

# To Predict or to Reject: Causal Effect Estimation with Uncertainty on Networked Data

Hechuan Wen<sup>1</sup>, Tong Chen<sup>1\*</sup>, Li Kheng Chai<sup>2</sup>, Shazia Sadiq<sup>1</sup>, Kai Zheng<sup>3</sup>, Hongzhi Yin<sup>1</sup>

<sup>1</sup>The University of Queensland, Australia

<sup>2</sup>Health and Wellbeing Queensland, Australia

<sup>3</sup>University of Electronic Science and Technology of China, China

{h.wen, tong.chen, h.yin1}@uq.edu.au, likheng.chai@hw.qld.gov.au, shazia@eecs.uq.edu.au, zhengkai@uestc.edu.cn

arXiv:2309.08165v1 [cs.LG] 15 Sep 2023

**Abstract**—Due to the imbalanced nature of networked observational data, the causal effect predictions for some individuals can severely violate the positivity/overlap assumption, rendering unreliable estimations. Nevertheless, this potential risk of individual-level treatment effect estimation on networked data has been largely under-explored. To create a more trustworthy causal effect estimator, we propose the uncertainty-aware graph deep kernel learning (GraphDKL) framework with Lipschitz constraint to model the prediction uncertainty with Gaussian process and identify unreliable estimations. To the best of our knowledge, GraphDKL is the first framework to tackle the violation of positivity assumption when performing causal effect estimation with graphs. With extensive experiments, we demonstrate the superiority of our proposed method in uncertainty-aware causal effect estimation on networked data. The code of GraphDKL is available at <https://github.com/uqhwen2/GraphDKL>.

**Index Terms**—causal effect estimation, networked data, uncertainty quantification, feature collapse

## I. INTRODUCTION

Estimating causal effect to support decision-making in high-stake domains such as healthcare, education, and e-commerce is crucial. With the prevalence of networked data, [1] has recently started exploring both the features of individuals (i.e., nodes) and their structural connectivity (i.e., edges) with graph neural networks (GNNs) for causal effect estimation.

Owing to the inherent nature of observational data, violation of positivity is inevitable yet potentially devastating for causal effect estimation at the individual level, as the low-confidence predictions on non-overlapping samples may suggest a wrong treatment or introduce false hope [2]. For networked data where individuals are mutually connected, the violation of positivity is further amplified because of the presence of additional structural information. As shown in Figure 1 (a), to predict the health status of older users (i.e., control group) after using dietary supplements, one may train a causal estimator based on observational data from younger users (i.e., treated groups). However, as Figure 1 (b) depicts, although decent counterfactual estimations can be made within the overlapping area, a higher risk exists when estimating in the non-overlapping area of a different group. In worse scenarios, the predicted treatment outcomes contradict the ground truth, leading to a false recommendation with adverse effect.

\*Tong Chen is the corresponding author.

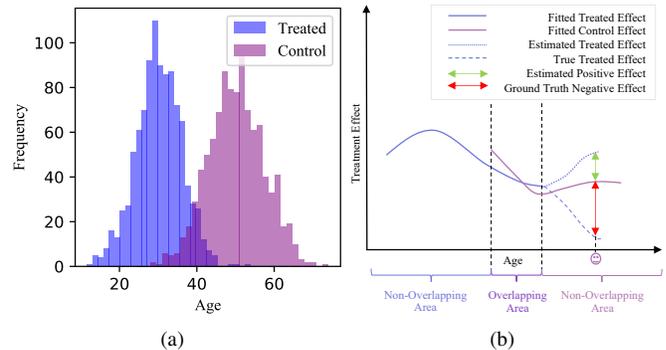


Fig. 1: (a) Histogram of two treatment groups on a one-dimensional toy dataset w.r.t. age. (b) The high risk of causal effect estimation in the non-overlapping area due to violation of positivity.

Thus, instead of blindly making recommendations based on low-confidence predictions on individual treatment effect, a more desirable capability of a causal estimator is to flag every highly uncertain estimation resulted from violation of positivity, which can be deferred for human inspections and used to guide improvements on the observational data collection process. However, existing solutions for measuring the uncertainty of each counterfactual prediction [3] are predominantly centered around tabular data without any inter-dependencies among samples. This renders existing uncertainty-aware methods unable to capture the nuanced divergence between samples in graph-structured data, given the combinatorial impact from not only individuals' own variables but also their connections with others in the network.

To fill the gap in uncertainty-aware causal estimation with networked data in the presence of positivity violation, we propose our Graph Deep Kernel Learning (GraphDKL) framework which offers uncertainty estimation to identify the likely unreliable counterfactual predictions. We introduce Gaussian process (GP) to the GNN architecture, so as to let the causal effect estimator benefit from the probabilistic nature of GP by referring to the derived prediction variance as a precise indicator of the estimator's confidence in each prediction it makes. To increase the scalability of GraphDKL on large graphs, we further design a sparse variational optimization process to replace the time-consuming covariate matrix in-

version in GP with a more computationally tractable learning objective, which significantly reduces the complexity from  $\mathcal{O}(N^3)$  to  $\mathcal{O}(M^2N)$  with  $N$  being the number of training samples ( $M \ll N$ ). Meanwhile, another notable obstacle with deep architectures used for causal effect estimation is the feature collapse issue [4], i.e., two distinct raw data points can share nearly identical representations after being mapped to the latent space via deep layers. Despite the richer information embedded after the graph convolution, the collapse of different individuals’ latent representations can seriously hinder the uncertainty quantification. For instance, a sample from the non-overlapping area is intuitively associated with a stronger uncertainty in its counterfactual prediction. However, in the networked data, if it is connected to one or more samples from the overlapping area, then its representation learned via GNN’s message passing is likely to possess high similarity with its neighbors’ representations. Consequently, uncertainty quantification based on the learned representations will assign the same individual from the non-overlapping area with a low uncertainty (i.e., high prediction confidence), which is misleading and thus undesirable. To mitigate the feature collapse, we constrain our GraphDKL model with Lipschitzness [5] to preserve the local distances in the latent space, such that the semantic manifold of the original variables is preserved in every intermediate latent space during the sequential, layer-by-layer neural mapping. Hence, predictions on high- and low-confidence samples can be effectively distinguished, making it possible for uncertainty-aware causal effect estimation on networked observational data.

