual is characterized by their features, with typically two possible outcomes considered: the outcome with the applied treatment and the outcome in the absence of the treatment. The objective is to estimate the difference between these outcomes (i.e., the treatment effect) based on the individual's features, enabling to customize treatment policies for that individual.

Accurate estimation of ITE is a challenging task, mainly because only one potential outcome for
 each individual is observable based on the applied action (i.e., the factual outcome). Inferring the
 unobserved outcome (i.e., the counterfactual outcome) from the outcomes observed in other indi viduals often leads to poor estimates due to selection bias (Vokinger et al., 2021) (also referred to
 as confounding bias). This bias stems from the assignment policy to treatment or control groups,

- typical in studies where randomized controlled trials (RCTs) are unavailable due to budget considerations, difficulty in recruiting patients, or ethical constraints. For example, suppose a medication is primarily given to patients with severe symptoms. Using the observed data from the control group —those not receiving the medication—to predict the counterfactual outcome of the treated group,
  i.e., their predicted outcome of lack of treatment, would result in overly optimistic estimates. Conversely, using the treated group to predict the outcome of treatment in the control group would likely result in pessimistic estimates. Furthermore, in observational data, we often lack direct insight into the mechanisms (i.e., the confounding variables) to infer the treatment policy.
- 062 Traditional methods approximate ITE by identifying nearest neighbors using matching techniques 063 to estimate counterfactuals (Ho et al., 2007; Gu & Rosenbaum, 1993; Dehejia & Wahba, 2002; 064 Schwab et al., 2018). Tree-based methods (Chipman et al., 2010; Green & Kern, 2012; Lu et al., 2018; Athey & Imbens, 2016; Wager & Athey, 2018) view forests as an adaptive neighborhood 065 metric and estimate treatment effects at the leaf nodes (Wager & Athey, 2018). Other approaches 066 use Gaussian processes for ITE estimation (Alaa & Van Der Schaar, 2017; Alaa & Schaar, 2018). 067 Representation learning has become central in ITE estimation by harnessing the power of latent 068 representations (Bengio et al., 2013). To tackle selection bias, these methods (Johansson et al., 069 2016; Shalit et al., 2017; Yao et al., 2018; Yoon et al., 2018; Guo et al., 2023; Du et al., 2021; Johansson et al., 2022) often aim to minimize covariate shift by balancing covariate representations 071 across groups using distance metric regularization (Johansson et al., 2016; Shalit et al., 2017; Yao 072 et al., 2018; Guo et al., 2023) or adversarial methods (Yoon et al., 2018; Du et al., 2021), ensuring 073 that counterfactual predictions are guided by the most reliable aspects of the data (Johansson et al., 074 2016). However, we assert that unlike classical domain adaptation, where covariate shift is treated 075 as an artifact to be mitigated (e.g., blurred vs. clear images), in ITE estimation, the covariate shift is inherent and directly impacts the causal treatment effect. Directly minimizing covariate shift can 076 inadvertently reduce predictive components that are essential for understanding the causal treatment 077 effect (Yoon et al., 2018; Yao et al., 2018; Du et al., 2021), and thus produce biased ITE estimate even in the limit of infinite data (Johansson et al., 2018). This is because the selection policy is typically 079 applied based on highly predictive features that often show significant imbalance, as doctors, for example, usually assign treatments according to these features. 081

This paper introduces a new approach that revisits the role of covariates and hypothesis functions in ITE estimation. We present a new upper bound of the ITE estimation error that shifts focus from 083 the covariate shift to the hypothesis function discrepancies. Following this result, unlike traditional 084 methods that minimize the covariate shift, our method, termed LITE (Laplacian Individual Treat-085 ment Effect), goes beyond this by integrating both covariate and hypothesis considerations into the ITE estimation process. To this end, we construct a graph within the latent learned representation 087 space, capturing the covariate shift and serving as our model's foundational structure. Utilizing 088 this graph, we compute the graph Laplacian, which in turn, is used for regularizing the hypothesis 089 functions to allow estimate ITE directly under the shift. Specifically, we use the graph Laplacian 090 quadratic form to facilitate the smoothness of the hypothesis functions with the geometry of the 091 learned representation across predicted factual and counterfactual outcomes. We demonstrate LITE 092 on a simulation and two leading benchmarks. We show that LITE outperforms both established and recent methods by a large margin, achieving state of the art results.

094 095

096

098

099

102

103

105

### Our main contributions are as follows:

- **Theoretical foundation:** We present a new upper bound of the ITE estimation error that redirects the focus from the covariate shift to the discrepancies between hypothesis functions.
- LITE: We present a new method for ITE estimation that unlike existing methods considers both the covariate shift and the hypothesis function discrepancies.
- **Geometry-aware representation:** We utilize Laplacian regularization not only to align the learned representation with the hypothesis function outcomes, but also to dynamically learn the graph structure itself through the optimization process.
- State of the art results: LITE demonstrates state of the art results on leading benchmarks.
- Broader impact: Our method is easily extended to handle multiple treatment scenarios, thus enhancing practicality in many fields, and specifically, in healthcare applications, where multiple treatments are often available.

### 2 RELATED WORK

109

110 Laplacian regularization has been widely used in various domains (Pang & Cheung, 2017; Liu et al., 111 2018; Ziko et al., 2020), and more particularly in semi-supervised learning, as demonstrated by 112 Belkin et al. (2006); Cabannes et al. (2021); Calder et al. (2023). However, to the best of our 113 knowledge, its application within the latent space for ITE estimation is novel. Traditional Laplacian regularization is often applied in the input space, where geometry can be obscured by non-relevant 114 features. Our approach contrasts with this by promoting geometry-aware structures within the la-115 tent space, focusing on predictive features relevant to the task. Furthermore, by incorporating this 116 Laplacian regularization in the learning objective, the latent representation and the resulting graph 117 are continuously refined through optimization, systematically aligning the learned representation 118 geometry with the intrinsic geometry of the underlying data structure. 119

120 While covariate balancing helps reduce the impact of selection bias, it might inadvertently reduce intrinsic differences that are important for accurate ITE estimation (Yoon et al., 2018; Johansson 121 et al., 2018; Yao et al., 2018; Du et al., 2021). There exist representation learning methods for ITE 122 estimation that aim to mitigate the adverse nature of covariate balancing. For instance, Yao et al. 123 (2018) proposed local similarity-preserved representations to prevent the potential loss of local sim-124 ilarity information during distribution balancing. Alternatively, Johansson et al. (2018) combined 125 re-weighting methods to alleviate predictive information loss. Du et al. (2021) employed mutual in-126 formation regularization to retain information that is highly predictive of the outcome. Despite such 127 mitigation strategies, these methods still apply direct minimization while regularized, which can 128 reduce the crucial intrinsic differences necessary for accurate ITE estimation. Moreover, these ap-129 proaches often neglect the critical role of hypothesis functions in addressing selection bias, focusing 130 instead solely on covariate balancing at the expense of the predictive capacity of these functions.

131

133

#### **3** INDIVIDUAL TREATMENT EFFECT BACKGROUND

134 3.1 PROBLEM FORMULATION

136 We adopt the framework of potential outcomes Pearl (2009), as originally formulated in Rubin (1974); Rosenbaum & Rubin (1983), to analyze the ITE. We follow the notations from Shalit et al. 137 (2017). The ITE problem aims to learn the potential effect of a treatment t based on individual 138 features x (also referred to as covariates). Within this framework, let  $\mathcal{X}$  be the input space,  $\mathcal{T}$  be the 139 action (treatment) space, and  $\mathcal{Y}$  be the outcome space. An individual is characterized by features 140  $x \in \mathcal{X}$ . In this paper, to simplify the presentation, we assume a binary treatment setting where each 141 individual is assigned an action  $t \in \mathcal{T} = \{0, 1\}$ , where 0 denotes the absence of treatment and 1 142 indicates the presence of treatment. The probability of assigning treatment, given a set of covariates, 143 is defined as the propensity score  $\pi(x) = p(t = 1|x)$  (Rosenbaum & Rubin, 1983), and reflects the 144 treatment assignment policy. 145

For each individual, two potential outcomes exist:  $Y_0$  without treatment and  $Y_1$  with treatment. However, in this setting, for each individual, only one of these potential outcomes, either  $Y_0$  or  $Y_1$ , is observed via y, depending on the applied treatment action:

$$y = \begin{cases} Y_0 & \text{if } t = 0\\ Y_1 & \text{if } t = 1 \end{cases}.$$

**Definition 1** (Individual treatment effect). The ITE  $\tau(x)$  (also termed the Conditional Average Treatment Effect) is defined as the expected difference in potential outcomes for an individual x:

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

The interest in ITE estimation lies in learning the function  $\tau(x)$ , which quantifies the difference between the potential outcomes  $Y_1$  and  $Y_0$  for a given individual with features x. Unlike traditional learning scenarios, this ITE  $\tau(x)$  is not directly observable in the training data. Specifically, for each individual, we only observe one potential outcome, dictated by the applied treatment action.

Definition 2 (PEHE). The expected Precision in Estimation of Heterogeneous Effect (PEHE) (Hill, 2011) loss is defined as:

149

150 151

152

153

154

$$\epsilon_{PEHE} = \int_{\mathcal{X}} (\hat{\tau}(x) - \tau(x))^2 p(x) dx, \qquad (2)$$

164 Minimizing PEHE enhances the accuracy of ITE estimates, which is crucial for developing targeted 165 interventions that are optimally tailored to individual characteristics. Our objective is to estimate 166 the function  $\tau(x)$  in a manner consistent with causal inference principles (Sauerbrei et al., 2014; 167 Cousens et al., 2011), using an observational dataset  $\mathcal{D}$ . This dataset typically consists of *n* inde-168 pendent samples of features, action, and outcome, represented by the tuple  $(x_i, t_i, y_i)$ .

<sup>169</sup> Two assumptions are commonly considered in ITE estimation.

