# Graph Neural Networks for Causal Inference Under Network Confounding

Michael P. Leung [†]         Pantelis Loupos[‡]

Marc

**Abstract.**

# 1   Introduction

Treatment assignment is said to be unconfounded if it is as good as random within subpopulations of observationally equivalent units. In settings where the stable unit treatment value assumption (SUTVA) is plausible, units with identical covariates are naturally considered observationally equivalent. However, when units are connected through a network, they may differ on other observed dimensions that may confound causal inference if SUTVA is violated and interference is mediated by the network. These dimensions include, for example, the number of type-**x** neighbors, the number of type-**x** neighbors with **m** neighbors of type **y**, and so on through higher-order neighbors.

Existing formulations of unconfoundedness only utilize a small subset of these dimensions. For example, a common set of controls used in the literature is the vector consisting of own covariates, number of neighbors, and average covariates of neighbors. This choice may be difficult to justify in practice due to a lack of behavioral models of selection. Neighbor covariates may influence selection into treatment in more complex ways not adequately captured by the mean. Furthermore, this choice of controls implies no confounding from higher-order neighbors, which we show rules out economically interesting sources of interference in treatment selection, such as endogenous peer effects.

In this paper, we study estimation and inference for treatment and spillover effects under a fully nonparametric formulation of unconfoundedness motivated by a model of selection. To allow for peer effects, selection is governed by the reduced form of a simultaneous-equations model, which is a function of the entirety of $\boldsymbol{X}$, the matrix of all units' covariates, and $\boldsymbol{A}$, the network adjacency matrix. As a result, it is not generally possible to summarize confounding by a simple low-dimensional function of these objects. Our unconfoundedness condition therefore considers units observationally equivalent if they occupy identical positions in the network, meaning that they match on all observed neighborhood and higher-order neighborhood dimensions.

Existing methods that rule out complex forms of interference in selection may result in biased estimates of treatment and spillover effects. For example, consider the causal effect of vaccine adoption on illness. With peer effects in vaccine adoption, vaccinated individuals tend to have more vaccinated direct and indirect social contacts, and a simple comparison of adopters and nonadopters may overstate vaccine

# GNNs for Network Confounding

utilize the model of approximate neighborhood interference (ANI) proposed by Leung (2022a), which posits that interference in the outcome stage decays with network path distance. Leung shows that ANI allows for endogenous peer effects but focuses on a setting with randomized assignment. In observational settings, it stands to reason that peer effects in selection may be a possibility. We theref

We provide conditions under which the doubly robust estimator is approximately normally distributed as the network size grows large. This type of result is well known for i.i.d. data (e.g. Farrell, 2018), but it is nontrivial to extend to our setting since we allow for a complex form of network dependence. For example, asymptotically linearizing the doubly robust estimator requires a new argument due to dependence, and application of an appropriate CLT requires verification of a high-level weak dependence condition under a nonparametric model with outcome and selection stages both governed by simultaneous-equations models. For inference, we utilize a network HAC estimator due to Kojevnikov et al. (2021) and propose a new bandwidth that adjusts for estimation error in the first-stage machine learners.

We substantiate the theory in a simulation study and empirical application to microfinance diffusion. The simulations demonstrate that the use of GNNs can substantially reduce bias relative to conventional choices of network controls even with shallow architectures. The empirical illustration revisits the microfinance diffusion application of He and Song (2024). We show how our estimands can capture complementary aspects of diffusion relative to their "average diffusion at the margin" measure. Our theoretical framework allows for more complex diffusion processes without requiring the econometrician to prespecify the maximum number of within-period rounds of diffusion. Finally, by including richer controls that account for network confounding, we find more attenuated diffusion effects.

## 1.2  Related Literature

There is a large literature on interference, much of which focuses on randomized control trials (e.g. Athey et al., 2018; Li and Wager, 2022; Toulis and Kao, 2013). We contribute to a growing recent literature on unconfoundedness, much of which operates in a partial interference setting where units are partitioned into disjoint groups with no interference across groups (e.g. Liu et al., 2019; Qu et al., 2022).

Studying a network interference setting, Veitch et al. (2019) propose to use "node embeddings" as network controls, which are learned functions of the graph. Since node embeddings can be obtained from a variety of methods, there remains the issue of justifying a particular choice of network controls. GNNs can be interpreted as a method of estimating node embeddings (see §3), and our behavioral model provides justification for their use. We defer to §2.1 a more detailed review of the literature

**i**'s $n$ $o$ and the elements of the same set of **K** $1$ as **i**'s $o$ neighbors. A unit **i**'s is $\mathbf{n}(\mathbf{i}, 1$, the number of neighbors.

# 2 Setup

Let $\mathcal{N}_{\mathbf{n}}$ $1, \ldots, \mathbf{n}$ be the set of units connected through the network $\boldsymbol{A}$. Each unit **i** $\mathcal{N}_{\mathbf{n}}$ is endowed with unobservables ( $_{\mathbf{i}}, _{\mathbf{i}}$ $\mathbb{R}^{\mathbf{d}_\varepsilon}$ $\mathbb{R}^{\mathbf{d}_\nu}$ and observables $\mathbf{X}_{\mathbf{i}}$ $\mathbb{R}^{\mathbf{d}_x}$. The model primitives determine outcomes and treatments according to

$$\mathbf{Y}_{\mathbf{i}} \quad \mathbf{g}_{\mathbf{n}}(\mathbf{i}, D, X, A, \varepsilon \quad \text{and} \quad \mathbf{D}_{\mathbf{i}} \quad \mathbf{h}_{\mathbf{n}}(\mathbf{i}, X, A, \nu, \tag{1}$$

respectively, where $\boldsymbol{X}$ $(\mathbf{X}_{\mathbf{i}} \, _{\mathbf{i}=1}^{\mathbf{n}}$ is the matrix with **i**th row equal to $\mathbf{X}_{\mathbf{i}}'$; $\boldsymbol{Y}$, $\boldsymbol{D}$, $\boldsymbol{\varepsilon}$, and $\boldsymbol{\nu}$ are similarly defined; and $(\mathbf{g}_{\mathbf{n}}, \mathbf{h}_{\mathbf{n}} \, _{\mathbf{n} \in \mathbb{N}}$ is a sequence of function pairs such that each $\mathbf{g}_{\mathbf{n}}($ has range $\mathbb{R}$ and $\mathbf{h}_{\mathbf{n}}($ has range $0, 1$. The econometrician observes $(\boldsymbol{Y}, \boldsymbol{D}, \boldsymbol{X}, \boldsymbol{A}$. Our analysis treats $(\boldsymbol{A}, \boldsymbol{X}, \boldsymbol{\varepsilon}, \boldsymbol{\nu}$ as random, but the asymptotic theory in §4 conditions on $(\boldsymbol{X}, \boldsymbol{A}$ to avoid imposing additional assumptions on its dependence structure.

We view the timing of the model as follows. First, nature draws the primitives $(\boldsymbol{A}, \boldsymbol{X}, \boldsymbol{\varepsilon}, \boldsymbol{\nu}$. Next, units select into treatment, potentially based on other units' decisions, and $\mathbf{h}_{\mathbf{n}}($ is the reduced-form outcome of that process. Finally, $\mathbf{g}_{\mathbf{n}}($ is the reduced form of the subsequent process that generates outcomes. Because $\mathbf{g}_{\mathbf{n}}($ and $\mathbf{h}_{\mathbf{n}}($ may depend on the primitives of all units, the setup allows $\mathbf{Y}_{\mathbf{i}}$ and $\mathbf{D}_{\mathbf{i}}$ to be outcomes of simultaneous-equations models with endogenous peer effects, as shown in the next examples.

**E a p** [1] (Linear-in-Means) Consider the outcome model

$$\mathbf{Y}_{\mathbf{i}} \quad + \quad \frac{\sum_{\mathbf{j}=1}^{\mathbf{n}} \mathbf{A}_{\mathbf{ij}} \mathbf{Y}_{\mathbf{j}}}{\sum_{\mathbf{j}=1}^{\mathbf{n}} \mathbf{A}_{\mathbf{ij}}} + \frac{\sum_{\mathbf{j}=1}^{\mathbf{n}} \mathbf{A}_{\mathbf{ij}} \mathbf{Z}_{\mathbf{j}}'}{\sum_{\mathbf{j}=1}^{\mathbf{n}} \mathbf{A}_{\mathbf{ij}}} + \mathbf{Z}_{\mathbf{j}}' + _{\mathbf{i}},$$

where $\mathbf{Z}_{\mathbf{i}}$ $(\mathbf{D}_{\mathbf{i}}, \mathbf{X}_{\mathbf{i}}'$ ' (Manski, 1993). The coefficient captures endogenous peer

denote the row-normalized adjacency matrix and $\mathbf{1}$ the $\mathbf{n}$-dimensional vector of ones, if $\boldsymbol{A}$ is connected, the reduced form of the model can be written in matrix form as

$$Y = \frac{}{1} \mathbf{1} + Z + \sum_{\mathbf{k}=0}^{\ddot{\mathbf{y}}} {}^{\mathbf{k}} \tilde{A}^{\mathbf{k}+1} Z + \sum_{\mathbf{k}=0}^{\ddot{\mathbf{y}}} {}^{\mathbf{k}} \tilde{A}^{\mathbf{k}} \varepsilon.$$

This characterizes $\mathbf{Y_i}$ as a function $\mathbf{g_n}(\mathbf{i}, D, X, A, \varepsilon)$.