## II. RELATED WORK

**Deterministic Model For Causal Effect Estimation.** So far, various neural methods [6]–[8] have been proposed due to the proliferation of deep learning (DL). These parametric models are good at modelling the individual-level causal effect and are applicable to unseen instances. Up to date, causal effect estimation has been extended to the graph domain [1], where the rich relational information is utilized to learn more robust deconfounded latent representations. However, all the above-mentioned models are deterministic, which can result in over-confident estimations [9] and is incapable of quantifying the prediction uncertainty to inform the causal estimation failure when the positivity assumption is violated.

**Probabilistic Model For Causal Effect Estimation.** It is noted that some attention has been paid to quantifying the predictive uncertainty in causal effect estimation with non-graph data. For example, the light-weight models BART [10] and CMGP [11] can offer predictive uncertainty for causal effect estimation, but they lack strong expressive power and fail to capture the complex relationship when modelling the high-dimensional data. To fix this issue, [3] and [12] leverage deep Bayesian methods to enhance the expressive power and become more capable than BART and CMGP. However, little attention has been paid to estimating the causal effect on network data with uncertainty.

## III. PRELIMINARIES

We aim to estimate individual treatment effect (ITE) on the networked data  $(\{\mathbf{x}_i, t_i, y_i\}_{i=1}^N, \mathbf{A})$ , where  $\mathbf{x}_i$ ,  $t_i$ ,  $y_i$  are respectively the raw variables, observed treatment, treatment outcome that correspond to the  $i$ -th individual, and  $\mathbf{A} \in \{0, 1\}^{N \times N}$  is the adjacency matrix indicating the connections between individuals, which can be obtained via consanguinity, doctor referrals, social networks, etc. The common practice is to learn a deconfounded latent representation  $\mathbf{z}_i$  for each individual with a GNN by aggregating its neighbour information [1], which is then used for counterfactual ITE prediction. To achieve this, three common assumptions are needed to lay the theoretical foundation.

*Assumption 1 (Stable Unit Treatment Value Assumption (SUTVA)):* For any individual  $i$ : (a) the potential outcomes for  $i$  do not vary with the treatment assigned to other individuals; and (b) there are no different forms or versions of each treatment that may lead to different potential outcomes.

*Assumption 2 (Unconfoundedness):* Treatment assignment is independent to the potential outcomes  $\{Y_{t=0}, Y_{t=1}\}$  given the latent covariate  $\mathbf{z}$ , i.e.,  $t \perp\!\!\!\perp \{Y_{t=0}, Y_{t=1}\} | \mathbf{z}$ . Note that the potential outcomes  $Y$  use a different notation w.r.t. the observed ones  $y$ .

*Assumption 3 (Positivity):* For every  $\mathbf{z}$ , the treatment assignment mechanism obeys:  $0 < p(t = 1 | \mathbf{z}) < 1$ .

### A. Hurdles in Quantifying Uncertainty with Graph Data

Based on the latent representation  $\mathbf{z}$  of each individual, we aim to estimate its treatment effect and assign an uncertainty to this estimation. Unfortunately, uncertainty quantification can be seriously poisoned by the feature collapse issue [4], especially for latent features extracted by deep neural networks (DNNs). Feature collapse describes the scenario where two distinct points in the original feature space  $\mathcal{X}$  can be mapped to two similar or even identical positions in the latent space  $\mathcal{Z}$ . Consequently, predictions on non-overlapping samples could be incorrectly assigned an uncertainty as low as predictions on overlapping ones due to their collapsed representations.

Despite the popularity of GNNs in learning individual representations for ITE estimation, little attention has been paid to the potential feature collapse issue. So far, the state-of-the-art GNN backbones, e.g, GraphSAGE [13], rely on non-linear mapping and are hence vulnerable to feature collapse. The message-passing scheme in GNNs potentially deteriorates the uniqueness of learned representations even further. For a proof-of-concept, we generate a toy graph with two-dimensional node features and four classes, as shown in Figure 2 (a). We train a 1-layer GraphSAGE with classes 0, 1 and 2, and nodes from class 3 are held out. The two-dimensional visualization in Figure 2 (b) shows that the representations generated by the trained GraphSAGE for class 3 nodes mostly collapse with class 0 in the latent space. Such collapse is different from the over-smoothing problem with GNNs, as only a shallow 1-layer structure is used and the node representations for the first three classes do not collide.

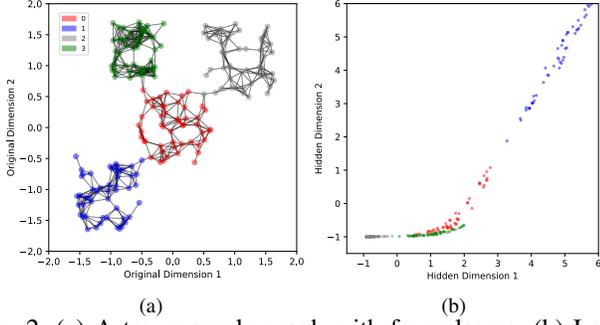


Fig. 2: (a) A toy example graph with four classes; (b) Latent representation from a 1-Layer GraphSAGE.

#### IV. METHODOLOGY

In this section, we present our graph deep kernel learning (GraphDKL) framework for handling the causal effect estimation with uncertainty on graph-structured data.