**Assumption 1** (Positivity). There exists a positive probability of receiving any treatment action, conditional on the individual features. Formally:  $\forall t \in \{0,1\}, \forall x \in \mathcal{X}, \quad 0 < P(t|x) < 1.$ 

This assumption ensures an overlap between the control and treatment groups. Violation of this positivity assumption implies that for some x, we lack any observation of one of the potential outcomes, making the counterfactual outcome estimation even more challenging.

**Assumption 2** (Ignorability). *The treatment assignment is conditionally independent of the potential outcomes, given the observed features. Formally:*  $\{Y_0, Y_1\} \perp t | x$ .

This condition is crucial for the identifiability of the ITE. The ignorability assumption (often referred to as "unconfoundedness") ensures that there are no hidden confounders affecting both treatment assignment and potential outcomes.

181 182 183

191

199

201

211 212

213

#### 3.2 THE SELECTION BIAS

Our objective is to estimate the ITE  $\tau(x)$  in Eq. equation 1 by learning two separate functions:  $m_1(x) = \mathbb{E}[Y_1|x]$  and  $m_0(x) = \mathbb{E}[Y_0|x]$ . To this end, access to the joint distribution function  $p(x, t, Y_0, Y_1)$ , defined on the input-action-potential outcome space, is essential. However, we have a sample with the factual outcome:  $(x_1, t_1, y_1), ..., (x_n, t_n, y_n)$ , where  $y_i \sim p(Y_1|x_i)$  if  $t_i = 1$ , and  $y_i \sim p(Y_0|x_i)$  if  $t_i = 0$ . This provides just a partial view of the joint distribution.

**Definition 3** (ITE estimate). For an individual x, let  $f : \mathcal{X} \times \mathcal{T} \to \mathcal{Y}$  be a function that maps individual and treatment pairs to outcomes. The ITE estimate of the hypothesis f is defined by:

$$\hat{\tau}(x) = f(x,1) - f(x,0),$$
(3)

where f(x, 1) and f(x, 0) are the estimates of  $m_1(x)$  and  $m_0(x)$ , respectively.

One may suggest learning two separate functions: f(x, 1) from individuals who received the treatment and f(x, 0) from those who did not. However, these functions are susceptible to the selection bias in the observational dataset (Vokinger et al., 2021). More specifically, f(x, 1) and f(x, 0) stem from different distributions, i.e., the treated and control distributions p(x|t = 1) and p(x|t = 0), respectively, and these distributions are shifted relative to the covariate marginal p(x).

Let  $L: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}^+$  be a loss function. To analyze the loss under the selection bias effect, we define the factual and counterfactual domains,  $p_F^t(Y_t, x) \triangleq p(Y_t, x|t)$  and  $p_{CF}^t(Y_t, x) \triangleq p(Y_t, x|1-t)$ , respectively, where both distributions are conditioned on the treatment assignment.

In the subsequent definitions, we denote the expected loss for an individual and treatment pair (x, t)as  $\ell_f(x,t) = \int_{\mathcal{Y}} L(Y_t, f(x,t)) p(Y_t|x) dY_t$  and the proportion of treated individuals as u = p(t = 1). We apply the ignorability assumption, implying that  $p(Y_t|x,t) = p(Y_t|x,1-t) = p(Y_t|x)$ , and obtain that the potential outcomes are conditionally independent of the treatment, given the covariates.

**Definition 4** (Factual loss). *The expected factual loss for a treatment assignment t is defined by:* 

$$\epsilon_F^t(f) = \int_{\mathcal{X} \times \mathcal{Y}} L(Y_t, f(x, t)) p_F^t(Y_t, x) dx dY_t = \int_{\mathcal{X}} \ell_f(x, t) p(x|t) dx$$

The factual loss (across all treatment assignments) is then: 215

$$\epsilon_F(f) = u\epsilon_F^{t=1}(f) + (1-u)\epsilon_F^{t=0}(f).$$
(4)

where  $\hat{\tau}(x)$  is the estimate of the ITE and p(x) is the p.d.f. over the input space  $\mathcal{X}$ .

216 Definition 5 (Counterfactual loss). The expected counterfactual loss for a treatment assignment t is
 217 defined by:
 218

 $\epsilon_{CF}^t(f) = \int_{\mathcal{X} \times \mathcal{Y}} L(Y_t, f(x, t)) p_{CF}^t(Y_t, x) dx dY_t = \int_{\mathcal{X}} \ell_f(x, t) p(x|1-t) dx$ 

The counterfactual loss (across all treatment assignments) is then:

$$\epsilon_{CF}(f) = (1 - u)\epsilon_{CF}^{t=1}(f) + u\epsilon_{CF}^{t=0}(f).$$
(5)

The following theorem, presented in Shalit et al. (2017), links the ITE estimation with the counterfactual error.

**Theorem 1.** The Precision in Estimation of Heterogeneous Effects (PEHE) is bounded by:

$$\rho_{EHE}(f) \le 2(\epsilon_{CF}(f) + \epsilon_F(f) - \sigma_V^2), \tag{6}$$

where  $\epsilon_F(f)$  and  $\epsilon_{CF}(f)$  denote the factual and counterfactual losses with respect to the squared loss, respectively, and  $\sigma_Y^2$  is the variance of the outcome variable  $Y_t$ .

This theorem highlights how ITE estimation is tightly linked to counterfactual errors, emphasizing the critical role of understanding and addressing these errors for accurate ITE estimation. An intuitive strategy involves minimizing errors with respect to the factual outcomes, as these are directly observable and quantifiable. While this empirical approach might seem effective at first glance, it overlooks the importance of counterfactual outcomes, crucial for ITE estimation. Consequently, while this intuitive strategy may perform well in terms of factual error, its performance is often suboptimal for ITE estimation. For more details, we refer the readers to Shalit et al. (2017).

238 239 240

241

219 220 221

222 223

228

229

230

231 232

233

234

235

236

237

## 4 PROPOSED APPROACH

To address the challenges associated with estimating ITE, we employ a representation-learning framework (Bengio et al., 2013), which is widely used in the ITE literature (Johansson et al., 2016; Shalit et al., 2017; Guo et al., 2023; Yao et al., 2018; Shi et al., 2019). This framework utilizes a representation function  $\Phi : \mathcal{X} \to \mathcal{R}$  that maps observed features x into a latent space  $\mathcal{R}$  via  $\Phi(x)$ . In this latent space, we use a hypothesis function h with two possible second arguments:  $h(\Phi(x), t = 0)$  and  $h(\Phi(x), t = 1)$ , for predicting the potential outcomes  $Y_0$  and  $Y_1$ , respectively. Our general hypothesis function is then expressed as  $f(x, t) = h(\Phi(x), t)$ .

To effectively capture the complex relationships between observed features x and potential out-249 comes, we parameterize  $\Phi(x)$  and  $h(\Phi,t)$  using deep neural networks. Specifically,  $h(\Phi,t)$  is 250 realized by two separate fully connected networks (Caron et al., 2022), one for each treatment 251  $t \in \{0, 1\}$ . In addition,  $\Phi$  is realized using fully connected layers. This parameterization choice lever-252 ages the advanced capabilities of neural networks to capture intricate patterns in the data. Moreover, 253 employing a shared representation  $\Phi(x)$  across treatment assignments leverages the commonalities 254 among treated and control groups, enhancing generalization and efficiency. While we use fully 255 connected layers to realize  $\Phi(x)$ ,  $h(\Phi, t = 0)$ , and  $h(\Phi, t = 1)$ , our approach is flexible and can 256 accommodate alternative network architectures to potentially enhance model performance.

257

Covariate balancing. Recent methods (Johansson et al., 2016; Shalit et al., 2017; Yao et al., 2018;
 Yoon et al., 2018; Guo et al., 2023; Du et al., 2021) emphasize covariate balancing within the latent representation space through direct covariate shift minimization to mitigate counterfactual error, an approach underpinned by established theoretical frameworks. This strategy is illustrated by the following theorem from Shalit et al. (2017).

Let G be a function family  $g: S \to \mathbb{R}$ . For a pair of distributions  $p_1, p_2$  over S, define the Integral Probability Metric (IPM) as follows:  $IPM_G(p_1, p_2) = \sup_{g \in G} \int_S g(s)(p_1(s) - p_2(s)) ds$ 

**Theorem 2** (PEHE upper bound by covariate discrepancy). Assuming a representation function  $\Phi: \mathcal{X} \to \mathcal{R}$ , and a hypothesis function  $h: \mathcal{R} \times \{0, 1\} \to \mathcal{Y}$ , the PEHE  $\epsilon_{PEHE}(h, \Phi)$  can be bounded by the IPM distance between the distribution of treated and control groups:

268  
269
$$\epsilon_{PEHE}(h,\Phi) \le 2(\epsilon_F^{t=0}(h,\Phi) + \epsilon_F^{t=1}(h,\Phi) + B_{\Phi} \cdot IPM_G(p_{\Phi}^{t=1}, p_{\Phi}^{t=0}) - 2\sigma_Y^2), \tag{7}$$

where  $p_{\Phi}$  is the distribution induced by  $\Phi$  over  $\mathcal{R}$  and  $B_{\Phi}$  is a constant bounding the loss functions.

270 This theorem supports an objective function that reduces the discrepancy between treated and control 271 group distributions in the latent space, aiming to avoid reliance on potentially unreliable data aspects 272 when generalizing from factual to counterfactual domains (Johansson et al., 2016). This discrepancy 273 is often quantified using the Wasserstein distance (Villani et al., 2009; Cuturi & Doucet, 2014) or 274 through adversarial methods (Yoon et al., 2018; Du et al., 2021), among other metrics (Yao et al., 2018; Guo et al., 2023; Gretton et al., 2012). Formally, the training objective integrates not only 275 the factual outcome errors but also a term to account for unobserved counterfactual outcomes by 276 reducing the divergence between the distributions in the latent space: 277

278

$$\mathcal{O}(\theta) = \epsilon_F(h, \Phi) + \alpha \cdot \operatorname{IPM}_G(p_{\Phi}^{t=1}, p_{\Phi}^{t=0}), \tag{8}$$

where  $\alpha$  is a hyperparameter balancing factual accuracy against distributional balance.