**E a p'** (Binary Game) Consider the binary analog of Example 1 but for selection into treatment:

$$D_i = \mathbf{1} \left\{ + \frac{\sum_{j=1}^{n} A_{ij} D_j}{\sum_{j=1}^{n} A_{ij}} + \frac{\sum_{j=1}^{n} A_{ij} Z_j'}{\sum_{j=1}^{n} A_{ij}} + Z_i' + {}_{i} \geq 0 \right\}. \tag{2}$$

Unlike Example 1, there may exist multiple equilibria. The equilibrium selection mechanism is a reduced-form mapping from the primitives $(X, A, \nu)$ to outcomes $D$ and therefore characterizes $\mathbf{D_i}$ as a function $\mathbf{h_n}(\mathbf{i}, X, A, \nu)$. This formulation corresponds to a game of complete information. In a game of incomplete information, as modeled by Xu (2018) for instance, a unit $\mathbf{i}$'s information set is $(_{i}, X, A)$. Here an analog of (2) holds with each $\mathbf{D_j}$ replaced with $_{j}(X, A)$, the equilibrium belief that $\mathbf{D_j} = 1$. This characterizes $\mathbf{D_i}$ as a function $\mathbf{h_n}(\mathbf{i}, X, A, _{i})$.

**E a p'** (Diffusion) He and Song (2024) study the following two-period diffusion model. Let $\mathbf{D_i}$ denote $\mathbf{i}$'s decision to adopt microfinance in period 0 and $\mathbf{Y_i}$ its decision in period 1. Their equations (2.4) and (3.6) posit that

$$\mathbf{Y_i} = \mathbf{g_n}(D_{\mathcal{N}(\mathbf{i}, \mathbf{K})}, _{i}) \quad \text{and} \quad \mathbf{D_i} = \mathbf{1}\{W_i' \geq _{i}\},$$

where $\mathbf{W_i}$ is a known function of $(X, A)$ and $\mathbf{K}$ is the maximum distance that adoption decisions can diffuse through the network between periods 0 and 1. We provide a more detailed comparison of our models in §7.

Given specification (1), we define potential outcomes as

$$\mathbf{Y_i}(d) = \mathbf{g_n}(\mathbf{i}, d, X, A, \varepsilon).$$

Confounding may arise first because $\mathbf{Y_i}(d)$ is potentially correlated with $\mathbf{D_i}$ due to

the high-dimensional observables $(\boldsymbol{X}, \boldsymbol{A}$ and second because of dependence between unobservables that drive outcomes $\boldsymbol{\varepsilon}$ and those that drive selection $\boldsymbol{\nu}$. We restrict the second source of confounding.

**Assu pt on** [1] (Unconfoundedness) $o$ $n_{y}$ $\mathsf{n}$ $\mathbb{N}, \boldsymbol{\varepsilon}$ $\boldsymbol{\nu}$ $\boldsymbol{X}, \boldsymbol{A}$

As discussed below, unconfoundedness conditions used in the existing literature additionally limit the first source of confounding to known summary statistics of $(\boldsymbol{X}, \boldsymbol{A}$. Ours is analogous to standard formulations of unconfoundedness under SUTVA ( $_\mathsf{i}$
$_\mathsf{i}$ $\mathsf{X_i}$) since we do not impose further restrictions on the nature of observed confounding.

Because the econometrician only observes a single network, a large-sample theory requires restrictions on interference in order to obtain some form of weak dependence. We next specify a nonparametric model of decaying interference that accommodates the previous examples. For any $\mathsf{S}$ $\mathcal{N}_\mathsf{n}$, let $\boldsymbol{D_S}$ $(\mathsf{D_i}\ _\mathsf{i\in S}$, and similarly define $\boldsymbol{X_S}$ and other such submatrices. Let $\boldsymbol{A_S}$

counterfactual **s**-neighborhood outcomes. The error from approximating the observed outcome with the **s**-neighborhood counterfactual is bounded by $\phi_n(s)$, which decays with the neighborhood radius **s**. This formalizes the idea that $Y_i$ is primarily determined by units relatively proximate to **i**, so that the further a unit is from **i**, the less it influences **i**'s outcome. The second equation imposes the analogous requirement on $D_i$.

**Example continued.** For the linear-in-means model in Example 1, an argument similar to Proposition 1 of Leung (2022a) shows that (3) holds with $\sup_n \phi_n(s) \leq C|\lambda|^s$ for some $C > 0$. For the binary game in Example 2, an argument similar to Proposition 2 of Leung (2022a) establishes (4) with $\sup_n \phi_n(s)$ decaying at an exponential rate with **s**. Finally, for the He and Song (2024) diffusion model in Example 3, $Y_i$ only depends on $D$ through $D_{\mathcal{N}(i,K)}$, so (3) holds with $\phi_n(s) \leq c\mathbf{1}\{s < K\}$ for some universal constant **c**. In their empirical application, they use own covariates as controls, so $W_i = X_i$, in which case (4) holds with $\phi_n(s) = 0$ for all **s**.

## 2.1 Related Literature

The standard SUTVA model and unconfoundedness condition correspond to

$$Y_i = g(D_i, X_i, \varepsilon_i) \quad \text{and} \quad \varepsilon_i \perp D_i \mid X_i. \tag{5}$$

To generalize this setup to allow for network interference, the typical approach in the existing literature is as follows. Define

$$T_i = f_n(i, D, A) \quad \text{and} \quad W_i = q_n(i, X, A) \tag{6}$$

where $f_n(\cdot)$ and $q_n(\cdot)$ are known vector-valued functions. The $n$ (Manski, 2013) or $o \, u$ $n$ (Aronow and Samii, 2017,

model and unconfoundedness condition

$$Y_i = g(T_i, W_i, \varepsilon_i) \quad \text{and} \quad \varepsilon_i \perp T_i \mid W_i, \tag{7}$$

which is a direct generalization of (5) (Emmenegger et al., 2022; Forastiere et al., 2021; Ogburn et al., 2022). Here $T_i$ entirely summarizes interference while $W_i$ summarizes confounding.

Common examples of $T_i$ and $W_i$ are

$$T_i = \left( D_i, \sum_{j=1}^{n} A_{ij} D_j \right) \quad \text{and} \quad W_i = \left( X_i, \sum_{j=1}^{n} A_{ij}, \frac{\sum_{j=1}^{n} A_{ij} X_j}{\sum_{j=1}^{n} A_{ij}} \right). \tag{8}$$

This choice of $T_i$ implies that $Y_i$ depends on $D$ only through two statistics: own treatment and the number of treated neighbors. Variation in the first component identifies a direct treatment effect and variation in the second a spillover effect. Like most exposure mappings used in the literature, this only depends on $D_{\mathcal{N}(i,1)}$, so the outcome model (7) implies no interference beyond the 1-neighborhood. Likewise, this choice of $W_i$ implies no confounding beyond 1-neighborhood covariates.

More generally, one could restrict the outcome model to depend only on the $K$-neighborhood treatments $D_{\mathcal{N}(i,K)}$ for some fixed threshold $K$. As shown by Leung (2022a), this rules out economically interesting forms of interference such as endogenous peer effects, which motivates the ANI condition (3). Furthermore, (7) assumes the econometrician can correctly specify the summary statistic $T_i$ in the outcome model, which may be difficult to justify (Sävje, 2024).

Whereas Leung (2022a) and Sävje (2024) focus on randomized experiments, we study observational data on economic agents that choose to select into treatment. It then becomes important to specify a behavioral model rationalizing the choice of controls $W_i$. Sánchez-Becerra (2022) is the first to provide such a model. Under neighborhood interference (7) and an exposure mapping similar to (8), he shows that it is sufficient to set $W_i = X_i$, that is, to solely control for own covariates. Since much of the literature utilizes controls such as (8), this raises the question of what model of selection justifies their use or more broadly the use of "network controls" that depend more generally on $X$ and $A$.