##### A. Lipschitz-constrained Graph Representation Learning

As a versatile framework, GraphDKL is agnostic to any GNNs. Without loss of generality, we leverage GraphSAGE [13] as the base GNN given its balance between efficiency and effectiveness, and the ability to scale with batch training. At the  $l$ -th layer of GraphSAGE, the core neural operation to learn the latent representation for individual/node  $i$  is:

$$\mathbf{h}_i^l = \sigma(\mathbf{W}_l \cdot \text{MEAN}(\{\mathbf{h}_i^{l-1}\} \cup \{\mathbf{h}_j^{l-1}, \forall j \in \mathcal{N}(i)\})), \quad (1)$$

where  $\sigma$  is the non-linear activation,  $\mathbf{W}_l$  is the weight matrix at layer  $l$ , while the mean aggregator  $\text{MEAN}(\cdot)$  is used to merge the representations of node  $i$  and its neighbours  $j \in \mathcal{N}(v)$  from layer  $l-1$ .

To facilitate uncertainty estimation, our proposed framework, shown in Figure 3, combines the GNN with deep kernel-based GP to get the best of both worlds – the deconfounded node representations containing both individual features and structural information are extracted via GNN first, and the learned representations are fed into two independent DNNs, with each of them mapping the graph-based representations to treated and control latent spaces for subsequent predictions. For notation simplicity, we omit the formulation of each DNN, which is a multi-layer perceptron (MLP) with  $L'$  layers, and takes the final-layer representation  $\mathbf{h}_i^{L'}$  from GNN as its input. Unless specified, the following descriptions on DNNs apply to both treatment branches  $t \in \{0, 1\}$ .

**Decoupling Collapsed Representations.** To alleviate feature collapse and ensure accurate uncertainty estimation, we propose to preserve the local distance among points after non-linear mapping. In GraphDKL, this constraint needs to be enforced in both the GNN for learning individual representations, as well as the two DNN branches that respectively model treated and control groups. In a nutshell, the distance  $\|\mathbf{x}_i - \mathbf{x}_j\|$  between any two points  $\mathbf{x}_i$  and  $\mathbf{x}_j$  from the raw feature space has a corresponding meaningful distance in the latent space. To achieve this desired property, we introduce the notion of Lipschitz constant. Specifically, for each given function  $\mathbf{s}' = f(\mathbf{s})$  with input  $\mathbf{s}$  and output  $\mathbf{s}'$ , then the

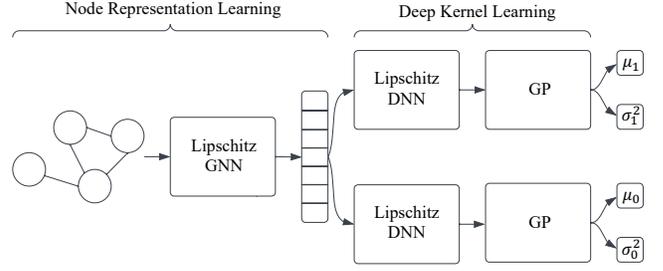


Fig. 3: Structure of GraphDKL framework. The Lipschitz prefix denotes the Lipschitz-constrained neural networks.

Lipschitz constant  $\text{Lip}(f)$  w.r.t.  $f(\cdot)$  satisfies that, for any pair of inputs  $(\mathbf{s}_1, \mathbf{s}_2)$ ,  $\|\mathbf{s}'_1 - \mathbf{s}'_2\| \leq \text{Lip}(f) \|\mathbf{s}_1 - \mathbf{s}_2\|$  holds. In other words,  $\text{Lip}(f) \geq \frac{\|\mathbf{s}'_1 - \mathbf{s}'_2\|}{\|\mathbf{s}_1 - \mathbf{s}_2\|}$  for any  $(\mathbf{s}_1, \mathbf{s}_2)$  pair. If  $\text{Lip}(f) \leq 1$ , then it essentially means that the difference in function values is controlled by the original pairwise distance obtained from the input space. This property ensures that small changes in the input result in small changes in the output, providing a sense of stability and predictability, and  $f(\cdot)$  is also termed 1-Lipschitz (local distance preserving). With the context given, we define the 1-Lipschitz GraphDKL below.

*Theorem 1 (1-Lipschitz GraphDKL):* GraphDKL has  $L$  layers of graph convolution  $\mathbf{H}^L = g_L(g_{L-1}(\dots g_1(\mathbf{X}, \mathbf{A})))$  in the GNN where  $\mathbf{X} \in \mathbb{R}^{N \times D}$ ,  $\mathbf{A} \in \{0, 1\}^{N \times N}$ , and  $\mathbf{H}^L \in \mathbb{R}^{N \times S}$  respectively denote the  $D$ -dimensional raw variables, adjacency matrix, and  $S$ -dimensional latent representations of  $N$  individuals. The GNN is followed by an  $L'$ -layer DNN  $\mathbf{Z}^{L'} = \phi_{L'}(\phi_{L'-1}(\dots \phi_1(\mathbf{H}^L)))$  in either the treated/control branch, with  $\mathbf{Z}^{L'} \in \mathbb{R}^{N \times S}$  being the  $N$  representations from the final layer. The entire representation learning pipeline in GraphDKL is 1-Lipschitz if:

$$\begin{cases} \text{Lip}(g_l) \leq 1, & \forall l \leq L \\ \text{Lip}(\phi_{l'}) \leq 1, & \forall l' \leq L' \end{cases}, \quad (2)$$

where  $\text{Lip}(g_l)$  and  $\text{Lip}(\phi_{l'})$  respectively denote the Lipschitz constant of a single GNN and DNN layer.