While this classical theorem provides valuable insights for ITE estimation, we assert that its direct application in minimizing the covariate shift might mitigate the significant impact of covariates on outcomes. In medical settings, for example, treatments are assigned based on predictive features, introducing the selection bias (Yoon et al., 2018). Therefore, merely reducing these discrepancies without recognizing their contributions to the causal structure might lead to models that misrepresent treatment effects and result in suboptimal outcomes.

Hypothesis balancing. We present a theorem that shifts the focus from traditional covariate shifts to discrepancies within hypothesis functions for assessing counterfactual error. This theorem underpins our approach by integrating considerations of both covariate and hypothesis function discrepancies. For proof and further details see Appendix A.

**Theorem 3** (PEHE upper bound by hypothesis function discrepancy). Let  $\Phi : \mathcal{X} \to \mathcal{R}$  be an invertible representation with  $\Psi$  its inverse, and let  $h : \mathcal{R} \times \{0,1\} \to \mathcal{Y}$  be a hypothesis function. Recall that  $\ell_f(x,t) = \int_{\mathcal{Y}} L(Y_t, f(x,t)) p(Y_t|x) dY_t$  is the expected loss for an individual and treatment pair (x,t). The PEHE  $\epsilon_{PEHE}(h, \Phi)$  is bounded by discrepancies within the hypothesis functions:

$$\epsilon_{PEHE}(h,\Phi) \le 2\left(2\epsilon_F(h,\Phi) + \int_{\mathcal{R}} \left|\left(\ell_{h,\Phi}(\Psi(r),1) - \ell_{h,\Phi}(\Psi(r),0)\right)\right| dr - \sigma_Y^2\right),\tag{9}$$

299 The term  $\int_{\mathcal{R}} |(\ell_{h,\Phi}(\Psi(r),1) - \ell_{h,\Phi}(\Psi(r),0))| dr$  reflects the difference between the expected loss 300 predictions for both potential outcomes relative to the learned representation, and is governed by the 301 hypothesis functions. This difference is primarily induced by the selection bias, as this term may en-302 compass estimation errors arising in sparsely-sampled regions where the unobserved counterfactual labels are insufficiently represented by the estimated functions, which rely on factual labels. While 303 selection bias is the principal issue, our ITE estimation scheme intentionally goes beyond covariate 304 shift minimization by integrating the hypothesis functions to estimate outcomes directly under the 305 shift. This approach advocates for incorporating hypothesis functions in ITE estimation instead of 306 solely relying on covariate shift minimization. We propose a way to minimize this difference by 307 demanding smoothness of the hypothesis functions with respect to the learned representation, in 308 regions with counterfactual relevance. This allows the model to infer counterfactual from factual 309 samples and, thus, to effectively reduce the difference between prediction losses, and consequently, 310 the ITE estimation error, as suggested by the new bound.

311

296 297 298

# 4.1 LITE: LAPLACIAN INDIVIDUAL TREATMENT EFFECT

We present LITE: Laplacian Individual Treatment Effect, a method that integrates covariate and hypothesis function discrepancies through a regularized-Laplacian representation. To this end, we construct a graph within the latent space that captures covariate shift and requires the functions to maintain smoothness over the defined graph geometry by regularizing the hypothesis functions relative to the geometry of the learned representation.

Consider an observational dataset  $\mathcal{D} = \{(\boldsymbol{x}_i, t_i, y_i)\}_{i=1}^N$ , where  $\boldsymbol{x}_i \in \mathbb{R}^d$  are the features, d denotes the number of features,  $t_i$  is the treatment indicator, and  $y_i$  is the factual outcome, all of the *i*th sample. In our deep learning framework, this dataset is batch-processed through a neural network to obtain a representation function  $\Phi(\boldsymbol{x})$  into a latent representation space  $\mathcal{R}$ . This representation is then input into the hypothesis function h, split into the two branches of treatment assignment  $h(\Phi, t = 0)$  and  $h(\Phi, t = 1)$ , which compute the respective potential outcomes. In the latent space, we construct a graph where each node corresponds to one sample from the batch, and the edges represent the affinities between these samples, quantified by a radial basis function (RBF). The adjacency matrix  $A \in \mathbb{R}^{b \times b}$  of the graph, where *b* denotes the batch size, is defined by:

$$\boldsymbol{A}_{ij} = \exp\left(-\frac{\|\Phi(\boldsymbol{x}_i) - \Phi(\boldsymbol{x}_j)\|^2}{2\sigma^2}\right),\tag{10}$$

where each entry  $A_{ij}$  in the matrix represents the weight of the edge between nodes *i* and *j*, capturing the strength of interaction based on the similarity in their latent representations. Here,  $\Phi(x) \in \mathbb{R}^r$  is the latent representation of each sample, where *r* denotes the latent dimension, and  $\sigma$ is a scale parameter set to the mean of the pairwise distances, up to some factor. For further details on the kernel type, scale, and distance metric, see Appendix B.1.

Using this adjacency matrix, we construct the Laplacian  $\mathcal{L}$  as follows:

$$\mathcal{L} = D - A, \tag{11}$$

where D is the diagonal matrix with  $D_{ii} = \sum_{j} A_{ij}$ . Then, our objective function is expressed as:

$$\mathcal{O}(\theta) = \epsilon_F(h, \Phi) + \alpha \cdot S_{\text{LITE}}(h, \Phi), \tag{12}$$

where the LITE term  $S_{\text{LITE}}(h, \Phi)$  is given by:

327 328

335

336 337

342 343 344

345 346 347

348 349

350

351

352

353 354

374 375

$$S_{\text{LITE}}(h,\Phi) = \frac{1}{b^2} \left( \mathbf{h}_0^T \mathcal{L} \mathbf{h}_0 + \mathbf{h}_1^T \mathcal{L} \mathbf{h}_1 \right),$$
(13)

where  $\mathbf{h}_t = [h(\Phi(\boldsymbol{x}_1), t), \dots, h(\Phi(\boldsymbol{x}_b), t)]^T$  for  $t \in \{0, 1\}$ , and the factual error  $\epsilon_F(h, \Phi)$  is:

$$\epsilon_F(h,\Phi) = \frac{1}{b} \sum_{i=1}^{b} L(h(\Phi(\boldsymbol{x}_i), t = t_i), y_i).$$
(14)

LITE is summarized in Algorithm 1. It is described for binary treatment for simplicity, however, our approach is flexible and can easily be extended to accommodate multiple treatments. See Appendix C. LITE is designed to handle large datasets efficiently and is implemented to allow for fast computation. For further details on scalability and computation time, see Appendix D.

Algorithm 1 The LITE Algorithm

Ι	<b>nit:</b> Initialize network parameters $\theta$	
1: 2: 3:	Feed batch from $\mathcal{D}$ into $\Phi$ to obtain $\Phi(\boldsymbol{x})$ Compute $h(\Phi(\boldsymbol{x}), t = 0)$ and $h(\Phi(\boldsymbol{x}), t = 1)$ Construct adjacency graph $\boldsymbol{A}$ Build the Lanlacian operator $\boldsymbol{C}$	<ul> <li>Latent representation</li> <li>Potential outcome prediction</li> <li>According to Eq. (10)</li> <li>According to Eq. (11)</li> </ul>
ч. 5:	Compute $S_{\text{LITE}}$	▷ According to Eq. (13
6:	Compute factual error $\epsilon_F(h, \Phi)$	$\triangleright$ According to Eq. (14)
7:	Calculate the objective:	
	$\mathcal{O}(\theta) = \epsilon_F(h, \Phi) + \alpha \cdot S_{LI}$	$_{ ext{TE}}(h,\Phi)$
8: e F	Update $\theta$ and validate on a hold-out set <b>nd while</b> Return $\theta$ with lowest validation objective value	

$${}^{T}\mathcal{L}\mathbf{h}_{t} = \sum_{i,j} \boldsymbol{A}_{ij} (h_{t}[i] - h_{t}[j])^{2}, \qquad (15)$$

where  $A_{ij}$  is defined in equation 10, and  $h_t[k]$  is the potential outcome of sample k with treatment assignment t. Minimizing this expression enforces predictions  $h_t[k]$  and  $h_t[l]$  are similar when samples k and l are close according to the geometry in the learned representation space defined by  $A_{ij}$ . Broadly, the adjacencies in A (based on  $||\Phi(x_i) - \Phi(x_j)||$ ) carry information on the covariate shift, while  $|h_t[i] - h_t[j]|$  induces the hypothesis function smoothness. The quadratic form minimization promotes predictions that live in the span of smooth eigenvectors of the Laplacian graph. Our proposed approach leverages available factual labels, ensuring that model predictions are grounded in observations, while counterfactuals are inferred from them under this smoothness requirement.

We now clarify the connection between Theorem 3 and the LITE method, which employs Laplacian 384 regularization, highlighting how the regularization enforces smoothness across samples and reduces 385 the difference between prediction losses for both potential outcomes. Theorem 3 establishes an 386 upper bound on the ITE error based on differences in expected loss predictions for both potential 387 outcomes. The graph Laplacian regularization enforces smoothness across samples by minimizing 388 the quadratic form of both factual and counterfactual samples in the latent space. We argue that smoothness across outcomes of different samples leads to reducing the difference between pre-389 diction losses. In particular, in regions without selection bias, both potential outcome predictions 390 would ideally have similar levels of predicted losses, and thus, the difference would be small. How-391 ever, due to the selection bias, and particularly in regions with severe selection bias, one potential 392 outcome prediction function is based on factual samples, while the other potential outcome pre-393 diction function is insufficiently represented by unobserved counterfactual outcomes. This leads to 394 higher predicted losses and larger differences between the predicted losses. The Laplacian addresses 395 the unobserved regions by enforcing smoothness across outcomes of different samples, allowing 396 the model to infer counterfactual from factual samples and, thus, effectively reducing the difference 397 between prediction losses. 398

We note that while our method, LITE, effectively utilizes Theorem 3, the framework of our Theorem is general and can accommodate other alternatives incorporating function handling rather than covariate minimization. Following this approach, we show that LITE achieves significant empirical improvements over existing methods.