Our model (1) provides an answer. The presence of complex interference in both the outcome and treatment stages induces selection on $(X, A)$, so that it is generally

with and without at least one treated neighbor, which captures a spillover effect. For $\mathbf{t} = (1, 0)$ and $\mathbf{t}' = (0, 0)$, it compares the average outcomes of treated and untreated units with no treated neighbors, which captures a treatment effect. For overlap, we need to exclude units with zero degree since a treated neighbor occurs with probability zero for such units. This is accomplished by choosing $\mathcal{M}_\mathbf{n}$ to be the subset of units whose degree $\mathbf{n}(\mathbf{i}, 1) = |\mathcal{N}(\mathbf{i}, 1)|$ lies in some desired set excluding zero. That is, choose some $\mathbf{\Gamma} \subseteq \mathbb{R}_+ \setminus \{0\}$ and

$$\mathcal{M}_\mathbf{n} = \left\{ \mathbf{i} : \left( \qquad \div \qquad \mid \right) \qquad \left( \quad \mid \quad \mid \quad \mid \right. \right.$$

can be a complex function of $(\boldsymbol{X}, \boldsymbol{A}, \boldsymbol{\nu})$, and both $\mathsf{T_i}$ and $\mathsf{Y_i}$ can be complex functions of $\boldsymbol{D}$, which makes it difficult to characterize the dependence structure necessary for the application of a central limit theorem without additional structure.

Identification results in Leung (2024) provide conditions under which $(\mathsf{t}, \mathsf{t}')$ has a causal interpretation. The focus of our paper is estimation, so we provide only a brief discussion here. Under the neighborhood interference model (7), $(\mathsf{t}, \mathsf{t}')$ has a transparent causal interpretation. In settings where (7) fails to hold, Leung (2024) shows that $(\mathsf{t}, \mathsf{t}')$ retains a causal interpretation under restrictions on interference either in potential outcomes or selection into treatment. For example, suppose treatment adoption follows a nonparametric game of incomplete information where $\nu_i$ is unit $\mathsf{i}$'s private information, so that $\mathsf{D_i} = \mathsf{h_n}(\mathsf{i}, \boldsymbol{X}, \boldsymbol{A}, \nu_i)$ (see Example 2). If private information is independent across units conditional on $(\boldsymbol{X}, \boldsymbol{A})$, as is typically assumed in structural analyses of the model (e.g. Lin and Vella, 2021; Xu, 2018), then by Theorem 1 of Leung (2024) $(\mathsf{t}, \mathsf{t}')$ can be written as a non-negatively weighted average of certain unit-level effects.

Returning to the vaccine adoption example in §1, recall that Adukx Aerelct

# GNNs for Network Confounding

where

$$\hat{\tau}_i(t, t') = \frac{\mathbb{1}[T_i = t](Y_i - \hat{\mu}_t(i, X, A))}{\hat{p}_t(i, X, A)} + \hat{\mu}_t(i, X, A)$$
$$- \frac{\mathbb{1}[T_i = t'](Y_i - \hat{\mu}_{t'}(i, X, A))}{\hat{p}_{t'}(i, X, A)} - \hat{\mu}_{t'}(i, X, A).$$

t = To estimate

## 3.1 Architecture

The standard GNN architecture consists of nested, paramete

GNNs for Network Confounding

aggregation no longer holds when the support of $\mathbf{X_i}$ is uncountable, so using multiple aggregators can result in more powerful architectures (Corso et al., 2020).

For an example of $\Gamma(\ )$, let $\boldsymbol{\mu}(\ ),\ (\ ),\Sigma(\ ),\min(\ )$, and $\max(\ )$ be respectively the mean, standard deviation, sum, min, and max functions, defined component-wise for multisets of vectors. Then setting $\Gamma(\ )\ \Gamma_1(\ )$ for

$$\Gamma_1(\ )\quad \begin{bmatrix} \boldsymbol{\mu}(\ ) & (\ ) & \Sigma(\ ) & \min(\ ) & \max(\ ) \end{bmatrix}$$

results in an architecture utilizing five aggregation functions.

The authors combine multiple aggregators with "degree scalars" that multiply each aggregation function by a function of the size of the multiset input $\mathbf{n}(\mathbf{i},1\ )$. The simplest example is the identity scalar, which maps any multiset to unity. This trivially multiplies each aggregation function in $\Gamma_1(\ )$ above, but it is useful to consider non-identity scalars. Let $|\ |$ be the function that takes as input a multiset and outputs its size. Corso et al. (2020) define logarithmic amplification and attenuation scalers

$$\mathbf{S}(\ ,\ )\quad \left[ \frac{\log(|\ |\ +1\ )}{\quad}\ \right]^{\ },\quad \frac{1}{\mathbf{n}}\sum_{\mathbf{i}=1}^{\ } \log\left(\sum_{\mathbf{j}=1}^{\ } \mathbf{A_{ij}}\ +1\right),\quad \in\ \{\ 1,1\ \}.$$

The choice of $\ $ defines whether the scalar "amplifies" ($\ 1$) or "attenuates" ($\ 1$) the aggregation function, and $\ 0$ is the identity scalar. The purpose of the logarithm is to prevent small changes in degree from amplifying gradients in an exponential manner with each successive GNN layer. Thus, an aggregation function that augments $\Gamma_1(\ )$ with logarithmic amplification and attenuation is

$$\Gamma_2(\ )\quad \mathbf{S}(\ )$$

# GNNs for Network Confounding

function ( :

$$\hat{\mathbf{f}}_{\mathrm{GNN}} \quad \underset{\mathbf{f} \in \mathcal{F}_{\mathrm{GNN}}(\mathbf{L})}{\mathrm{argmin}} \sum_{\mathbf{i} \in} \ddot{\mathbf{y}}$$

As discussed in §3.1, any $\mathbf{f} \in \mathcal{F}_{\mathrm{GNN}}(\mathbf{L})$ is invariant, so by using GNNs to esti-

that does not depend on **i**, so evaluating **i**'s propensity score is now only a matter of supplying the correct **i**-specific inputs ( $_i(X$ , $_i(A$ .

We close this section with a result demonstrating that invariance is an extremely weak requirement. In particular, it holds under minimal exchangeability conditions on the structural primitives.

**ropos t on** [1]     $u$   $o$    $o$    $n_l$  **n**   $\mathbb{N}$   $n$          $u$      $on$   ,

$$\mathbf{f_n}(\mathbf{i}, D, A \qquad \mathbf{f_n}( \ (\mathbf{i} \ , \ (D \ , \ (A \ ,$$
$$\mathbf{g_n}(\mathbf{i}, D, X, A, \varepsilon \qquad \mathbf{g_n}( \ (\mathbf{i} \ , \ (D \ , \ (X \ , \ (A \ , \ (\varepsilon \ , \qquad n$$

conditions imposed below. Define

$$
\psi_{t,t'}(i) \quad \frac{1\{T_i \quad t\}(Y_i \quad \mu_t(i, X, A)}{p_t(i, X, A)} \quad +\mu_t(i, X, A)
$$

$$
\frac{1\{T_i \quad t'\}(Y_i \quad \mu_{t'}(i, X, A)}{p_{t'}(i, X, A)} \quad \mu_{t'}(i, X, A) \qquad (t, t' ,
$$

whose average over $i \quad \mathcal{M}_n$ is the doubly robust moment condition, and let

$$
\sigma_n^2 \quad \text{Var} \quad \frac{1}{m_n} \sum_{i \in \mathcal{M}_n} \tilde{\psi}_{t,t'}(i \; X, A .
$$

**Assumption** (Moments) $\quad M \quad n \; p \quad 4 \; u \quad o \quad n_{\text{f}}$
$n \quad \mathbb{N}, \; i \quad \mathcal{N}_n, \quad n \; d \quad 0, 1^n, \; E||Y_i(d \; |^p \quad X, A \quad M$
$|\_,^- \; (0, 1 \quad u \quad \hat{p}_t(i, X, A , p_t(i, X, A \quad |\_,^- \quad n$

# GNNs for Network Confounding

The next assumption is used to show that $_{t,t'}(i\ _{i=1}^{n}$ is  -dependent (see Definition C.1, which is due to Kojevnikov et al., 2021) to apply m728(o)-313.551(a)-.823-5.1.(o)-313.551(a)-.823-5.1.(

Parts (a) and (b) are used to establish that $\phi_{t,t'}(i)_{i=1}^n$ is $\psi$-dependent. Part (b) is a Lipschitz condition that holds if potential outcomes are uniformly bounded. In particular we can take $\Lambda_n(s) = 2M$ where $M$ is the uniform bound on the ranges

strengthens Assumption 5(b) but only mildly since we nonparametrically estimate both nuisance functions. Since it does not require uniform convergence, it is more readily verifiable for machine learning estimators. Part (c) is Assumption 4.1(iii) of Kojevnikov et al. (2021). Parts (d)–(f) correspond to Assumptions 7(b)–(d) of Leung (2022a), which are used to characterize the bias of the variance estimator. We discuss verification of (c)–(f) in §B.2; the derivations there show that (f) is closely related to (c).

or $n$ ~$_{t,t^1}($i $j$ $n$ $(t, t'$ $n$ $n$ $on$ $o$ $_{t,t^1}($i $_i(t, t'$

$E|Y_i$ $T_i$ $t, X,$

case $\hat{}^2$ would be asymptotically conservative. This can be formalized under additional weak dependence conditions on the superpopulation as in §A of Leung (2022b).