*Proof.* We denote the hidden representation of individual  $i$  at the  $l$ -th GNN layer as  $\mathbf{h}_i^l$ , and that of the same  $i$  at the  $l'$ -th DNN layer as  $\mathbf{z}_i^{l'}$ . Note that the raw feature  $\mathbf{x}_i$  is the input to the first-layer GNN, whose final-layer representation  $\mathbf{h}_i^L$  is the input to the first-layer DNN. Then, for any pair of instances  $(i, j)$ , we have:

$$\begin{aligned} \frac{\|\mathbf{z}_i^{L'} - \mathbf{z}_j^{L'}\|}{\|\mathbf{x}_i - \mathbf{x}_j\|} &= \frac{\|\mathbf{z}_i^{L'} - \mathbf{z}_j^{L'}\|}{\|\mathbf{z}_i^{L'-1} - \mathbf{z}_j^{L'-1}\|} \times \dots \times \frac{\|\mathbf{z}_i^1 - \mathbf{z}_j^1\|}{\|\mathbf{h}_i^L - \mathbf{h}_j^L\|} \times \\ &\quad \frac{\|\mathbf{h}_i^L - \mathbf{h}_j^L\|}{\|\mathbf{h}_i^{L-1} - \mathbf{h}_j^{L-1}\|} \times \dots \times \frac{\|\mathbf{h}_i^1 - \mathbf{h}_j^1\|}{\|\mathbf{x}_i - \mathbf{x}_j\|} \\ &\leq \text{Lip}(\phi_{L'}) \times \dots \times \text{Lip}(\phi_1) \times \text{Lip}(g_L) \times \dots \times \text{Lip}(g_1). \end{aligned} \quad (3)$$

As every  $\text{Lip}(g_l)$  and  $\text{Lip}(\phi_{l'})$  is no larger than 1 by premise,  $\frac{\|\mathbf{z}_i^{L'} - \mathbf{z}_j^{L'}\|}{\|\mathbf{x}_i - \mathbf{x}_j\|} \leq 1$ . We thus conclude that the entire neural mapping from  $\mathbf{x}_i$  to  $\mathbf{z}_i^{L'}$  in GraphDKL is 1-Lipschitz.  $\square$

To ensure the local distance is preserved for each neural mapping layer, *spectral normalization* has been proven rigorous [5] for enforcing 1-Lipschitz. Taking the weight matrix  $\mathbf{W}_l$  at the  $l$ -th graph convolution layer in (1) as an example, the spectral normalization states:

$$\text{Lip}(g_l) \leq 1, \quad \text{if } \|\mathbf{W}_l\|_2 \leq 1, \quad (4)$$

where  $\|\cdot\|_2$  denotes the spectral norm, i.e.,  $L_2$  matrix norm of  $\mathbf{W}_l$ . Compared with using the spectral norm as a penalization term for *regularization* purpose, we formulate a *normalization* process that strictly bounds the spectral norm to a designated value, and this ensures obedience of 1-Lipschitz at all layers and thus benefit the measurement of uncertainty.

As  $\|\mathbf{W}_l\|_2$  corresponds to the largest singular value of matrix  $\mathbf{W}_l$  which is known to be time-consuming to compute exactly, we perform power iteration [14] over  $\mathbf{W}_l$  to obtain an approximation  $\tau$  of the spectral norm, which is a lower bound on the largest singular value  $\|\mathbf{W}_l\|_2$ . Then, the weight matrix is normalized as:

$$\overline{\mathbf{W}}_l = \frac{1}{\tau} \mathbf{W}_l, \quad (5)$$

which empirically makes  $\text{Lip}(g_l) \leq 1$  consistent across all scenarios by rescaling  $\mathbf{W}_l$  [5]. Analogously, the same spectral normalization is adopted on all DNN layers' weight matrices.

## B. Deep Kernel Learning

In each treatment effect prediction branch, the deep kernel learning (DKL) module passes the latent representation  $\mathbf{z}_i^{L'}$  from the final DNN layer into a Gaussian process (GP) for causal effect estimation with uncertainty. From now on, we let  $\mathbf{z}_i = \mathbf{z}_i^{L'}$  for better clarity when there is no ambiguity.

**Standard GP.** A standard GP is a finite number of random variables which have a joint Gaussian distribution [15]. Mathematically, it is denoted as  $\mathcal{GP}$ , with mean function  $\mu(\cdot) : \mathcal{X} \rightarrow \mathbb{R}$  and covariance function  $k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  over the real-valued stochastic function  $f(\cdot) : \mathcal{X} \rightarrow \mathbb{R}$  whose input is the  $D$ -dimensional variable vector  $\mathbf{x} \in \mathbb{R}^D$ , namely,

$$f(\mathbf{x}) \sim \mathcal{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')). \quad (6)$$

By evaluating the GP at  $N$  samples  $\mathbf{Z} = \{\mathbf{z}_i\}_{i=1}^N$  (any subset from domain  $\mathcal{X}$ ), we end up with  $N$  multivariate Gaussian distributions  $\mathbf{f} = \{\mathbf{f}_i\}_{i=1}^N$  as follows:

$$\mathbf{f} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}), \quad (7)$$

where  $\boldsymbol{\mu} \in \mathbb{R}^N$  is the variance vector and  $\mathbf{K} \in \mathbb{R}^{N \times N}$  is the covariance matrix. With  $i, j$  for indexing,  $\boldsymbol{\mu}[i] = \mu(\mathbf{x}_i)$  is  $i$ 's mean,  $\mathbf{K}[i, j] = k(\mathbf{x}_i, \mathbf{x}_j)$  is the covariance between  $i$  and  $j$ .

**Ramping up Expressiveness.** Given the limited capacity of standard GP in learning the latent distributions [16], [17], recent frameworks have been expanding the expressiveness of the standard GP. For instance, deep Gaussian process [16] stacks a series of GPs, and deep kernel [17] utilizes the latent variables produced from a deep learning method for the GP. In this paper, we investigate the deep kernel framework since it is a natural extension to our neural architecture, as well as its superiority in expressiveness and computational efficiency.