402 403 404

405

412

#### 5 EXPERIMENTS

We evaluate the performance of LITE through a simulation and two leading benchmarks and compare it to both recent and established methods. While both potential outcomes are available for evaluation, in all the experiments, the optimization process does not involve counterfactual labels. The source code will be made available on GitHub upon acceptance. For additional details on the experimental setup, hyperparameter configurations, and other supplementary information, see Appendix B.

413 **Performance Metrics.** We report the ITE estimation error,  $\epsilon_{PEHE}$ \_  $\frac{1}{n}\sum_{i=1}^{n} \left( (h(x_i, 1) - h(x_i, 0)) - (m_1(x_i) - m_0(x_i)) \right)^2, \text{ and the absolute error in the estimated average treatment effect, } \epsilon_{ATE} = \left| \frac{1}{n} \sum_{i=1}^{n_1} (h(x_i, 1) - h(x_i, 0)) - \frac{1}{n} \sum_{i=1}^{n} (m_1(x_i) - m_0(x_i)) \right|,$ 414 415 416 where  $m_i = \mathbb{E}[Y_i|x]$ . We consider two evaluation tasks. (i) *In-sample* evaluation is conducted on 417 subsets used during optimization, including training and validation sets. Unlike traditional machine 418 learning scenarios, this task is challenging because counterfactual labels are not available during the 419 optimization process. (ii) Out-of-sample evaluation is conducted on the test set and involves unseen 420 data, where neither factual nor counterfactual labels were available during optimization.

421 422

423

### 5.1 Illustrative simulation

424 We demonstrate the impact of LITE in mitigating the counterfactual error in a simulation. To this 425 end, we generate a synthetic dataset with 600 individuals, each characterized by a scalar covariate 426 x within the range (-2.5, 2.5). The simulated potential outcomes  $m_1(x)$  and  $m_0(x)$  are depicted 427 as dashed lines in Fig. 1(a). The treatment assignment is biased and modeled by the propensity score function:  $\pi(x) = p(t = 1|x) = (1 + e^{-(-0.1 + 0.9x)})^{-1}$ . This selection bias is shown in Fig. 428 1(a), where treated individuals (blue dots) are more likely to have higher x values, while control 429 individuals (green dots) are more likely to have lower x values, resulting in a skewed distribution 430 across x. For each individual, we only observe the factual outcome according to the treatment 431 assignment policy, while the unobserved counterfactual outcomes are depicted as gray dots. In

441

442

443

444

445

446

447 448



Figure 1: Illustrative example of ITE estimation using LITE. (a) The expected outcomes  $m_0$  and  $m_1$  (shown as dashed lines) with factual labels for treated and control (in blue and green dots, respectively) and with unobserved counterfactual labels (in gray). (b) The corresponding ITE function (dotted line) with the unobserved ITE samples. (c) Factual, counterfactual, and ITE error (PEHE) for varying values of LITE regularization coefficient. (d) Predicted outcomes for the best  $\alpha$  value in terms of PEHE, compared to without using LITE ( $\alpha = 0$ ). (e) The corresponding ITE estimation.

449 Fig. 1(b), the ITE function  $m_1 - m_0$  is shown, with the unobserved ITE samples representing the 450 difference between factual and counterfactual outcomes.

451 In Fig. 1(c), we observe the factual, counterfactual, and PEHE errors for different values of the 452 regularization term  $\alpha$ . At  $\alpha = 0$ , optimization is based solely on factual error without the LITE 453 term. As  $\alpha$  increases, we see a significant reduction in counterfactual error and consequently a 454 smaller ITE error (PEHE). These results are averaged over 100 realizations. See more details in 455 Appendix B. In Fig. 1(d), the predicted outcomes for one realization are presented. We see that 456 without LITE ( $\alpha = 0$ ), the model struggles to capture the ground truth in regions with severe 457 selection bias due to the influence of a few factual labels, creating inconsistent trends in small 458 sample regimes. In contrast, with LITE ( $\alpha = 0.1$ ), while not informed by counterfactual labels, 459 the model ensures consistency in those regions by aligning counterfactual predictions with factual labels. This refinement lowers the predicted slope in these regions, enabling the model to capture 460 the true trend more accurately. In Fig. 1(e), we present the corresponding ITE estimation, further 461 demonstrating the effectiveness of LITE. 462

463 464

465

# 5.2 IHDP AND NEWS BENCHMARKS

**IHDP.** The IHDP dataset Hill (2011) is arguably the most used benchmark for ITE estimation. It 466 examines the impact of specialist home visits on cognitive test scores and consists of 747 units (608 467 control, 139 treated) with 25 covariates related to the children and their mothers. These features and 468 treatment assignments were extracted from a real-world clinical trial, with selection bias introduced 469 by selectively removing a subset of the patients. We report  $\sqrt{\epsilon_{\text{PEHE}}}$  and  $\epsilon_{\text{ATE}}$  for both in-sample and 470 out-of-sample evaluations. We compare our method to 20 other methods, including both established 471 and recent state-of-the-art models, whose performance has been reported in the literature.

472 Tab. 1 presents the results. We see that our method achieves the best ITE estimation, demonstrating 473 a significant margin compared to state-of-the-art methods. While obtaining the best ITE estimation, 474 our method achieves the second-best results in terms of Average Treatment Effect (ATE), which are 475 comparable to the best results obtained by ABCEI. Yet, ABCEI yields much inferior ITE estimation.

476 477 478

479

**News.** The News dataset comprises 5,000 New York Times articles, each represented by a 3,477word vocabulary, analyzed for consumer opinions on different devices. For more details, see Appendix B.1. The results are presented in Tab. 2. We see that our method achieves the best results in terms of both ITE estimation and ATE by a large margin.

484

#### 6 CONCLUSION, LIMITATIONS, AND FUTURE WORK

In this paper, we presented a new approach for ITE estimation that diverges from traditional meth-485 ods by considering, rather than minimizing, covariate shifts, as well as discrepancies between hy-

486	
487	Table 1: Results on the IHDP dataset (1000 iterations). The best results are in bold, and the second-
488	best results are underlined. Values are as reported in the literature (see Table 5 in the appendix for
489	the references). 'n.r.' denotes values not reported.

			$\sqrt{\epsilon_{\text{PEHE}}}$		$\epsilon_{ m ATE}$	
Group		Method	In-sample	Out-of-sample	In-sample	Out-of-sample
Classic ML regression		OLS <sub>1</sub> OLS <sub>2</sub>	$\begin{array}{c} 5.8\pm0.3\\ 2.4\pm0.1\end{array}$	$\begin{array}{c} 5.8\pm0.3\\ 2.5\pm0.1\end{array}$	$\begin{array}{c} 0.73 \pm 0.04 \\ 0.14 \pm 0.01 \end{array}$	$\begin{array}{c} 0.94 \pm 0.06 \\ 0.31 \pm 0.02 \end{array}$
Matching		k-NN PSM PM	$\begin{array}{c} 2.1 \pm 0.1 \\ 4.92 \pm 0.312 \\ \text{n.r.} \end{array}$	$\begin{array}{c} 4.1 \pm 0.2 \\ 4.92 \pm 0.312 \\ 0.84 \pm 0.61 \end{array}$	0.14 ± 0.01 n.r. n.r.	$\begin{array}{c} 0.79 \pm 0.05 \\ 0.78 \pm 0.03 \\ 0.24 \pm 0.20 \end{array}$
Tree-based		BART R. For. C. For.	$\begin{array}{c} 2.1 \pm 0.1 \\ 4.2 \pm 0.2 \\ 3.8 \pm 0.2 \end{array}$	$\begin{array}{c} 2.3 \pm 0.1 \\ 6.6 \pm 0.3 \\ 3.8 \pm 0.2 \end{array}$	$\begin{array}{c} 0.23 \pm 0.01 \\ 0.73 \pm 0.05 \\ 0.18 \pm 0.01 \end{array}$	$\begin{array}{c} 0.34 \pm 0.02 \\ 0.96 \pm 0.06 \\ 0.40 \pm 0.03 \end{array}$
Gaussian processes		CMGP NSGP	$\begin{array}{c} 0.61 \pm 0.011 \\ 0.51 \pm 0.013 \end{array}$	$\begin{array}{c} 0.76 \pm 0.012 \\ 0.64 \pm 0.030 \end{array}$	$\frac{0.11\pm0.10}{\text{n.r.}}$	$\begin{array}{c} 0.13 \pm 0.12 \\ 0.23 \pm 0.01 \end{array}$
	General	TARNET CEVAE	$\begin{array}{c} 0.88\pm0.02\\ 2.7\pm0.1\end{array}$	$\begin{array}{c} 0.95 \pm 0.02 \\ 2.6 \pm 0.1 \end{array}$	$\begin{array}{c} 0.26 \pm 0.01 \\ 0.34 \pm 0.01 \end{array}$	$\begin{array}{c} 0.28 \pm 0.01 \\ 0.46 \pm 0.02 \end{array}$
Representation learning	Balanced	BLR BNN CFR MMD CFR WASS SITE MitNet GANITE ABCEI	$5.8 \pm 0.3$ 2.2 ± 0.1 0.73 ± 0.01 0.71 ± 0.02 0.69 ± 0.0 n.r. 1.9 ± 0.4 0.71 ± 0.0	$5.8 \pm 0.3 \\ 2.1 \pm 0.1 \\ 0.78 \pm 0.02 \\ 0.76 \pm 0.02 \\ 0.75 \pm 0.0 \\ 0.60 \pm 0.03 \\ 2.4 \pm 0.4 \\ 0.73 \pm 0.0$	$\begin{array}{c} 0.72 \pm 0.04 \\ 0.37 \pm 0.03 \\ 0.30 \pm 0.01 \\ 0.25 \pm 0.01 \\ 0.22 \pm 0.01 \\ \text{n.r.} \\ 0.43 \pm 0.05 \\ \textbf{0.09 \pm 0.01} \end{array}$	$\begin{array}{c} 0.93 \pm 0.05 \\ 0.42 \pm 0.03 \\ 0.31 \pm 0.01 \\ 0.27 \pm 0.01 \\ 0.24 \pm 0.01 \\ 0.25 \pm 0.01 \\ 0.49 \pm 0.05 \\ \textbf{0.09} \pm \textbf{0.01} \end{array}$
	Geometric	LITE (Our Method)	$\textbf{0.35} \pm \textbf{0.004}$	$\textbf{0.37} \pm \textbf{0.005}$	$\underline{0.11\pm0.003}$	$\underline{0.12\pm0.003}$

Table 2: Results on the News dataset (50 iterations). See Table 6 in the appendix for the references.