# 5   Approximate Sparsity

As discussed in §3, the number of layers **L** in a GNN determines its                          ,
the neighborhood $(\boldsymbol{X}_{\mathcal{N}(\mathbf{i},\mathbf{L})}$

$(W_i)_{i=1}^n$ and $\nu \equiv (\nu_i)_{i=1}^n$, define

$$V_i(W, \nu; \cdot) \equiv \cdot + \frac{\sum_{j=1}^n A_{ij} W_j}{\sum_{j=1}^n A_{ij}} + \frac{\sum_{j=1}^n A_{ij} X_j}{\sum_{j=1}^n A_{ij}} + X_j + \nu_i + \frac{\sum_{j=1}^n A_{ij} \nu_j}{\sum_{j=1}^n A_{ij}}$$

where $\cdot \equiv (\cdot, \cdot, \cdot, \cdot)$. We generate $(Y_i)_{i=1}^n$ from the linear-in-means model, where $Y_i \equiv V_i(Y, \varepsilon; \cdot_y)$ and $\cdot_y \equiv (0.5, 0.8, 10, \cdot 1)$. We generate $(D_i)_{i=1}^n$ according to Example 2, so that $D_i \equiv 1\{V_i(D, \nu; \cdot_d) \geq 0\}$ with $\cdot_d \equiv (\cdot 0.5, 1.5, 1, \cdot 1)$. The equilibrium selection mechanism is myopic best-response dynamics starting from the initial condition $(D_i^0)_{i=1}^n$ for $D_i^0 \equiv 1\{V_i(0, \nu; \cdot_d) \geq 0\}$.

The design induces a greater degree of dependence than what our assumptions allow. The error term $\nu_i + \sum_{j=1}^n A_{ij} \nu_j \sum_{j=1}^n A_{ij}$ is not conditionally independent across units unlike what Assumption 6(a) requires. Also, back-of-the-envelope calculations indicate that peer effects are sufficiently large in magnitude that Assumption 6(d) is violated.

We use the estimand in Example 5 whose true value is zero. About 57 percent of units select into treatment, so the effective sample size used to estimate the outcome regressions is around $n\cdot 2$ since $E[Y_i \mid T_i \equiv t, X, A]$ is estimated only with observations for which $T_i \equiv t$. We report results for $n \in \{1000, 2000, 4000\}$.

## 6.2 Nonparametric Estimators

The GNNs use the PNA architecture in Example 9 with aggregator $\Gamma_2(\cdot)$ defined in the example and $L \equiv 1$,

(2022) and Forastiere et al. (2021). For these, we estimate the nuisance functions using GLMs (logistic and linear regression) with polynomial sieves of order 1, 2, or 3. Recall that a GNN with $\mathbf{L}$ 1 corresponds to a receptive field that only encompasses the ego's 1-neighborhood. This is the same as the implied receptive field of the GLM estimators.

Table 1: Simulation results for random geometric graph

| | $L = 1$ | | | $L = 2$ | | | $L = 3$ | | |
|---|---|---|---|---|---|---|---|---|---|
| $n$ | | | | | | | | | |
| $H$ | | | | | | | | | |
| $\hat{\tau}(1,0)$ | | | | | | | | | |
| CI | | | | | | | | | |
| CI | | | | | | | | | |
| $W\ \hat{\tau}(1,0)$ | | | | | | | | | |
| $W$ CI | | | | | | | | | |
| $W$ | | | | | | | | | |
| CI | | | | | | | | | |

si u ations,  he esti  and is  $(1,0) = 0$,   treated  $\approx$

A choice of $L = 2$ is not unusual in the literature. Zhou et al. (2021) compute the prediction error of GNNs on several different datasets with $L = 2, 4, 8, \ldots$ and find that $L = 2$ has the best performance across several architectures. The fact that GNN performance often fails to improve (and indeed can worsen) with larger $L$ is well known in the GNN literature, and we survey different expl

study or later.

## 7.1 Comparison with He and Song (2024)

For the first (second) choice, $(\mathbf{t}, \mathbf{t}'$ measures the effect of going from 0 to 1 (more than 1) adopting neighbor(s). This sheds light on a different dimension of diffusion relative to the ADM. Whereas the ADM measures how many others are affected by the ego's adoption, our estimands quantify the effect of having multiple adopting neighbors on the ego's adoption. We find below that having multiple adopting neighbors has a much larger effect than having only one.

As previously stated, He and Song (2024) define the treatment in two ways. One is a binary indicator for having a leader in the household, the idea being that all leaders were initially informed about microfinance and told to spread the word. However, not all leaders adopted in the first trimester, which perhaps

# GNNs for Network Confounding

To compute the estimates, we concatenate the village networks into a single adjacency matrix of size $\mathbf{n}$ 4413. For the GNN and GLM estimates, we trim observations with propensity scores outside of $[0.01, 0.99]$. Standard errors are obtained from the network HAC variance estimator defined in §2.3.

Table 3: Exposure Mapping $\mathsf{T}_i^{(1)}$

| | ADM | GNN | | | GLM | | |
|---|---|---|---|---|---|---|---|
| | | Layer | Layer | Layer | Order | Order | Order |
| **Leader case** | | | | | | | |
| $\mathbf{G}_{ee}$ | | | | | | | |
| $\mathbf{G}_{sc}$ | | | | | | | |
| $\mathbf{G}_{all}$ | | | | | | | |
| **Leader adopter case** | | | | | | | |
| $\mathbf{G}_{ee}$ | | | | | | | |
| $\mathbf{G}_{sc}$ | | | | | | | |
| $\mathbf{G}_{all}$ | | | | | | | |
| **Adopter case** | | | | | | | |
| $\mathbf{G}_{ee}$ | | | | | | | |
| $\mathbf{G}_{sc}$ | | | | | | | |
| $\mathbf{G}_{all}$ | | | | | | | |

$\mathbf{n}$ = 4413

the leader case.

First consider the estimand using $T_i^{(1)}$, which contrasts microfinance adoption rates for units with 1 versus 0 initially adopting neighbors. The GNN results are consistent across $L$. For the leader case, we obtain precise zeros for almost all estimates, including the ADM. For the leader-adopter case, the GNN estimates are substantially smaller in magnitude than the ADM with an effect size of at most 10 percentage points compared to the smallest ADM estimate of 20 percentage points. This may be attributed to the use of richer network controls. For the adopter case, the contrast is even starker. Our estimates are an order of magnitude smaller than the corresponding ADM estimates. The GLM estimates are typically slightly smaller than the GNN estimates except for the order-3 polynomials, which are outliers in terms of magnitude.

The estimand using $T_i^{(2)}$ contrasts units with 2 + versus 0 initially adopting neighbors. The estimates for the leader case are similar to those of $T_i^{(1)}$. We find sizeable effects for the leader-adopter case, almost of the same order as ADM, but the robustness of the result is partly tempered by the large amount of trimming discussed below. The adopter case sees estimates of around 20 percentage points, whereas the ADM is double or triple that. Once again the GLM estimates are often slightly smaller relative to the GNN estimates, except for the order-3 polynomials.

The number of observations trimmed for $T_i^{(1)}$ is negligible in the leader and adopter cases. In the leader-adopter case, more units are trimmed since fewer units are leader-adopters, but trimming never drops more than 200 observations. The story is quite different for $T_i^{(2)}$, as reported in Table 5. The first three columns of the table report the number of observations for which the number of initially adopting neighbors

$$N_i \qquad \sum_{j=1}^{n} A_{ij} D_j$$

in the lasso literature, which posit that a high-dimensional regression function is well-approximated by a function of a relatively small number of covariates.

# A    Additional Results on GNNs

Primitive conditions for Assumption 5 appear to be beyond the scope of the existing GNN literature, but we provide some potentially useful intermediate results. Consider the problem of establishing a rate of convergence for the propensity score:

$$\frac{1}{m_n} \sum_{i \in \mathcal{M}_n} \left( \hat{p}_t(i, X, A) - \breve{p}_t(i, X, A) \right)^2 = o_p(n^{-1/2}).$$

Under network approximate sparsity (Definition 1), the problem simplifies to showing

$$\frac{1}{m_n} \sum_{i \in \mathcal{M}_n} \left( \hat{p}_t(i, X, A) - \breve{p}_t(i, X_{\mathcal{N}(i,L)}, A_{\mathcal{N}(i,L)}) \right)^2 = o_p(n^{-1/2}). \tag{A.1}$$

Since $\hat{p}_t(i, X, A)$ is an $L$-layer GNN, which only uses information from $(X_{\mathcal{N}(i,L)}, A_{\mathcal{N}(i,L)})$, this should well approximate $\breve{p}_t(i, X_{\mathcal{N}(i,L)}, A_{\mathcal{N}(i,L)})$ under appropriate conditions, so (A.1) should be more feasible to verify directly.