Specifically, by replacing the raw variables  $\mathbf{x}$  with the latent output  $\mathbf{z}$  from GraphDKL's neural mapping, (7) is updated with  $\boldsymbol{\mu}[i] = \mu(\mathbf{z}_i)$  and  $\mathbf{K}[i, j] = k(\mathbf{z}_i, \mathbf{z}_j)$ .

In GP, a mean of zero is normally assumed, i.e.,  $\boldsymbol{\mu} = \mathbf{0}$ , and we consider the infinitely smooth radial basis function as the kernel for computing the covariance, i.e.,  $k(\mathbf{z}_i, \mathbf{z}_j) = \sigma_{ker}^2 \exp(-\frac{(\mathbf{z}_i - \mathbf{z}_j)^2}{2l^2})$ , where  $\boldsymbol{\theta}_{ker} = \{\sigma_{ker}, l\}$  is a parameter set of the GP kernel to be optimized. With  $N$  latent representations  $\mathbf{Z} \in \mathbb{R}^{N \times S}$  learned from  $\mathbf{X}$  and their corresponding real-valued labels  $\mathbf{y} = \{y_i\}_{i=1}^N$ , we can obtain the joint marginal likelihood w.r.t. the updated (7) as follows:

$$p(\mathbf{y}|\mathbf{Z}, \boldsymbol{\theta}_{ker}) = \int \prod_{i=1}^N p(y_i|\mathbf{f}_i, \mathbf{z}_i) p(\mathbf{f}_i|\mathbf{z}_i) d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}), \quad (8)$$

where the optimal kernel parameters are obtained by finding its maximum through gradients.

**ITE Prediction with Uncertainty Quantification.** To evaluate the model at an arbitrary test point  $\mathbf{x}_* \in \mathbb{R}^D$  with corresponding latent representation  $\mathbf{z}_*$ , we leverage the property that the joint distribution of the training labels  $\mathbf{y}$  and the test label  $y_*$  is still Gaussian:

$$\begin{pmatrix} y_* \\ \mathbf{y} \end{pmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} k(\mathbf{z}_*, \mathbf{z}_*) & \mathbf{k}_*^T \\ \mathbf{k}_* & \mathbf{K} \end{bmatrix}\right), \quad (9)$$

where column vector  $\mathbf{k}_* \in \mathbb{R}^N$  measures the covariance between  $\mathbf{z}_*$  and all  $N$  training samples, i.e.,  $\mathbf{k}_*[i] = k(\mathbf{z}_*, \mathbf{z}_i)$ . Thus, the posterior label distribution is:

$$y_* | (\mathbf{z}_*, \mathbf{Z}, \mathbf{y}, \boldsymbol{\theta}_{ker}) \sim \mathcal{N}(\mu_*, \sigma_*^2), \quad (10)$$

which has the following closed-form solution [15]:

$$\mu_* = \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{y}, \quad \sigma_*^2 = k(\mathbf{z}_*, \mathbf{z}_*) - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_*, \quad (11)$$

where the mean  $\mu_*$  serves as the prediction of  $y_*$  w.r.t. a treatment  $t \in \{0, 1\}$ , and the variance  $\sigma_*^2$  of the prediction is used as a direct indicator of the uncertainty w.r.t. sample  $\mathbf{x}_*$ . Essentially, the value of  $\sigma_*^2$  holds a positive correlation with the uncertainty.

## C. Sparse Variational Optimization

If the exact Gaussian process were applied to our proposed GraphDKL framework, the model would suffer from a  $\mathcal{O}(N^3)$  complexity due to the inversion of the covariance matrix  $\mathbf{K} \in \mathbb{R}^{N \times N}$ , which is computationally prohibitive when handling large graphs. To increase scalability, we adopt a sparse GP [15] with stochastic variational inference (SVI) [18] to our setting, building a computationally tractable GraphDKL framework with the ability to scale better.

We start by assuming a set of latent inducing points  $\mathbf{M} = \{\mathbf{m}_i\}_{i=1}^M$  in the same latent space as  $\mathbf{Z}$ . For the stochastic function  $f(\cdot)$  with Gaussian prior in (6), we obtain the corresponding outputs  $\mathbf{v} = f(\mathbf{Z})$  and  $\mathbf{u} = f(\mathbf{M})$  w.r.t. the distributions of  $\mathbf{Z}$  and  $\mathbf{M}$  in the same space, respectively. To form a tractable objective, we derive the evidence lower bound  $\mathcal{L}$  as follows:

$$\begin{aligned} \log p(\mathbf{y}) &= \log \int p(\mathbf{y}, \mathbf{v}, \mathbf{u}) d\mathbf{v} d\mathbf{u} = \log \int \frac{p(\mathbf{y}, \mathbf{v}, \mathbf{u})}{q(\mathbf{v}, \mathbf{u})} q(\mathbf{v}, \mathbf{u}) d\mathbf{v} d\mathbf{u} \\ &= \log \mathbb{E}_{q(\mathbf{v}, \mathbf{u})} \left[ \frac{p(\mathbf{y}, \mathbf{v}, \mathbf{u})}{q(\mathbf{v}, \mathbf{u})} \right] \geq \mathcal{L} = \mathbb{E}_{q(\mathbf{v}, \mathbf{u})} \left[ \log \frac{p(\mathbf{y}, \mathbf{v}, \mathbf{u})}{q(\mathbf{v}, \mathbf{u})} \right], \end{aligned} \quad (12)$$

where the SVI process approximates the posterior  $q(\mathbf{v}, \mathbf{u})$  by minimizing the Kullback-Leibler divergence  $\text{KL}(q||p)$  between the variational posterior  $q$  and the prior  $p$  [18]. As we essentially aim to perform SVI with a set of global variables, we let  $\mathbf{u}$  take this role with a variational distribution  $q(\mathbf{u})$ , and follow the widely accepted variational posterior [18]  $q(\mathbf{v}, \mathbf{u}) = p(\mathbf{v}|\mathbf{u})q(\mathbf{u})$ . We set  $q(\mathbf{u}) \sim \mathcal{N}(\mathbf{u}|\boldsymbol{\mu}_u, \mathbf{K}_u)$  with mean  $\boldsymbol{\mu}_u \in \mathbb{R}^M$  and covariance  $\mathbf{K}_u \in \mathbb{R}^{M \times M}$  to be learned. To this end, the evidence lower bound (ELBO)  $\mathcal{L}$  can be further decomposed into the following form, with an additional constraint on  $q(\mathbf{u})$ :