			$\sqrt{\epsilon_{\text{PEHE}}}$		$\epsilon_{\mathrm{ATE}}$	
Group		Method	In-sample	Out-of-sample	In-sample	Out-of-sample
Classic ML regression		LASSO <sub>1</sub> LASSO <sub>2</sub>	$\begin{array}{c} 4.23 \pm 0.17 \\ 2.03 \pm 0.08 \end{array}$	$\begin{array}{c} 4.25 \pm 0.17 \\ 2.31 \pm 0.16 \end{array}$	$\begin{array}{c} 2.5 \pm 0.07 \\ 0.33 \pm 0.02 \end{array}$	$\begin{array}{c} 2.5 \pm 0.07 \\ 0.34 \pm 0.03 \end{array}$
Gaussian processes		CMGP	n.r.	$2.21\pm0.05$	n.r.	n.r.
	General	TARNet CEVAE	$\begin{array}{c} 1.81 \pm 0.05 \\ \text{n.r.} \end{array}$	$\begin{array}{c} 1.93 \pm 0.06 \\ 3.74 \pm 0.18 \end{array}$	$\begin{array}{c} 0.32\pm0.04\\ \text{n.r.} \end{array}$	0.30 ± 0.04 n.r.
Representation learning	Balanced	CFR WASS SITE ABCEI NOFELITE	$\begin{array}{c} 1.83 \pm 0.05 \\ 2.20 \pm 0.07 \\ \underline{1.63 \pm 0.05} \\ \text{n.r.} \end{array}$	$\begin{array}{c} 1.98 \pm 0.06 \\ 2.44 \pm 0.09 \\ \underline{1.81 \pm 0.07} \\ 2.18 \pm 0.05 \end{array}$	$\begin{array}{c} 0.34 \pm 0.04 \\ \underline{0.18 \pm 0.02} \\ \underline{0.18 \pm 0.03} \\ n.r. \end{array}$	$\begin{array}{c} 0.37 \pm 0.04 \\ \underline{0.22 \pm 0.03} \\ 0.23 \pm 0.04 \\ \text{n.r.} \end{array}$
	Geometric	LITE (Our Method)	$\textbf{1.24} \pm \textbf{0.04}$	$\textbf{1.44} \pm \textbf{0.05}$	$\textbf{0.15} \pm \textbf{0.02}$	$\textbf{0.16} \pm \textbf{0.02}$

pothesis functions. Building on this approach, we proposed LITE (Laplacian Individual Treatment Effect), a new method that incorporates both covariate shift and hypothesis function into a Laplacian-regularized representation. We showed that LITE outperforms established and recent SOTA methods on leading benchmarks.

LITE is primarily demonstrated on binary treatment frameworks for simplicity. As described in the
 paper, it is effective for discrete treatment categories and has been explicitly extended to handle
 multiple treatment conditions by generalizing the Laplacian term for multi-treatment scenarios.

However, many practical applications, especially in healthcare, involve treatments with continuous variables, e.g. dosage levels (Schwab et al., 2020). The treatment effect may depend on the amount of a drug administered, which requires understanding the dose-response relationship to optimize outcomes. Future work will explore integrating into LITE methodologies that handle continuous treatment variables. This extension will involve developing new geometric embeddings or adapting the existing regularization to accommodate a continuous treatment space. Such advancements could enhance the model's utility in precision medicine by enabling nuanced analyses of optimal dosing strategies.

# 540 7 ETHICS STATEMENT

541 542

Our research adheres to the ICLR Code of Ethics, ensuring responsible research conduct and commitment to high scientific standards. In this research, we exclusively use publicly available datasets for individual treatment effect estimation. Our work does not involve human subjects, personal data, or sensitive information. We are committed to transparency, reproducibility, and the ethical use of machine learning techniques. Upon acceptance, we will make the source code available.

This statement does not count towards the page limit as per ICLR guidelines.

548 549

547

# 8 REPRODUCIBILITY STATEMENT

550 551

To ensure the reproducibility of our research, detailed explanations of the methodologies and experimental settings are provided. The full set of assumptions and complete proofs for all theoretical results are documented in Appendix A.The LITE algorithm is thoroughly summarized in Algorithm 1 in the main text. Commonly used evaluation criteria are clearly specified within the paper to ensure clarity and adherence to standard practices. All experimental details, including data processing steps, model configurations, dataset splitting, optimization methods, and hyperparameters, are comprehensively described in Appendix B.

We are committed to transparency and will make the source code publicly available upon acceptance,
allowing other researchers to replicate and verify our results using the same methodologies and data.
The experimental setup and computing resourced are disclosed in Appendix B.1. For additional
details on the scalability and computation time of our LITE algorithm, see Appendix D.

- 563 This statement does not count towards the page limit as per ICLR guidelines.
- 564 565

566

573

574

575

576

583

584

585

References

- Alberto Abadie and Guido W Imbens. Matching on the estimated propensity score. *Econometrica*, 84(2):781–807, 2016.
- Ahmed Alaa and Mihaela Schaar. Limits of estimating heterogeneous treatment effects: Guidelines for practical algorithm design. In *International Conference on Machine Learning*, pp. 129–138.
  PMLR, 2018.
  - Ahmed M Alaa and Mihaela Van Der Schaar. Bayesian inference of individualized treatment effects using multi-task gaussian processes. *Advances in neural information processing systems*, 30, 2017.
- Susan Athey and Guido Imbens. Recursive partitioning for heterogeneous causal effects. *Proceed- ings of the National Academy of Sciences*, 113(27):7353–7360, 2016.
- Joshua Beemer, Kelly Spoon, Juanjuan Fan, Jeanne Stronach, James P Frazee, Andrew J Bohonak, and Richard A Levine. Assessing instructional modalities: Individualized treatment effects for personalized learning. *Journal of Statistics Education*, 26(1):31–39, 2018.
  - Mikhail Belkin, Partha Niyogi, and Vikas Sindhwani. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of machine learning research*, 7 (11), 2006.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- Leo Breiman. Random forests. *Machine learning*, 45:5–32, 2001.
- Vivien Cabannes, Loucas Pillaud-Vivien, Francis Bach, and Alessandro Rudi. Overcoming the curse of dimensionality with laplacian regularization in semi-supervised learning. *Advances in Neural Information Processing Systems*, 34:30439–30451, 2021.

594 Jeff Calder, Dejan Slepčev, and Matthew Thorpe. Rates of convergence for laplacian semi-595 supervised learning with low labeling rates. *Research in the Mathematical Sciences*, 10(1):10, 596 2023. 597 Alberto Caron, Gianluca Baio, and Ioanna Manolopoulou. Estimating individual treatment effects 598 using non-parametric regression models: A review. Journal of the Royal Statistical Society Series A: Statistics in Society, 185(3):1115–1149, 2022. 600 601 Peipei Chen, Wei Dong, Xudong Lu, Uzay Kaymak, Kunlun He, and Zhengxing Huang. Deep rep-602 resentation learning for individualized treatment effect estimation using electronic health records. 603 Journal of biomedical informatics, 100:103303, 2019. 604 Hugh Chipman and Robert McCulloch. Bayestree: Bayesian additive regression trees. R package 605 version 0.3-1.4, 7, 2016. 606 607 Hugh A Chipman, Edward I George, and Robert E McCulloch. Bart: Bayesian additive regression 608 trees. 2010. 609 Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network 610 learning by exponential linear units (elus). arXiv preprint arXiv:1511.07289, 2015. 611 612 Simon Cousens, J Hargreaves, Chris Bonell, B Armstrong, J Thomas, BR Kirkwood, and R Hayes. 613 Alternatives to randomisation in the evaluation of public-health interventions: statistical analysis and causal inference. Journal of Epidemiology & Community Health, 65(7):576-581, 2011. 614 615 Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. Advances in neural 616 information processing systems, 26, 2013. 617 Marco Cuturi and Arnaud Doucet. Fast computation of wasserstein barycenters. In International 618 conference on machine learning, pp. 685–693. PMLR, 2014. 619 620 Rajeev H Dehejia and Sadek Wahba. Propensity score-matching methods for nonexperimental 621 causal studies. Review of Economics and statistics, 84(1):151-161, 2002. 622 623 Vincent Dorie. Npci: Non-parametrics for causal inference. URL: https://github. com/vdorie/npci, 624 11:23, 2016. 625 Xin Du, Lei Sun, Wouter Duivesteijn, Alexander Nikolaev, and Mykola Pechenizkiy. Adversarial 626 balancing-based representation learning for causal effect inference with observational data. Data 627 Mining and Knowledge Discovery, 35(4):1713–1738, 2021. 628 629 Thomas A Glass, Steven N Goodman, Miguel A Hernán, and Jonathan M Samet. Causal inference in public health. Annual review of public health, 34:61–75, 2013. 630 631 Donald P Green and Holger L Kern. Modeling heterogeneous treatment effects in survey experi-632 ments with bayesian additive regression trees. Public opinion quarterly, 76(3):491-511, 2012. 633 634 Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. The Journal of Machine Learning Research, 13(1):723–773, 2012. 635 636 Xing Sam Gu and Paul R Rosenbaum. Comparison of multivariate matching methods: Structures, 637 distances, and algorithms. Journal of Computational and Graphical Statistics, 2(4):405-420, 638 1993. 639 Xingzhuo Guo, Yuchen Zhang, Jianmin Wang, and Mingsheng Long. Estimating heterogeneous 640 treatment effects: mutual information bounds and learning algorithms. In International Confer-641 ence on Machine Learning, pp. 12108-12121. PMLR, 2023. 642 643 Jennifer L Hill. Bayesian nonparametric modeling for causal inference. Journal of Computational 644 and Graphical Statistics, 20(1):217-240, 2011. 645 Daniel E Ho, Kosuke Imai, Gary King, and Elizabeth A Stuart. Matching as nonparametric prepro-646 cessing for reducing model dependence in parametric causal inference. Political analysis, 15(3): 647 199-236, 2007.