Farrell et al. (2021) provide a bound analogous to (A.1) for MLPs, which, were it applicable to our setting, would be of the form

$$\frac{1}{n} \sum_{i=1}^{n} \left( \hat{p}_t(i, X, A) - \breve{p}_t(i, X_{\mathcal{N}(i,L)}, A_{\mathcal{N}(i,L)}) \right)^2 \leq C \left( \frac{WL \log R}{n} \log n + \frac{\log \log n + \zeta}{n} \right) + \epsilon^2. \tag{A.2}$$

with probability at least $1 - e^{-\zeta}$. Here $W$ is the number of GNN parameters, $C$ is a constant that does not depend on $n$, $R$ depends on the architecture through the number of hidden neurons, and $\epsilon$ is the function approximation error, a measure of the ability of the neural network to approximate any function in a desired class. Establishing a corresponding result for GNNs requires an analog of Lemma 6 of Farrell et al. (2021), which is a bound on the pseudo-dimension of the GNN class, and concentration inequalities for $\psi$-dependent data. Jegelka (2022) surveys the few available complexity and generalization bounds for GNNs. These are not sufficiently general for our setup and only apply to settings where the sample consists of many independent networks.

To use a bound of the form (A.2) to verify Assumption 5, we require knowledge of how varies with key aspects of the architecture, such as $\mathbf{W}, \mathbf{R}, \mathbf{L}, \mathbf{n}$. As a first step toward obtaining such a result, it is necessary to characterize the function class that GNNs can approximate. Our next result, which draws heavily from existing results in the GNN literature, shows that an additional shape restriction on the function class beyond invariance (§3.3) is required.

## A.1 WL Function Class

MLPs can approximate any measurable function (Hornik et al., 1989), so given the discussion in §3.3, a natural question is whether GNNs can approximate any measurable, $n$ $n$ function of graph-structured inputs. In other words, is it enough to require invariance (and regularity conditions), or are stronger restrictions on the function class necessary? For reasons related to the graph isomorphism problem, it turns out stronger restrictions are necessary. We next motivate the need for such restrictions and then state our function approximation result.

Chen et al. (2019) show that, for a function class such as GNNs to approximate any invariant function, some element of the class must be able to separate any pair of non-isomorphic graphs. By "separate," we mean that for any non-isomorphic "labeled graphs" $(\boldsymbol{X}, \boldsymbol{A}$, $(\boldsymbol{X}', \boldsymbol{A}'$, the function $\mathbf{f}$ satisfies $\mathbf{f}(\boldsymbol{X}, \boldsymbol{A} \neq \mathbf{f}(\boldsymbol{X}', \boldsymbol{A}'$. Hence, a function with separating power of this sort solves the graph isomorphism problem, a problem for which no known polynomial-time solution exists (Kobler et al., 2012; Morris et al., 2021). Since GNNs can be computed in polynomial time, this suggests that approximating any invariant function is too demanding of a requirement.

To define the subclass of invariant functions that GNNs can approximate, we need to take a detour and discuss graph isomorphism tests. The subclass will be defined by a weaker graph separation criterion than solving the graph isomorphism problem, in particular one defined by the $W$ $n$ $W$ . This is a (generally imperfect) test for graph isomorphism on which almost all practical graph isomorphism solvers are based (Morris et al., 2021).

Given a labeled graph $(\boldsymbol{X}, \boldsymbol{A}$, the WL test outputs a graph coloring (a vector of

labels for each unit) according to the following recursive procedure, whose definition

$\mathbf{H}^2$ such that

$$(\mathbf{h}, \mathbf{h}' \quad (\mathcal{F} \quad \text{if and only if} \quad \mathbf{f}(\mathbf{h} \quad \mathbf{f}(\mathbf{h}' \quad \text{for all} \quad \mathbf{f} \quad \mathcal{F}.$$

For any two sets of functions $\mathcal{E}, \mathcal{F}$ with domain $\mathbf{H}$, we say that $\mathcal{E}$ is $\quad o$ $\quad n \quad \mathcal{F}$ if $(\mathcal{F} \quad (\mathcal{E}$ .

This is essentially Definition 2 of Azizian and Lelarge (2021). Intuitively, if $\mathcal{E}$ is at most as separating as $\mathcal{F}$, the latter is more complex in the sense that some function in $\mathcal{F}$ can separate weakly more elements of $\mathbf{H}$ than any function in $\mathcal{E}$.

Let $\mathbf{f}_{L,\mathbf{L}}$ denote the function of $(\boldsymbol{X}, \boldsymbol{A}$ with range $\Sigma^{\mathbf{n}}$ that outputs the vector of node colorings from the WL test run for $\mathbf{L}$ iterations. Let $\mathcal{C}(\mathbf{H}$ be the set of continuous functions with domain $\mathbf{H}$. For any $\mathbf{L} \quad \mathbb{N}$, define the $W \quad \iota n \quad on$

$$\mathcal{F}_{\mathrm{L}}(\mathbf{L} \quad \mathbf{f} \quad \mathcal{C}(\mathbf{H} : (\mathbf{f}_{L,\mathbf{L}} \quad (\mathbf{f} \quad .$$

This is the set of continuous functions of $(\boldsymbol{X}, \boldsymbol{A}$ that are at most as separating as the WL test with $\mathbf{L}$ iterations.

The next result says that $\mathbf{p_t}($ and $\boldsymbol{\mu_t}($ can be approximated by

In other words, any function in the class $\mathcal{F}_L(\mathbf{L}$ can be approximated by $\mathbf{L}$-layer GNNs in $\mathcal{F}_{\text{GNN}}(\mathbf{L}$. The result is a consequence of a Stone-Weierstrauss theorem due to Azizian and Lelarge (2021) and a version of the Morris et al. (2019) and Xu et al. (2018) result on the equivalent separation power of GNNs and the WL test. The proof is given below.

The result is essentially Theorem 4 of Azizian and Lelarge (

$\mathsf{T}$ in place of $\mathsf{L}$) and add an additional MLP layer that implements their equation (26). Here we use the additional MLP layer added to the output of our architecture (see the paragraph prior to the statement of Theorem A.1). In particular, for any $\mathsf{f} \in \mathcal{F}_{\text{GNN}} (\mathsf{L}$, consider the mapping

$$(X, A \underbrace{\Big( \overset{\sim}{\overset{\ddot{y}}{\mathsf{f}}} (\mathsf{i}, X, A ,\ldots, \overset{\ddot{y}}{\mathsf{f}} (\mathsf{i}, X, A}_{\mathsf{n} \text{ ti es}} \Big) \in \mathbb{R}^{\mathsf{n}}$$

(their (26) in our notation) corresponds to adding a linear output layer $\mathsf{L} +1$ that is implementable by an MLP of the form $_{\mathsf{L}+1}(\mathsf{h}_{\mathsf{i}}^{(\mathsf{L})}, \mathsf{h}_{\mathsf{j}}^{(\mathsf{L})} : \mathsf{j} \in \mathcal{N}_{\mathsf{n}}$ . The mapping remains an element of $\mathcal{F}_{\text{GNN}} (\mathsf{L}$, which completes the argument for (A.5).

By Theorems VIII.1 and VIII.4 of Grohe (2021), which use finiteness of the support of $\mathsf{X}_{\mathsf{i}}$, $(\mathcal{F}_{\text{GNN}} (\mathsf{L}$ $(\mathsf{f}$ $_{\mathsf{L},\mathsf{L}}$ . That is, $\mathsf{L}$-layer GNNs have the same separation power as the WL test run for $\mathsf{L}$ iterations. ∎

## A.2 Disadvantages of Depth

The receptive field is the main consideration when selecting $\mathsf{L}$, but Theorem A.1 provides a second consideration, which is imposing a weaker implicit shape restriction. It shows that, for GNNs to approximate a target function well, the target must satisfy a shape restriction stronger than invariance, namely that it is at most as separating as the WL test with $\mathsf{L}$ iterations. The larger the choice of $\mathsf{L}$, the weaker the shape restriction imposed. However, there are several reasons why shallow architectures remain preferable.

**ow r turns to ` pt** A natural question is how many iterations are required for the WL test to converge for a given graph, which corresponds to the choice of $\mathsf{L}$ for which the shape restriction is weakest. Unfortunately, the answer is not generally known, being determined by the topology of the input graph in a complex manner. However, there is a range of results bounding the number of iterations required for convergence. For instance, Kiefer and McKay (2020) construct graphs for which the WL test requires $\mathsf{n}$ 1 iterations to converge, so such graphs require $\mathsf{n}$ 1 layers to obtain the weakest shape restriction. This makes the estimation problem extremely high-dimensional, requiring substantially more layers than what is typically required

than the standard architecture (13) (Dwivedi et al., 2022). These disadvantages may partly explain the common use in practice of the standard architecture with few layers.