$$\mathcal{L} = \mathbb{E}_{q(\mathbf{v})} [\log p(\mathbf{y}|\mathbf{v})] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})), \quad (13)$$

where  $p(\mathbf{u})$  is a priori. Since posterior  $q(\mathbf{u})$  is Gaussian, with the analytically achievable  $p(\mathbf{v}|\mathbf{u})$  analogous to (9) and (10) by conditioning on the prior  $p(\mathbf{u})$ , the variational posterior  $q(\mathbf{v})$  can be analytically obtained as follows:

$$q(\mathbf{v}) = \int p(\mathbf{v}|\mathbf{u})q(\mathbf{u})d\mathbf{u} = \mathcal{N}(\mathbf{v}|\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}), \quad (14)$$

where  $\tilde{\boldsymbol{\mu}}$  and  $\tilde{\boldsymbol{\Sigma}}$  are parameterized w.r.t.  $\boldsymbol{\mu}_u$  and  $\mathbf{K}_u$ . Therefore, the first term in the simplified  $\mathcal{L}$  can be calculated with Monte Carlo sampling since posterior  $q(\mathbf{v})$  is available with the unknown parameters  $\{\boldsymbol{\mu}_u, \mathbf{K}_u\}$  to be learned. Notably, the time complexity of sparse variational GraphDKL is significantly reduced due mainly to the smaller  $M \times M$  matrix to be inverted, bringing a non-dominant  $\mathcal{O}(M^3)$  complexity. Consider the matrix multiplication in deriving  $\tilde{\boldsymbol{\Sigma}}$  in 13 for the  $q(\mathbf{v})$  to sample  $N$  times in order to calculate  $\mathcal{L}$ , the asymptotic time complexity is capped to  $\mathcal{O}(M^2N)$  [15], [18]. With  $M \ll N$  in our case, handling large graph-structured data with GraphDKL is tractable.

By optimizing the tractable objective  $\mathcal{L}$  derived by SVI, we solve the unknown parameters such that the variational posterior  $q(\mathbf{v})$ , which only depends on the input  $\mathbf{x}$ , can be optimized to fit the data. Eventually, with the fully trained GraphDKL, given a test point  $\mathbf{x}_*$  (with latent representation  $\mathbf{z}_*$ ), we can use the much smaller  $\mathbf{K}' \in \mathbb{R}^{M \times M}$  given by priori,  $\mathbf{z}_*$ 's covariance with those  $M$  prior samples  $\mathbf{k}'_* \in \mathbb{R}^M$ , and the optimized posterior  $q(\mathbf{u})$  to obtain prediction (i.e., mean  $\mu_*$  and variance (i.e., uncertainty)  $\sigma_*^2$ ):

$$\mu_* = \Gamma \mathbf{u}_u, \quad \sigma_*^2 = k(\mathbf{z}_*, \mathbf{z}_*) - \Gamma(\mathbf{K}' - \mathbf{K}_u)\Gamma^T, \quad (15)$$

where  $\Gamma = \mathbf{k}'_*{}^T \mathbf{K}'^{-1}$ .

To conclude, the final ITE estimation and its associated prediction uncertainty w.r.t. test sample  $\mathbf{x}_*$  have the following approximations:

$$\begin{aligned} \text{ITE}_* &= \mathbb{E}[Y_{t=1} - Y_{t=0} | \mathbf{z}_*] \simeq \mu_{*,t=1} - \mu_{*,t=0}, \\ \text{Uncertainty}_* &= \mathbb{E}[(Y_{t=1} - Y_{t=0})^2 | \mathbf{z}_*] \simeq \sigma_{*,t=1}^2 + \sigma_{*,t=0}^2. \end{aligned} \quad (16)$$

We present our experimental analysis in this section.

### A. Experimental Setup

1) *Dataset*: We adopt two public benchmarks with networked observational data: **Blogcatalog** [1] and **Flickr** [1]. Both BlogCatalog and Flickr datasets are processed and simulated in the same practice in [1]. For both datasets, three settings are created with  $k = 0.5, 1, \text{ and } 2$ , respectively, where  $k$  denotes the magnitude of the imbalance in the semi-synthetic dataset. The higher the  $k$  value, the more imbalanced the dataset is. In total, we evaluate the model performance in six different scenarios.

2) *Metric*: We use precision in estimation of heterogeneous effect (PEHE) [6], a well-established metric defined as  $\sqrt{\epsilon_{\text{PEHE}}} = \sqrt{\sum_{i=1}^N ((Y_{i,t=1} - Y_{i,t=0}) - (\mu_{i,t=1} - \mu_{i,t=0}))^2 / N}$  for measuring the treatment estimation accuracy at the individual level. The lower the  $\sqrt{\epsilon_{\text{PEHE}}}$ , the better the performance.

3) *Baselines*: Note, that our proposed GraphDKL is the first model to handle causal effect estimation with uncertainty on graph data. To obtain better comparisons, we share the learned node representation with the other baselines which can only be operated on the non-graph data: BART [10], BCFRMMD [3], BCEVAE [3], and CMGP [11].