648 Fredrik Johansson, Uri Shalit, and David Sontag. Learning representations for counterfactual infer-649 ence. In International conference on machine learning, pp. 3020–3029. PMLR, 2016. 650 Fredrik D Johansson, Nathan Kallus, Uri Shalit, and David Sontag. Learning weighted representa-651 tions for generalization across designs. arXiv preprint arXiv:1802.08598, 2018. 652 653 Fredrik D Johansson, Uri Shalit, Nathan Kallus, and David Sontag. Generalization bounds and rep-654 resentation learning for estimation of potential outcomes and causal effects. Journal of Machine 655 Learning Research, 23(166):1-50, 2022. 656 657 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint 658 arXiv:1412.6980, 2014. 659 Aurélie Lemmens and Sunil Gupta. Managing churn to maximize profits. Marketing Science, 39 660 (5):956-973, 2020. 661 662 Weifeng Liu, Xueqi Ma, Yicong Zhou, Dapeng Tao, and Jun Cheng. p-laplacian regularization for 663 scene recognition. IEEE Transactions on Cybernetics, 49(8):2927–2940, 2018. 664 Christos Louizos, Uri Shalit, Joris M Mooij, David Sontag, Richard Zemel, and Max Welling. 665 Causal effect inference with deep latent-variable models. Advances in neural information pro-666 cessing systems, 30, 2017. 667 668 Min Lu, Saad Sadiq, Daniel J Feaster, and Hemant Ishwaran. Estimating individual treatment effect 669 in observational data using random forest methods. Journal of Computational and Graphical 670 Statistics, 27(1):209-219, 2018. 671 Jiahao Pang and Gene Cheung. Graph laplacian regularization for image denoising: Analysis in the 672 continuous domain. IEEE Transactions on Image Processing, 26(4):1770–1785, 2017. 673 674 Judea Pearl. Causal inference in statistics: An overview. 2009. 675 676 Lutz Prechelt. Early stopping-but when? In Neural Networks: Tricks of the trade, pp. 55-69. 677 Springer, 2002. 678 Mattia Prosperi, Yi Guo, Matt Sperrin, James S Koopman, Jae S Min, Xing He, Shannan Rich, 679 Mo Wang, Iain E Buchan, and Jiang Bian. Causal inference and counterfactual prediction in 680 machine learning for actionable healthcare. Nature Machine Intelligence, 2(7):369–375, 2020. 681 682 Abbavaram Gowtham Reddy and Vineeth N Balasubramanian. Nester: An adaptive neurosymbolic 683 method for causal effect estimation. In Proceedings of the AAAI Conference on Artificial Intelli-684 gence, volume 38, pp. 14793-14801, 2024. 685 Paul R Rosenbaum and Donald B Rubin. The central role of the propensity score in observational 686 studies for causal effects. Biometrika, 70(1):41-55, 1983. 687 688 Donald B Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies. 689 *Journal of educational Psychology*, 66(5):688, 1974. 690 691 Willi Sauerbrei, Michal Abrahamowicz, Douglas G Altman, Saskia le Cessie, James Carpenter, 692 and STRATOS initiative. Strengthening analytical thinking for observational studies: the stratos initiative. Statistics in medicine, 33(30):5413-5432, 2014. 693 694 Nicholas J Schork. Personalized medicine: time for one-person trials. *Nature*, 520(7549):609–611, 695 2015. 696 697 Stefan Schrod, Fabian Sinz, and Michael Altenbuchinger. Adversarial distribution balancing for counterfactual reasoning. arXiv preprint arXiv:2311.16616, 2023. 699 Patrick Schwab, Lorenz Linhardt, and Walter Karlen. Perfect match: A simple method for 700 learning representations for counterfactual inference with neural networks. arXiv preprint arXiv:1810.00656, 2018.

702 703 704	Patrick Schwab, Lorenz Linhardt, Stefan Bauer, Joachim M Buhmann, and Walter Karlen. Learning counterfactual representations for estimating individual dose-response curves. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 34, pp. 5612–5619, 2020.
705 706 707 708	Uri Shalit, Fredrik D Johansson, and David Sontag. Estimating individual treatment effect: general- ization bounds and algorithms. In <i>International conference on machine learning</i> , pp. 3076–3085. PMLR, 2017.
709 710	Claudia Shi, David Blei, and Victor Veitch. Adapting neural networks for the estimation of treatment effects. <i>Advances in neural information processing systems</i> , 32, 2019.
712 713	Toon Vanderschueren, Jeroen Berrevoets, and Wouter Verbeke. Noflite: Learning to predict individ- ual treatment effect distributions. <i>Transactions on Machine Learning Research</i> , 2023.
714	Cédric Villani et al. Optimal transport: old and new, volume 338. Springer, 2009.
715 716 717	Kerstin N Vokinger, Stefan Feuerriegel, and Aaron S Kesselheim. Mitigating bias in machine learn- ing for medicine. <i>Communications medicine</i> , 1(1):25, 2021.
718 719	Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using random forests. <i>Journal of the American Statistical Association</i> , 113(523):1228–1242, 2018.
721 722 723	Yan Wang, Itai Kloog, Brent A Coull, Anna Kosheleva, Antonella Zanobetti, and Joel D Schwartz. Estimating causal effects of long-term pm2. 5 exposure on mortality in new jersey. <i>Environmental health perspectives</i> , 124(8):1182–1188, 2016.
724 725 726	Liuyi Yao, Sheng Li, Yaliang Li, Mengdi Huai, Jing Gao, and Aidong Zhang. Representation learning for treatment effect estimation from observational data. <i>Advances in neural information processing systems</i> , 31, 2018.
727 728 729 730	Jinsung Yoon, James Jordon, and Mihaela Van Der Schaar. Ganite: Estimation of individualized treatment effects using generative adversarial nets. In <i>International conference on learning representations</i> , 2018.
731 732 733	Imtiaz Ziko, Jose Dolz, Eric Granger, and Ismail Ben Ayed. Laplacian regularized few-shot learning. In <i>International conference on machine learning</i> , pp. 11660–11670. PMLR, 2020.
734	
735	
736	
737	
738	
739	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
754	
755	

# A PROOF OF THEOREM 3

**Definition 6.** Let  $p^{t=1}(x) := p(x \mid t = 1)$ , and  $p^{t=0}(x) := p(x \mid t = 0)$  denote respectively the treatment and control distributions.

**760 Definition 7.** For a representation function  $\Phi : \mathcal{X} \to \mathcal{R}$ , and for a distribution p defined over  $\mathcal{X}$ , let **761**  $p_{\Phi}$  be the distribution induced by  $\Phi$  over  $\mathcal{R}$ . Define  $p_{\Phi}^{t=1}(r) := p_{\Phi}(r \mid t = 1), p_{\Phi}^{t=0}(r) := p_{\Phi}(r \mid t = 0)$ , to be the treatment and control distributions induced over  $\mathcal{R}$ .

**Definition 8.** Let  $\Phi : \mathcal{X} \to \mathcal{R}$  be a representation function. Let  $h : \mathcal{R} \times \{0, 1\} \to \mathcal{Y}$  be a hypothesis defined over the representation space  $\mathcal{R}$ . The expected loss for the unit and treatment pair (x, t) is:

$$\ell_{h,\Phi}(x,t) = \int_{\mathcal{Y}} L(Y_t, h(\Phi(x),t)) p(Y_t \mid x) dY_t$$

**Assumption 3.** The representation function  $\Phi$  is one-to-one. Without loss of generality, we will assume that  $\mathcal{R}$  is the image of  $\mathcal{X}$  under  $\Phi$ , and define  $\Psi : \mathcal{R} \to \mathcal{X}$  to be the inverse of  $\Phi$ , such that  $\Psi(\Phi(x)) = x$  for all  $x \in \mathcal{X}$ .