# B  Verifying §8 Assumptions

Leung (2022a), §A, verifies analogs of Assumptions 6(d) and 7(c) from an older working paper version of Kojevnikov et al. (2021). This section repeats the exercise for Assumptions 6(d) and 7(c) and (d). We assume throughout that max $_n(s\ 2\ ,\ _n(s$ exp( $c(1\quad 4\ p\ ^{-1}s$ for some $c\quad 0$ and $p$ in Assumption 4(a). As in Leung (2022a), we say a sequence of networks exhibits polynomial neighborhood growth if

$$\sup_n \max_{i \in \mathcal{N}_n} |\mathcal{N}_A(i, s\ |\quad Cs^d$$

for some $C\quad 0, d > 1$

and $D'_{\mathbf{B}}$ $(D'_{\mathbf{j}})_{\mathbf{j}\in\mathbf{B}}$ for any $\mathbf{B} \subseteq \mathcal{N}_{\mathbf{n}}$. Using the first equality of (C.3),

$$p_t(\mathbf{i}, X, A) = P\left(D'_{\mathbf{i}} + (D_{\mathbf{i}} - D'_{\mathbf{i}}) \mid a, b, V'_{\mathbf{i}} + (V_{\mathbf{i}} - V'_{\mathbf{i}}) \mid , X, A\right)$$
$$= P\left(D'_{\mathbf{i}} \mid a, b +, V'_{\mathbf{i}} \mid , + X, A\right)$$
$$+ \underbrace{P(\ldots D_{\mathbf{i}} - D'_{\mathbf{i}} \ldots X, A \ldots + P(\ldots V_{\mathbf{i}} - V'_{\mathbf{i}} \ldots X, A)}_{\mathbf{R}_0}.$$

By (C.3), the right-hand side equals

$$P\left(D'_{\mathbf{i}} \mid a, b\right)$$

Combining (C.6) and (C.7) and using the law of iterated expectations,

$$|\mathbf{p_t}(\mathbf{i}, \boldsymbol{X}, \boldsymbol{A}) - \mathbf{p_t}(\mathbf{i}, \boldsymbol{X}_{\mathcal{N}(\mathbf{i}, \mathbf{r}_\lambda(\mathbf{s}+1))}, \boldsymbol{A}_{\mathcal{N}(\mathbf{i}, \mathbf{r}_\lambda(\mathbf{s}+1))})| \leq \delta_\mathbf{n}(\mathbf{s}+1) + 2\mathbf{R}_0.$$

**Proo o (C.2)** Noting that $\boldsymbol{\mu_t}(\mathbf{i}, \boldsymbol{X}, \boldsymbol{A}) = E[Y_i \mathbf{1_i}(t) \mid \boldsymbol{X}, \boldsymbol{A}] \mathbf{p_t}(\mathbf{i}, \boldsymbol{X}, \boldsymbol{A})$, we first bound the numerator. For $\mathbf{B} \in \mathcal{N}(\mathbf{i}, \mathbf{s})$, define $Y_i' \equiv \mathbf{g}_{\mathbf{n}(\mathbf{i},\mathbf{s})}(\mathbf{i}, D'_\mathbf{B}, \boldsymbol{X_\mathbf{B}}, \boldsymbol{A_\mathbf{B}}, \varepsilon_\mathbf{B})$. By Lemma C.2,

$$|E[Y_i \mathbf{1_i}(t) \mid \boldsymbol{X}, \boldsymbol{A}] - E[Y_i' \mathbf{1_i}(t) \mid \boldsymbol{X}, \boldsymbol{A}]| \leq \delta_\mathbf{n}(\mathbf{s}) + \Lambda_\mathbf{n}(\mathbf{i}, \mathbf{s}) \eta(\mathbf{i}, \mathbf{s}) \delta_\mathbf{n}(\mathbf{s})$$

where, using Assumption 4(b),

$$|\mathbf{R}_1| \quad {}_{\underset{\sim}{n}}(2\mathbf{s} +2 \ {}_{n}(\mathbf{s} +\Lambda_{\mathbf{n}}(\mathbf{i},\mathbf{s} \ n(\mathbf{i},\mathbf{s} \ {}_{\underset{\sim}{n}}(\mathbf{s} +\mathbf{C}'(1 +n(\mathbf{i},1 \ {}_{n}(\mathbf{s} ,$$

$$|\mathbf{R}_2| \quad \mathbf{C} \ {}_{n}(2\mathbf{s} +(1 +\mathbf{n}(\mathbf{i},1 \ {}_{n}(2\mathbf{s} \ 1 , \quad \text{and}$$

$$|\mathbf{R}_3| \quad \mathbf{C}''(|\mathbf{R}_1| +|\mathbf{R}_2|$$

for some universal $\mathbf{C}''$ 0. Substituting $\mathbf{s}$ 2 for $\mathbf{s}$ yields the result. ∎

$$\mathbf{a} \ \mathbf{C} \ ' \quad n \ \mathbf{B_i} \quad \mathcal{N}(\mathbf{i},\mathbf{s} \ , \mathbf{D_j'} \quad \mathbf{h_{n(j,s)}}(\mathbf{j}, X_{\mathbf{B}_j}, A_{\mathbf{B}_j}, \nu_{\mathbf{B}_j} \ , D_{\mathbf{B}_i}' \quad (\mathbf{D_j'}$$

for $|\mathbf{R}_1| \quad {}_n(\mathbf{s} \quad +\Lambda_\mathbf{n}(\mathbf{i}, \mathbf{s} \; \mathbf{n}(\mathbf{i}, \mathbf{s} \quad {}_n(\mathbf{s}$ . ∎

$\mathbf{a}\ \mathbf{C} \qquad n \quad \mathbf{Y}'_\mathbf{i}, \mathbf{D}'_\mathbf{i} \qquad n \qquad n \quad \mathbf{1}_\mathbf{i}(\mathbf{t}\ ' \quad 1\ \mathbf{D}'_\mathbf{i} \quad \mathbf{d}, \quad \sum_{\mathbf{j}=1}^\mathbf{n} \mathbf{A}_{\mathbf{ij}} \mathbf{D}'_\mathbf{j}$

$\Delta \qquad n \quad A \; \boldsymbol{u} \qquad on \quad , \quad , \quad , \quad n \qquad , \qquad \mathbf{C} \qquad 0 \; \boldsymbol{u} \qquad o \quad n_{\!\boldsymbol{y}}$

$\mathbf{n} \quad \mathbb{N}, \mathbf{i} \quad \mathcal{N}_\mathbf{n}, \quad n \quad \mathbf{s} \succ 0,$

$$\mathbf{E}|\mathbf{Y}_\mathbf{i}|\mathbf{1}_\mathbf{i}(\mathbf{t} \quad \mathbf{1}_\mathbf{i}(\mathbf{t}\ '|\ X, A \quad \mathbf{C}\ (1\ +\mathbf{n}(\mathbf{i}, 1 \quad {}_\mathbf{n}(\mathbf{s} \ .$$

**Proof.** Recall the definition of $\mathbf{a}, \mathbf{b}, \; , \; ,$ prior to (C.3). Define $\mathbf{V}_\mathbf{i} \quad \sum_{\mathbf{j}=1}^\mathbf{n} \mathbf{A}_{\mathbf{ij}} \mathbf{D}_\mathbf{j}$, $\mathbf{V}'_\mathbf{i} \quad \sum_{\mathbf{j}=1}^\mathbf{n} \mathbf{A}_{\mathbf{ij}} \mathbf{D}'_\mathbf{j}$, and $\mathcal{C} \quad |\mathbf{D}_\mathbf{i} \quad \mathbf{D}'_\mathbf{i}| \quad , |\mathbf{V}_\mathbf{i} \quad \mathbf{V}'_\mathbf{i}| \qquad$ . Then

$$\mathbf{E}|\mathbf{Y}_\mathbf{i}|\mathbf{1}_\mathbf{i}(\mathbf{t} \quad \mathbf{1}_\mathbf{i}(\mathbf{t}\ '|\ X \quad x, A \quad a$$

$$\mathbf{E}|\mathbf{Y}_\mathbf{i}|\mathbf{1}_\mathbf{i}(\mathbf{t} \quad \mathbf{1}_\mathbf{i}(\mathbf{t}\ '|\ \mathcal{C}, X \quad x, A \quad a \quad +\mathbf{C}\ P(\mathcal{C}^\mathbf{c} \quad X \quad x, A \quad a \quad \text{(C.11)}$$

for some universal $\mathbf{C} \quad 0$ by Assumptions 1 and 4(a). By Assumption 3,

$$\mathbf{1}_\mathbf{i}(\mathbf{t} \quad 1\ \mathbf{D}_\mathbf{i} \quad |\mathbf{a}, \mathbf{b}\ , \; \mathbf{V}_\mathbf{i} \quad |\ , \qquad \text{and} \quad \mathbf{1}_\mathbf{i}(\mathbf{t}\ ' \quad 1\ \mathbf{D}'_\mathbf{i} \quad |\mathbf{a}, \mathbf{b}\ , \; \mathbf{V}'_\mathbf{i} \quad |\ , \quad .$$