4) *Evaluation Scheme on Uncertainty Quantification*: We randomly split each dataset into train/val/test with a 3/1/1 ratio. To evaluate the most effective uncertainty-aware method (a.k.a. rejection method), we reject the estimations with the highest uncertainty and calculated the  $\sqrt{\epsilon_{\text{PEHE}}}$  over the retained samples: the lower the retained  $\sqrt{\epsilon_{\text{PEHE}}}$ , the better the rejection method. As setting an uncertainty threshold for rejection can be domain-specific in real-world cases, here we use the specific uncertainty threshold given by each method that rejects a certain proportion of the top most-uncertain test samples. We test on an increasing proportion  $\{0\%, 5\%, 10\%, 15\%, 20\%, 25\%, 30\%, 50\%, 70\%, 90\%\}$  in our experiments, where the  $\sqrt{\epsilon_{\text{PEHE}}}$  scores are reported for all methods from the same amount of retained test samples.

### B. Rejection Policy Performance

Since the main task of this paper is to explore the pivot of positivity assumption with the uncertainty-aware model for causal effect estimation on graph data. We compare our proposed GraphDKL with various rejection methods with the main results shown in Table I. When compared to other rejection methods, our method GraphDKL always initializes with a lower  $\sqrt{\epsilon_{\text{PEHE}}}$  at 0% rejection rate, even though all the other baselines designed for independent data leverage the same learned node representations from the GraphSAGE convolution. Furthermore, GraphDKL outperforms all the other models over the retained test set in terms of the following key performance: (1) it keeps rejecting the bad estimation while preserving the lowest  $\sqrt{\epsilon_{\text{PEHE}}}$  on both datasets under most settings; (2) it has the fastest error convergence with an increased rejection rate.

Table I: The proportion in each column represents the fixed percentage of the test samples rejected by each method. Thus, we calculate the  $\sqrt{\epsilon_{PEHE}}$  over the same-size retained test samples by averaging results from 10 simulations for each setting.

Dataset		BlogCatalog										Flickr									
$k$	Method	0%	5%	10%	15%	20%	25%	30%	50%	70%	90%	0%	5%	10%	15%	20%	25%	30%	50%	70%	90%
0.5	BART	10.15	10.18	10.21	10.23	10.28	10.30	10.31	10.16	9.73	9.27	8.86	8.86	8.86	8.86	8.87	8.87	8.88	8.89	8.90	8.85
	BCFRMMD	7.93	6.96	6.24	5.53	5.04	4.56	4.25	3.72	3.66	3.83	48.65	44.86	43.01	41.67	40.82	39.60	38.97	35.02	28.40	3.47
	BCEVAE	42.76	37.7	33.29	29.06	25.59	23.07	21.67	20.79	21.25	21.38	53.81	38.13	33.54	31.18	29.26	28.05	27.41	25.43	22.54	9.60
	CMGP	10.89	10.18	9.41	9.31	9.32	9.15	8.99	8.73	9.14	9.33	10.37	5.86	4.97	4.38	4.05	3.77	3.63	3.28	2.99	2.86
	GraphDKL	<b>4.31</b>	<b>4.21</b>	<b>3.98</b>	<b>3.80</b>	<b>3.67</b>	<b>3.48</b>	<b>3.31</b>	<b>2.90</b>	<b>2.64</b>	<b>2.21</b>	<b>3.92</b>	<b>3.24</b>	<b>3.10</b>	<b>3.01</b>	<b>2.95</b>	<b>2.92</b>	<b>2.92</b>	<b>2.84</b>	<b>2.69</b>	<b>2.46</b>
1	BART	12.97	12.89	12.94	12.96	13.01	13.07	13.15	13.00	12.16	11.63	15.70	15.70	15.7	15.7	15.62	15.63	15.64	15.46	15.41	15.51
	BCFRMMD	10.45	9.29	8.58	7.93	7.37	6.90	6.46	6.13	6.33	6.69	10.87	7.59	6.12	5.20	4.59	4.14	<b>3.81</b>	<b>3.37</b>	<b>3.40</b>	<b>3.61</b>
	BCEVAE	36.75	33.11	29.97	27.51	24.93	22.99	21.85	20.83	21.40	21.82	24.63	16.09	12.99	11.36	10.67	10.33	10.11	9.91	10.03	10.06
	CMGP	11.46	9.84	9.18	8.69	8.44	8.01	7.70	7.06	6.73	6.76	18.51	10.15	8.28	7.23	6.58	6.17	5.87	5.28	5.15	6.41
	GraphDKL	<b>5.08</b>	<b>4.79</b>	<b>4.48</b>	<b>4.32</b>	<b>4.17</b>	<b>4.00</b>	<b>3.90</b>	<b>3.76</b>	<b>3.63</b>	<b>2.97</b>	<b>7.29</b>	<b>4.49</b>	<b>4.23</b>	<b>4.11</b>	<b>4.04</b>	<b>3.97</b>	3.93	3.91	3.91	3.71
2	BART	30.96	31.32	31.51	31.65	31.97	31.63	31.48	31.78	30.10	29.50	28.80	28.80	28.81	28.75	28.76	28.76	28.73	28.53	28.8	28.38
	BCFRMMD	26.83	23.31	21.09	19.11	17.55	16.12	15.06	14.15	14.56	15.14	18.80	13.20	10.59	9.05	7.95	7.12	6.50	5.71	5.83	6.22
	BCEVAE	42.76	37.70	33.29	29.06	25.59	23.07	21.67	20.79	21.25	21.38	20.95	13.76	11.31	10.09	9.49	9.20	8.96	8.55	8.60	8.82
	CMGP	28.20	25.00	23.87	23.06	22.70	22.32	22.10	21.19	21.39	23.20	34.66	20.24	17.05	14.84	13.54	12.77	12.48	12.91	15.09	23.36
	GraphDKL	<b>10.32</b>	<b>8.92</b>	<b>8.13</b>	<b>7.56</b>	<b>7.23</b>	<b>6.82</b>	<b>6.64</b>	<b>5.96</b>	<b>5.49</b>	<b>5.33</b>	<b>12.38</b>	<b>6.73</b>	<b>6.26</b>	<b>6.00</b>	<b>5.85</b>	<b>5.76</b>	<b>5.71</b>	<b>5.67</b>	<b>5.58</b>	<b>5.38</b>