**Theorem 3** (PEHE upper bound by hypothesis function discrepancy). Let  $\Phi : \mathcal{X} \to \mathcal{R}$  be an invertible representation with  $\Psi$  its inverse, and a hypothesis function  $h : \mathcal{R} \times \{0,1\} \to \mathcal{Y}$ , the PEHE  $\epsilon_{PEHE}(h, \Phi)$  can be bounded by discrepancies within the hypothesis functions:

=

$$\epsilon_{PEHE}(h,\Phi) \le 2\left(2\epsilon_F(h,\Phi) + \int_{\mathcal{R}} \left| \left(\ell_{h,\Phi}(\Psi(r),1) - \ell_{h,\Phi}(\Psi(r),0)\right) \right| dr - \sigma_Y^2 \right), \quad (16)$$

776 Proof.

$$\epsilon_{CF}(h,\Phi) - \epsilon_F(h,\Phi) = \left[ (1-u) \cdot \epsilon_{CF}^{t=1}(h,\Phi) + u \cdot \epsilon_{CF}^{t=0}(h,\Phi) \right]$$

$$- \left[ (1-u) \cdot \epsilon_F^{t=0}(h,\Phi) + u \cdot \epsilon_F^{t=1}(h,\Phi) \right]$$
(17)

$$(1-u) \cdot \left[\epsilon_{CF}^{t=1}(h,\Phi) - \epsilon_{F}^{t=0}(h,\Phi)\right]$$
(18)

+ 
$$u \cdot \left[\epsilon_{CF}^{t=0}(h,\Phi) - \epsilon_{F}^{t=1}(h,\Phi)\right]$$

$$= (1-u) \int_{\mathcal{X}} p^{t=0}(x) \left( \ell_{h,\Phi}(x,1) - \ell_{h,\Phi}(x,0) \right) dx$$
(19)  
+  $u \int_{\mathcal{X}} p^{t=1}(x) \left( \ell_{h,\Phi}(x,0) - \ell_{h,\Phi}(x,1) \right) dx$ 

Figure 1787
Equality in 17 and 19 is by Definitions 5 and, 4. Then by changing of variables using Definition 7 we can get:

$$\epsilon_{CF}(h,\Phi) - \epsilon_{F}(h,\Phi) = (1-u) \int_{\mathcal{R}} p_{\Phi}^{t=0}(r) \left( \ell_{h,\Phi}(\Psi(r),1) - \ell_{h,\Phi}(\Psi(r),0) \right) dr$$
(20)

$$+ u \int_{\mathcal{R}} p_{\Phi}^{t=1}(r) \left( \ell_{h,\Phi}(\Psi(r), 0) - \ell_{h,\Phi}(\Psi(r), 1) \right) dr$$

$$\leq (1 - e^{1/2} \int_{\mathcal{R}} e^{t=0}(r) \left( \ell_{h,\Phi}(\Psi(r), 1) - \ell_{h,\Phi}(\Psi(r), 0) \right) dr$$

$$\leq (1-u) \left| \int_{\mathcal{R}} p_{\Phi}^{t=0}(r) \left( \ell_{h,\Phi}(\Psi(r), 1) - \ell_{h,\Phi}(\Psi(r), 0) \right) dr \right|$$

$$+ u \left| \int p_{\Phi}^{t=1}(r) \left( \ell_{h,\Phi}(\Psi(r), 0) - \ell_{h,\Phi}(\Psi(r), 1) \right) dr \right|$$
(21)

$$+ u | \int_{\mathcal{R}} p_{\Phi}^{t}(r) (\ell_{h,\Phi}(\Psi(r), 0) - \ell_{h,\Phi}(\Psi(r), 1)) dr |$$

$$\leq (1 - u) \int_{\mathcal{R}} p_{\Phi}^{t=0}(r) | (\ell_{h,\Phi}(\Psi(r), 1) - \ell_{h,\Phi}(\Psi(r), 0)) | dr$$

$$+ u \int_{\mathcal{R}} n^{t=1}(r) | (\ell_{h,\Phi}(\Psi(r), 0) - \ell_{h,\Phi}(\Psi(r), 0)) | dr$$

$$+ u \int_{\mathcal{R}} n^{t=1}(r) | (\ell_{h,\Phi}(\Psi(r), 0) - \ell_{h,\Phi}(\Psi(r), 0)) | dr$$

$$+ u \int_{\mathcal{R}} n^{t=1}(r) | (\ell_{h,\Phi}(\Psi(r), 0) - \ell_{h,\Phi}(\Psi(r), 0)) | dr$$

$$+ u \int_{\mathcal{R}} n^{t=1}(r) | (\ell_{h,\Phi}(\Psi(r), 0) - \ell_{h,\Phi}(\Psi(r), 0)) | dr$$

$$+ u \int_{\mathcal{R}} n^{t=1}(r) | (\ell_{h,\Phi}(\Psi(r), 0) - \ell_{h,\Phi}(\Psi(r), 0)) | dr$$

803  
804  
805  

$$+ u \int_{\mathcal{R}} p_{\Phi}^{-}(r) |(\ell_{h,\Phi}(\Psi(r), 0) - \ell_{h,\Phi}(\Psi(r), 1))| dr$$

$$\leq (1-u) \int_{\mathcal{R}} |(\ell_{h,\Phi}(\Psi(r), 1) - \ell_{h,\Phi}(\Psi(r), 0))| dr$$
(23)

$$+ u \int_{\mathcal{P}} \left| \left( \ell_{h,\Phi}(\Psi(r),0) - \ell_{h,\Phi}(\Psi(r),1) \right) \right| dr$$

808  
809 
$$= \int_{\mathcal{R}}^{\mathcal{T}} \left| \left( \ell_{h,\Phi}(\Psi(r), 1) - \ell_{h,\Phi}(\Psi(r), 0) \right) \right| dr$$
(24)

The transition in 21 utilizes the triangle inequality for integration, while 22 employs the Cauchy-Schwarz inequality. In 23, we use the fact that the probability terms  $p_{\Phi}^{t}(r)$  are less than or equal to one. Thus we get:

$$\epsilon_{CF}(h,\Phi) \le \epsilon_F(h,\Phi) + \int_{\mathcal{R}} \left| \left( \ell_{h,\Phi}(\Psi(r),1) - \ell_{h,\Phi}(\Psi(r),0) \right) \right| dr$$

and combining this result into equation 1, we get:

$$\epsilon_{\text{PEHE}}(h,\Phi) \le 2(2\epsilon_F(h,\Phi) + \int_{\mathcal{R}} \left| \left( \ell_{h,\Phi}(\Psi(r),1) - \ell_{h,\Phi}(\Psi(r),0) \right) \right| dr - \sigma_Y^2 \right).$$
(25)

We note that the inequalities could be applied directly in the input space in Eq. 19 rather than in the representation space, meaning that the smoothness could also be employed in the input space. We assert that the input space may be obscured by irrelevant aspects of the data, whereas the representation space consists of relevant features, making it more suitable the smoothness requirement.

#### **B** EXPERIMENTS

#### B.1 MORE DETAILS ON THE EXPERIMENTAL SETTINGS

833 Network architecture. Our representation network  $\Phi(x)$  and the hypothesis networks  $h(\Phi, t = 0)$ 834 and  $h(\Phi, t = 1)$  are realized using fully connected layers with ELU activation functions (Clevert 835 et al., 2015). While our framework accommodates more complex architectures, the current imple-836 mentation achieves robust ITE estimation, demonstrating its effectiveness even when compared to 837 other methods (Yoon et al., 2018; Du et al., 2021; Louizos et al., 2017) that employ more sophis-838 ticated architectures. Following Shalit et al. (2017), we normalize the representation layer. In our 839 case, this methodology prevents the optimization from favoring solutions that minimize the Lapla-839 cian term by trivially setting  $\Phi(x) = 0$ .

840 841

825

826

827 828 829

830 831

832

**Dataset splitting and optimization.** The benchmark datasets are divided into 63/27/10 splits for training, validation, and testing, in alignment with Johansson et al. (2016); Shalit et al. (2017). The same test realizations are used to maintain consistency with previous studies. Optimization employs the Adam optimizer (Kingma & Ba, 2014), using the default parameters:  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , and  $\epsilon = 10^{-8}$ . We implement an exponential decay schedule for the learning rate, decreasing it by 0.95 every 50 iterations. The training process includes early stopping (Prechelt, 2002) based on the LITE objective, as defined in 12, evaluated on the validation set, with a maximum of 10,000 iterations and a patience of 2,000 iterations for early stopping.

849

Hyperparameter Selection. We follow the commonly used practice in the literature for hyperparameter selection (Johansson et al., 2016) based on the PEHE metric on the validation set, while the
hyperparameters are fixed across multiple realizations. Although not applicable to real-world data,
this approach validates the robustness of the selected parameters and prevents overfitting.

854 **Baselines.** The comparison encompasses traditional ML regression methods such as Ordinary 855 Least Squares with treatment as a feature (OLS<sub>1</sub>), linear regression with separate regressors for each 856 treatment group (OLS<sub>2</sub>), and the Least Absolute Shrinkage and Selection Operator with treatment 857 as a feature (LASSO<sub>1</sub>), and separate regressors for each treatment group (LASSO<sub>2</sub>). We also con-858 sider matching methods like k-nearest neighbor (k-NN) Ho et al. (2007), propensity-score match-859 ing (PSM) Dehejia & Wahba (2002) and perfect match (PM)Schwab et al. (2018). Additionally, 860 we evaluate tree-based algorithms including Bayesian Additive Regression Trees (BART)Chipman 861 et al. (2010); Chipman & McCulloch (2016), Random Forests (R. For.) Breiman (2001), and Causal Forests (C. For.) Wager & Athey (2018), as well as Gaussian processes such as Causal Multi-862 task Gaussian Process (CMGP) Alaa & Van Der Schaar (2017) and non-stationary Gaussian Pro-863 cess (NSGP)Alaa & Schaar (2018). We also include various representation learning methods in 864 our comparison. General representation learning methods are Treatment-Agnostic Representation 865 Network (TARNET) Shalit et al. (2017), Causal Effect Variational Autoencoder (CEVAE) Louizos 866 et al. (2017). Balanced representation learning methods are Balancing Linear Regression (BLR) 867 Johansson et al. (2016), Balancing Neural Network (BNN) Johansson et al. (2016), CounterFactual 868 Regression with Maximum Mean Discrepancy (CFR MMD) Shalit et al. (2017), CounterFactual Regression with WASSerstein distance (CFR WASS) Shalit et al. (2017), Similarity preserved Individual Treatment Effect (SITE) Yao et al. (2018), Mutual Information Treatment Network (MitNet) 870 Guo et al. (2023), Adversarial Nets for inference of Individualized Treatment Effects (GANITE) 871 Yoon et al. (2018), Adversarial Balancing-based Representation learning for Causal Effect Infer-872 ence (ABCEI) Du et al. (2021), and NOrmalizing FLows Individual Treatment Effect (NOFLITE) 873 Vanderschueren et al. (2023). 874

875

**Kernel type, scale, and distance metric.** The selection of the kernel type, scale, and distance 876 metric may influence the perfomance of our approach. While our method is compatible with var-877 ious kernels, we chose to use the Gaussian (RBF) kernel paired with the Euclidean metric in our 878 experiments, because it is the common practice in manifold learning, kernel methods, and classifi-879 cation tasks. The focus of our study was to highlight the innovative aspects of our method rather 880 than the specifics of kernel selection, which is a standard consideration across kernel-based and 881 manifold learning techniques. Additionally, since our kernel is built on the latent space and is part 882 of the loss function, the RBF kernel exhibits properties such as differentiability and smoothness, 883 which are crucial for stable optimization. To determine the most suitable kernel scale, we employed 884 cross-validation to identify the optimal bandwidth, as detailed in the Appendix.