Under event $\mathcal{C}$,

$$1\ \mathbf{D}_\mathbf{i} \quad |\mathbf{a}, \mathbf{b}\ , \; \mathbf{V}_\mathbf{i} \quad |\ , \quad {}^\big( \quad 1\ \mathbf{D}'_\mathbf{i} \quad +(\mathbf{D}_\mathbf{i} \quad \mathbf{D}'_\mathbf{i} \quad |\mathbf{a}, \mathbf{b}\ , \; \mathbf{V}'_\mathbf{i} \quad +(\mathbf{V}_\mathbf{i} \quad \mathbf{V}'_\mathbf{i} \quad {}_\vee \; {}_\mathbf{i}| \; \mathbf{D}'$$

$$1\ \mathbf{D}'_\mathbf{i} \quad |\mathbf{a} \quad , \mathbf{b}\ + \; , \; \mathbf{V}'_\mathbf{i} \quad |$$

Lipschitz functions on $\mathbb{R}^{\mathbf{d}}$, Lip($\mathbf{f}$ be the Lipschitz constant of $\mathbf{f}$ $\mathcal{L}_{\mathbf{d}}$, and

$$\mathcal{P}_{\mathbf{n}}(\mathbf{h}, \mathbf{h}'; \mathbf{s} \qquad (\mathbf{H}, \mathbf{H}' : \mathbf{H}, \mathbf{H}' \quad \mathcal{N}_{\mathbf{n}}, \ \mathbf{H} \quad \mathbf{h}, \ \mathbf{H}' \quad \mathbf{h}', \ _{A}(\mathbf{H}, \mathbf{H}' \ \succ \mathbf{s} \ .$$

**D** $\mathbf{n}, \mathbf{s}$ **on** $\mathbf{C}^{\mathbf{1}}$ A triangular array $\mathbf{Z_i} \ _{\mathbf{i=1}}^{\mathbf{n}}$ is *on* *on* $\mathbf{s}$ *n* *n* *n* $\mathcal{F}_{\mathbf{n}}$ if there exist $\mathbf{C}$ (0, and an $\mathcal{F}_{\mathbf{n}}$-measurable sequence $_{\mathbf{n}}(\mathbf{s} \ _{\mathbf{s},\mathbf{n}\in\mathbb{N}}$ with $_{\mathbf{n}}(0 \qquad 1$ for all $\mathbf{n}$ such that

$$|\text{Cov}(\mathbf{f}(\mathbf{Z_H} \ , \mathbf{f}'(\mathbf{Z_{H^1}} \ | \quad \mathbf{Chh}'(\|\mathbf{f}\|_\infty \ + \text{Lip}(\mathbf{f} \ (\|\mathbf{f}'\|_\infty \ + \text{Lip}(\mathbf{f}' \ _{\mathbf{n}}(\mathbf{s} \quad \text{a.s.} \quad (\text{C.12})$$

for all $\mathbf{n}, \mathbf{h}, \mathbf{h}' \quad \mathbb{N}; \mathbf{s} \quad 0; \mathbf{f} \quad \mathcal{L}_{\mathbf{h}}; \mathbf{f}' \quad \mathcal{L}_{\mathbf{h^1}};$ and $(\mathbf{H}, \mathbf{H}' \quad \mathcal{P}_{\mathbf{n}}(\mathbf{h}, \mathbf{h}'; \mathbf{s} \ .$ We call $_{\mathbf{n}}(\mathbf{s}$ the $n$ $n$ $o$ $n$ of $\mathbf{Z_i} \ _{\mathbf{i=1}}^{\mathbf{n}}.$

$\qquad$ **a C** $\quad n \quad A \quad \mathbf{u} \quad on$ , , , $\quad n \quad , \quad n \quad n \quad , \quad o \quad n\mathbf{s}$
$\mathbf{t}, \mathbf{t}' \quad \mathcal{T}, \quad _{\mathbf{t},\mathbf{t^1}}(\mathbf{i} \ _{\mathbf{i=1}}^{\mathbf{n}} \quad on \quad on \quad \mathbf{s} \qquad n \quad n \quad n \ (\mathbf{X}, \mathbf{A} \qquad n \quad on$
$\qquad n \quad n \quad o \quad n \quad _{\mathbf{n}}(\mathbf{s} \qquad n \quad n \ (15)$

**Proof.** Let $\mathcal{F}_{\mathbf{n}}$ be the -algebra generated by $(\mathbf{X}, \mathbf{A} \ , (\mathbf{h}, \mathbf{h}' \quad \mathbb{N} \quad \mathbb{N}, (\mathbf{f}, \mathbf{f}' \ \mathcal{L}_{\mathbf{h}} \quad \mathcal{L}_{\mathbf{h^1}}, \mathbf{s} \quad 0,$ and $(\mathbf{H}, \mathbf{H}' \quad \mathcal{P}_{\mathbf{n}}(\mathbf{h}, \mathbf{h}'; \mathbf{s} \ .$ Define $\mathbf{Z_i} \qquad _{\mathbf{t},\mathbf{t^1}}(\mathbf{i} \ , \mathbf{Z_H} \qquad (\mathbf{Z_i} \ _{\mathbf{i}\in\mathbf{H}}, \ \mathbf{f}(\mathbf{Z_H} \ , \qquad \mathbf{f}'(\mathbf{Z_{H^1}} \ ,$ and

$$\mathbf{D}_{\mathbf{i}}^{(\mathbf{s})} \qquad \mathbf{h_{n(i,s)}}(\mathbf{i}, \mathbf{X}_{\mathcal{N}(\mathbf{i},\mathbf{s})}, \mathbf{A}_{\mathcal{N}(\mathbf{i},\mathbf{s})}, \boldsymbol{\nu}_{\mathcal{N}(\mathbf{i},\mathbf{s})} \ .$$

For $\mathbf{D}_{\mathcal{N}(\mathbf{i},\mathbf{s^1})}^{(\mathbf{s})} \quad (\mathbf{D}_{\mathbf{j}}^{(\mathbf{s})} \ _{\mathbf{j}\in\mathcal{N}(\mathbf{i},\mathbf{s^1})}),$ let

$$\mathbf{1}_{\mathbf{i}}^{(\mathbf{s})}(\mathbf{t} \qquad 1 \quad \mathbf{f_{n(i,s/2)}}(\mathbf{i}, \mathbf{D}_{\mathcal{N}(\mathbf{i},\mathbf{s}/2)}^{(\mathbf{s}/2)}, \mathbf{A}_{\mathcal{N}(\mathbf{i},\mathbf{s}/2)} \qquad \mathbf{t} \ ,$$

$$\mathbf{Y}_{\mathbf{i}}^{(\mathbf{s})} \qquad \mathbf{g_{n(i,s/2)}}(\mathbf{i}, \mathbf{D}_{\mathcal{N}(\mathbf{i},\mathbf{s}/2)}^{(\mathbf{s}/2)}, \mathbf{X}_{\mathcal{N}(\mathbf{i},\mathbf{s}/2)}, \mathbf{A}_{\mathcal{N}(\mathbf{i},\mathbf{s}/2)}, \boldsymbol{\varepsilon}_{\mathcal{N}(\mathbf{i},\mathbf{s}/2)} \ ,$$

$$\mathbf{Z}_{\mathbf{i}}^{(\mathbf{s})} \qquad \frac{\mathbf{1}_{\mathbf{i}}^{(\mathbf{s})}(\mathbf{t} \ (\mathbf{Y}_{\mathbf{i}}^{(\mathbf{s})} \quad \boldsymbol{\mu}_{\mathbf{t}}(\mathbf{i}, \mathbf{X}, \mathbf{A}}{\mathbf{p}_{\mathbf{t}}(\mathbf{i}, \mathbf{X}, \mathbf{A}} \quad + \boldsymbol{\mu}_{\mathbf{t}}(\mathbf{i}, \mathbf{X}, \mathbf{A}$$

$$\frac{\mathbf{1}_{\mathbf{i}}^{(\mathbf{s})}(\mathbf{t}' \ (\mathbf{Y}_{\mathbf{i}}^{(\mathbf{s})} \quad \boldsymbol{\mu}_{\mathbf{t^1}}(\mathbf{i}, \mathbf{X}, \mathbf{A}}{\mathbf{p}_{\mathbf{t^1}}(\mathbf{i}, \mathbf{X}, \mathbf{A}} \qquad \boldsymbol{\mu}_{\mathbf{t^1}}(\mathbf{i}, \mathbf{X}, \mathbf{A} \qquad _{\mathbf{i}}(\mathbf{t}, \mathbf{t}' \ .$$