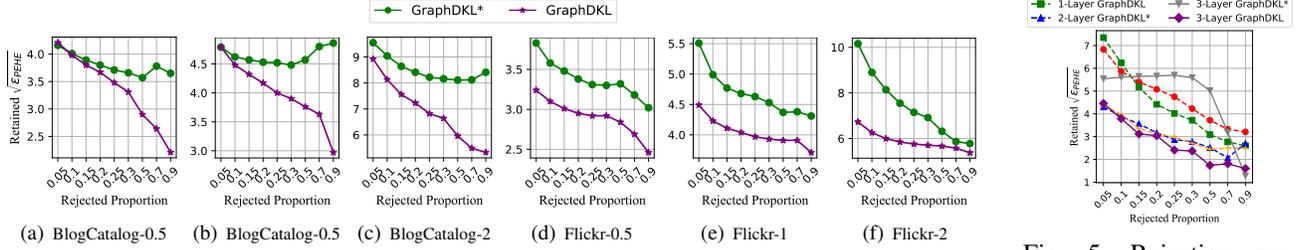


Fig. 4: Ablation study on Lipschitz constraint of the neural mapping. Lipschitz-constrained GraphDKL has a clear performance gain over GraphDKL\* without such constraint.

### C. Ablation Study on Lipschitzness

We conduct a detailed ablation study on the Lipschitz constraint. We use GraphDKL and GraphDKL\* to respectively denote variants with and without this constraint. As shown in 4, GraphDKL is superior to GraphDKL\* across all the scenarios. Note, that the base model’s performance on BlogCatalog datasets has a bouncing-back retained  $\sqrt{\epsilon_{PEHE}}$  when rejecting more samples on the test set. Additionally, we compare its influence to GNNs with varying capacity by using different graph convolution layers. Based on the results in Figure 5, the Lipschitz-constrained 3-Layer GraphDKL has the best rejection performance as shown in Figure 5 by decoupling the collapsed representation to get more accurate uncertainty of each estimation, while the proposed spectral norm can effectively bring performance gain for different GNN variants in uncertainty-aware counterfactual prediction.

## VI. CONCLUSION

We investigate the violation of the positivity assumption for causal effect estimation on graph data and take a novel perspective to create a safer causal estimator on graph data – quantifying the estimation uncertainty. Extensive experiments on the two widely used semi-synthetic graph datasets show the superiority of our proposed Lipschitz GraphDKL over the other baselines in terms of identifying high-risk estimations.

## VII. ACKNOWLEDGEMENT

This work is supported by the Australian Research Council under the streams of Future Fellowship (No. FT210100624), Discovery Project (No. DP190101985), Discovery Early Career Researcher Award (No. DE230101033), and Industrial Transformation Training Centre (No. IC200100022).

## REFERENCES

- [1] R. Guo, J. Li, and H. Liu, “Learning individual causal effects from networked observational data,” in *WSDM*, 2020.
- [2] A. Jesson *et al.*, “Quantifying ignorance in individual-level causal-effect estimates under hidden confounding,” in *ICML*, 2021.
- [3] A. Jesson, S. Mindermann, U. Shalit, and Y. Gal, “Identifying causal-effect inference failure with uncertainty-aware models,” in *NIPS*, 2020.
- [4] J. van Amersfoort, L. Smith, A. Jesson, O. Key, and Y. Gal, “On feature collapse and deep kernel learning for single forward pass uncertainty,” *arXiv preprint arXiv:2102.11409*, 2021.
- [5] J. Behrmann, W. Grathwohl, R. T. Chen, D. Duvenaud, and J.-H. Jacobsen, “Invertible residual networks,” in *ICML*, 2019.
- [6] U. Shalit, F. D. Johansson, and D. Sontag, “Estimating individual treatment effect: generalization bounds and algorithms,” in *ICML*, 2017.
- [7] C. Shi, D. Blei, and V. Veitch, “Adapting neural networks for the estimation of treatment effects,” in *NIPS*, 2019.
- [8] S. Zheng, H. Yin, T. Chen, Q. V. H. Nguyen, W. Chen, and L. Zhao, “Dream: Adaptive reinforcement learning based on attention mechanism for temporal knowledge graph reasoning,” *SIGIR*, 2023.
- [9] D.-B. Wang, L. Feng, and M.-L. Zhang, “Rethinking calibration of deep neural networks: Do not be afraid of overconfidence,” in *NIPS*, 2021.
- [10] H. A. Chipman, E. I. George, and R. E. McCulloch, “Bart: Bayesian additive regression trees,” *The Annals of Applied Statistics*, 2010.
- [11] A. Alaa and M. Van Der Schaar, “Bayesian inference of individualized treatment effects using multi-task gaussian processes,” in *NIPS*, 2017.
- [12] Y. Zhang *et al.*, “Learning overlapping representations for the estimation of individualized treatment effects,” in *AISTATS*, 2020.
- [13] W. Hamilton, Z. Ying, and J. Leskovec, “Inductive representation learning on large graphs,” in *NIPS*, 2017.
- [14] H. Gouk, E. Frank, B. Pfahringer, and M. Cree, “Regularisation of neural networks by enforcing lipschitz continuity,” *Machine Learning*, 2021.
- [15] C. K. Williams and C. E. Rasmussen, *Gaussian processes for machine learning*. MIT press Cambridge, MA, 2006.
- [16] A. Damianou and N. D. Lawrence, “Deep gaussian processes,” in *AISTATS*, 2013.
- [17] A. G. Wilson, Z. Hu, R. Salakhutdinov, and E. P. Xing, “Deep kernel learning,” in *AISTATS*, 2016.
- [18] H. Salimbeni and M. Deisenroth, “Doubly stochastic variational inference for deep gaussian processes,” in *NIPS*, 2017.