885

894

895

896 897

898

899 900 901

902 903

904

**Betails on the News dataset.** The News dataset simulates consumer opinions on news items viewed on different devices, using 5,000 randomly sampled articles from the New York Times. Each news item is represented by word counts from a 3,477-word vocabulary. The simulated outcome is the reader's opinion, influenced by whether the news is viewed on a desktop (t = 0) or a mobile device (t = 1). Bias in the "treatment" assignment is based on the similarity between the topic distribution of the news items and two centroids, indicating a consumer preference for certain topics on mobile. This dataset enables the analysis of how the device impacts the reader's experience. For more details, see Johansson et al. (2016).

**Details on the IHDP dataset.** We follow Shalit et al. (2017) and use the same simulated outcomes from the NPCI package Dorie (2016). This dataset includes 1000 realizations for robust evaluation.

**Computing resources.** All the experiments were performed using Python on NVIDIA DGX A100 systems, each equipped with A100 GPUs and 512 GB of RAM.

**B.2** Hyperparameter selection

In the illustrative example, we utilized a learning rate of 1e-2 with a batch size of 100. The kernel scale factor was set to 1. Both the hypothesis and representation layers were configured with 4 layers each, and all layers were dimensioned at 25.

905 906 907

Table 3: Hyperparameter grids for	IHDP and News Benchmarks.
-----------------------------------	---------------------------

Parameter	IHDP	News	
Learning rate	1e-2,	1e-3, 1e-4	
Batch size	100	200, 500, 1000	
Num. Representation layers	5, 7, 9	2, 3, 4	
Dim. Representation layers	25, 50, 75	100, 150, 200	
Num. Hypothesis layers	5, 7, 9	4, 5, 6	
Dim. Hypothesis layers	25, 50, 75	100, 150, 200	
LITE Reg term $(\alpha)$	0:	0.1:10	
Kernel scale ( $\sigma$ )	0.05:0.05:1		

4.	. Scielled parameter settings for INDP and News									
	Parameter	IHDP	News							
	Learning rate	1e-3	1e-2							
	Batch size	100	500							
	Num. Representation layers	9	2							
	Dim. Representation layers	25	150							
	Num. Hypothesis layers	9	5							
	Dim. Hypothesis layers	25	100							
	LITE Reg term $(\alpha)$	9.6	1.9							
	Kernel scale ( $\sigma$ )	0.1	0.35							

# Table 4: Selected parameter settings for IHDP and News datasets.

B.3 REFERENCES FOR REPORTED PERFORMANCE METRICS.

The results of the competing methods in Tables 1 and 2 were sourced from the original papers where available, as detailed in Tables 5 and 6 respectively.

Table 5: References for reported performance metrics on the IHDP dataset (1000 iterations). 'n.r.' denotes values not reported.

		1					
938				$\sqrt{\epsilon_{\rm F}}$	EHE	$\epsilon_{\rm ATE}$	
939	Group	Method		Within-sample	Out-of-sample	Within-sample	Out-of-sample
940	Classic ML regression		OLS <sub>1</sub> OLS <sub>2</sub>	Shalit et al. (2017) Shalit et al. (2017)			
941 942	Matching		k-NN PSM PM	Shalit et al. (2017) Alaa & Schaar (2018) n.r.	Shalit et al. (2017) Alaa & Schaar (2018) Schwab et al. (2018)	Shalit et al. (2017) n.r. n.r.	Shalit et al. (2017) Chen et al. (2019) Schwab et al. (2018)
943	Tree-based		BART R. For. C. For.	Shalit et al. (2017) Shalit et al. (2017) Shalit et al. (2017)	Shalit et al. (2017) Shalit et al. (2017) Shalit et al. (2017)	Shalit et al. (2017) Shalit et al. (2017) Shalit et al. (2017)	Shalit et al. (2017) Shalit et al. (2017) Shalit et al. (2017)
945	Gaussian processes		CMGP NSGP	Alaa & Schaar (2018) Alaa & Schaar (2018)	Alaa & Schaar (2018) Alaa & Schaar (2018)	Reddy & Balasubramanian (2024) n.r.	Reddy & Balasubramanian (2024) Guo et al. (2023)
946		General	TARNET CEVAE	Shalit et al. (2017) Louizos et al. (2017)			
947			BLR BNN	Shalit et al. (2017) Shalit et al. (2017)			
948	Representation learning	Balanced	CFR MMD CFR WASS	Shalit et al. (2017) Shalit et al. (2017)			
949			SITE MitNet	Du et al. (2021)	Du et al. (2021) Guo et al. (2023)	Du et al. (2021)	Du et al. (2021) Guo et al. (2023)
950			GANITE ABCEI	Yoon et al. (2018) Du et al. (2021)	Yoon et al. (2018) Du et al. (2021)	Reddy & Balasubramanian (2024) Du et al. (2021)	Reddy & Balasubramanian (2024) Du et al. (2021)

Table 6: References for reported performance metrics on the News dataset (50 iterations). 'n.r.' denotes values not reported.

				$\sqrt{\epsilon_{\text{PEHE}}}$	$\epsilon_{\mathrm{ATE}}$		
Group		Method	In-sample	Out-of-sample	In-sample	Out-of-sample	
Classic ML regression		LASSO <sub>1</sub> LASSO <sub>2</sub>	Schrod et al. (2023) Schrod et al. (2023)	Schrod et al. (2023) Schrod et al. (2023)	Schrod et al. (2023) Schrod et al. (2023)	Schrod et al. (2023) Schrod et al. (2023)	
Gaussian processes		CMGP	n.r.	Vanderschueren et al. (2023)	n.r.	n.r.	
	General	TARNet CEVAE	Schrod et al. (2023) n.r.	Schrod et al. (2023) Vanderschueren et al. (2023)	Schrod et al. (2023) n.r.	Schrod et al. (2023) n.r.	
Representation learning	Balanced	CFR WASS SITE ABCEI NOFELITE	Schrod et al. (2023) Schrod et al. (2023) Schrod et al. (2023) n.r.	Schrod et al. (2023) Schrod et al. (2023) Schrod et al. (2023) Vanderschueren et al. (2023)	$\frac{\text{Schrod et al. (2023)}}{\frac{\text{Schrod et al. (2023)}}{\frac{\text{Schrod et al. (2023)}}{\text{n.r.}}}$	Schrod et al. (2023) Schrod et al. (2023) Schrod et al. (2023) n.r.	

# C EXTENSION TO MULTI-TREATMENT FRAMEWORK

 LITE is outlined in Algorithm 1. While presented for binary treatment for clarity, the method **969** is adaptable and can be readily extended to handle multiple treatments. Specifically, extensions **970** to multi-treatment scenarios involve adding hypothesis functions for each treatment condition. **971** These functions are optimized through the Laplacian framework by generalizing LITE as follows:  $S_{\text{LITE}}(h, \Phi) = \frac{1}{b^2} \sum_{t=0}^{m} \mathbf{h}_t^T \mathcal{L} \mathbf{h}_t$ , where *m* is the number of treatments. This extension ability presents an advantage over many existing methods that cannot be extended easily. For example,
 the covariate balancing metric in classical methods is typically designed only for binary cases.

# D LITE SCALABILITY AND COMPUTATION TIME

UTE is designed to handle large datasets efficiently and is implemented to allow for fast computation.

While some methods struggle with scalability on high-dimensional data, LITE mini-batching train-ing enhances scalability, allowing us to handle large datasets efficiently.

The computation involves two key steps. First, the computation of predicted potential outcomes involves calculating both factual and counterfactual samples by simply forwarding them through the network. Second, the computation of the Laplacian operator requires calculating pairwise distances between samples within the representation layer for kernel construction. This step encompasses the primary computational burden. The kernel construction in the representation space is common also to covariate shift minimization methods like CFRNET (Shalit et al., 2017) employing MMD or Wasserstein distance for shift minimzation. The Wasserstein distance, which generally offers better performance over the MMD, requires additional computational steps involving Sinkhorn-Knopp (Cuturi, 2013) iterations to compute the Wasserstein distance. 

For the early stopping phase of optimization, we employ the entire validation set to ensure a comprehensive geometric inference. This validation set for the News dataset includes 1,350 samples in a single batch, each with a latent space dimension of 150. We optimize pairwise distances using vectorized operations, significantly reducing processing time to just 0.002 seconds on a GPU.

To further illustrate the efficiency of LITE compared to CFRNET Wasserstein distance (using the POT package), we created a Colab notebook, which is available here(Please note: the reported time from the first run may be inaccurate due to server GPU initialization. It is recommended to run it twice).

The results show that LITE takes 0.0012 seconds on a GPU, while the CFRNET Wasserstein Distance takes 0.017 seconds, making LITE 14 times faster – an order of magnitude difference. On a T4 GPU, LITE remains highly efficient, taking less than 2 seconds for 1,000 optimization iterations, with the potential for even better performance on advanced GPUs.