Finally, let $^{(\mathbf{s})} \quad \mathbf{f}((\mathbf{Z}_{\mathbf{i}}^{(\mathbf{s})} \ _{\mathbf{i}\in\mathbf{H}}$ and $^{(\mathbf{s})} \quad \mathbf{f}'((\mathbf{Z}_{\mathbf{i}}^{(\mathbf{s})} \ _{\mathbf{i}\in\mathbf{H^1}} \ .$

By Assumption 6(a), $(\mathbf{Z_i}^{(s/2,\ )}\ _{\mathbf{i} \in \mathbf{H}}\quad (\mathbf{Z_j}^{(s/2,\ )}\ _{\mathbf{j} \in \mathbf{H}^1}\quad \mathcal{F_n}$, so

$$|\text{Cov}(\ ,\quad \mathcal{F_n}|\quad |\text{Cov}(\quad^{(s/2)},\quad \mathcal{F_n}|\ +|\text{Cov}(\quad^{(s/2)},\quad^{(s/2)}\ \mathcal{F_n}|$$
$$2\|\mathbf{f'}\|_\infty \mathbf{E}\|\quad^{(s/2)}|\ \mathcal{F_n}\ +2\|\mathbf{f}\|_\infty \mathbf{E}\|\quad^{(s/2)}|\ \mathcal{F_n}$$
$$2\ \mathbf{h}\|\mathbf{f'}\|_\infty \text{Lip}(\mathbf{f}\quad +\mathbf{h'}\|\mathbf{f}\|_\infty \text{Lip}(\mathbf{f'}\quad \max_{\mathbf{i} \in \mathcal{N}}$$

$\mathcal{L}_{\mathbf{h}}$ $\mathcal{L}_{\mathbf{h}^1}$, $\mathbf{s}$ $0$, $(\mathbf{H}, \mathbf{H}'$ $\mathcal{P}_{\mathbf{n}}(\mathbf{h}, \mathbf{h}'; \mathbf{s}$ ,

$$\mathsf{Y}_{\mathbf{i}}^{(\mathbf{s})} \quad \mathfrak{g}_{\mathbf{n}(\mathbf{i},\mathbf{s})}(\mathbf{i}, D_{\mathcal{N}(\mathbf{i},\mathbf{s})}, X_{\mathcal{N}(\mathbf{i},\mathbf{s})}, A_{\mathcal{N}(\mathbf{i},\mathbf{s})}, \varepsilon_{\mathcal{N}(\mathbf{i},\mathbf{s})} ,$$

$\mathbf{f}((\mathsf{Y}_{\mathbf{i}} {}_{\mathbf{i} \in \mathbf{H}}$ , $\mathbf{f}'((\mathsf{Y}_{\mathbf{i}} {}_{\mathbf{i} \in \mathbf{H}^1}$ , ${}^{(\mathbf{s})}$ $\mathbf{f}((\mathsf{Y}_{\mathbf{i}}^{(\mathbf{s})} {}_{\mathbf{i} \in \mathbf{H}}$ , and ${}^{(\mathbf{s})}$ $\mathbf{f}'((\mathsf{Y}_{\mathbf{i}}^{(\mathbf{s})} {}_{\mathbf{i} \in \mathbf{H}^1}$ . By Assumption 6(a),

$$|\mathrm{Cov}( , \quad \mathcal{F}'_{\mathbf{n}} | \quad |\mathrm{Cov}( {}^{(\mathbf{s}/2)}, \quad \mathcal{F}'_{\mathbf{n}} | + |\mathrm{Cov}( {}^{(\mathbf{s}/2)}, \quad {}^{(\mathbf{s}/2)} \quad \mathcal{F}'_{\mathbf{n}} |$$
$$2\|\mathbf{f}'\|_{\infty}\mathbf{E}\| \quad {}^{(\mathbf{s}/2)}| \quad \mathcal{F}'_{\mathbf{n}} + 2\|\mathbf{f}\|_{\infty}\mathbf{E}\| \quad {}^{(\mathbf{s}/2)}| \quad \mathcal{F}'_{\mathbf{n}}$$
$$2 \mathbf{h}\|\mathbf{f}'\|_{\infty}\mathrm{Lip}(\mathbf{f} + \mathbf{h}'\|\mathbf{f}\|_{\infty}\mathrm{Lip}(\mathbf{f}' \max_{\mathbf{i} \in \mathcal{N}_n} \mathbf{E}\|\mathsf{Y}_{\mathbf{i}} \quad \mathsf{Y}_{\mathbf{i}}^{(\mathbf{s}/2)}| \quad \mathcal{F}'_{\mathbf{n}}$$
$$2 \mathbf{h}\|\mathbf{f}'\|_{\infty}\mathrm{Lip}(\mathbf{f} + \mathbf{h}'\|\mathbf{f}\|_{\infty}\mathrm{Lip}(\mathbf{f}' \quad {}_{\mathbf{n}}(\mathbf{s} 2 ,$$

the last line using Assumption 2. Given -dependence, the claim follows from Corollary A.2 of Kojevnikov et al. (

$\hat{p}_t(i, X, A$ . For some universal constants $C, C'$ $0$, $E|R_{1t}^2$ equals

$$\frac{1}{m_n} \ddot{y} \ddot{y}_{i \in \mathcal{M}_n \, j \in \mathcal{M}_n} E'' E \,|\, (Y_i \quad \mu_i \, (Y_j \quad \mu_j \quad D, X, A \quad 1_i(t \, 1_j(t \, ($$

$\hat{p}_t(i, X, A$ , and

$$\Delta_i(t \quad (\hat{\mu}_t(i \quad \mu_t(i \quad {}^{p_t(i} \quad 1_i(t$$

# References

**A on an ` E Ya av**, "On the Bottleneck of Graph Neural Networks and its Practical Implications," in "International Conference on Learning Representations" 2021.

**Aronow an ` C a**, "Estimating Average Causal Effects Under General Interference, with Application to a Social Network Experiment," *Ann o A* , 2017, (4), 1912–1947.

**At D E s an ` G I ns**, "Exact **p**-Values for Network Interference," *o n o A n A o on*, 2018, (521), 230–240.

**G I ns J t r an ` E unro**, "Using Wasserstein Generative Adversarial Networks for the Design of Monte Carlo Simulations," *o n o ono* , 2021.

/BPC1 ID EIQ 0Nu.601000cmB,Tc/R20211.9552Tf 1001101.76489.96Tm (,)-390.203(G)0.743337(.)-380.153(5.6 0.2.4m):
b25(a)-2(n)875((L)-0.63102-0.73312,75(o)-2.37)1.38865(o)-2(u)2.92127(r)0.861963..0916(t)S(m)-3-2.nm

**Bronst n** , "Do We Need Deep Graph Neural Networks?," `https //towardsdatascience com/do-we-need-deep-graph-neural-networks -be 2d3ecpcp` 2020. Accessed: 2022-07-02.

**J Bruna Co n an ov** , "Geometric Deep Learning: Grids, Groups, Graphs, Geodesics, and Gauges," $X$ $n$ $X$ , 2021.

**C n Z ar C n an J Bruna**, "On the Equivalence Between Graph Isomorphism Testing and Function Approximation with GNNs," in "Advances in Neural Information Processing Systems," Vol. 32 2019.

**C rno u ov D C tv ov D r r E Du o C Hans n w an J o ns**, "Double/Debiased Machine Learning for Treatment and Structural Parameters," *ono ou n* , 2018, , C1–C68.

**Corso G Cava D B n , an ov** , "Principal Neighbourhood Aggregation for Graph Nets," in "Advances in Neural Information Processing Systems," Vol. 33 2020, pp. 13260–13271.

**D Groot** , "Reaching a Consensus," *ou n o A n A o on*, 1974, (345), 118–121.

**D ra a F C Gar a J no Donovan an A an B rra**, "Identifying Causal Effects in Experiments with Spillovers and Non-compliance," *ou n o ono* , 2023, *5* (2), 1589–1624.

**Dw v C Jos aur nt Y B n o an X Br sson**, "Benchmarking Graph Neural Networks," $X$ $n$ $X$ , 2022.

**E n r C po n an B ann**, "Treatment Effect Estimation from Observational Network Data Using Augmented Inverse Probability Weighting and Machine Learning," $X$ $n$ $X$ 5 , 2022.

**Farr** , "Robust Inference on Average Treatment Effects with Possibly more Covariates than Observations," $X$ $n$ $X$ , 2018.

**an an sra**, "Deep Neural Networks for Estimation and Inference," *ono* , 2021, (1), 181–213.

**F      an ` J     nss n**, "Fast Graph Representation Learning with PyTorch Geometric," *X       n   X    0 0      *, 2019.

**Forast r    E A rd ` an ` F     a  **, "Identification and Estimation of Treatment and Interference Effects in Observational Studies on Networks," *ou n  o    A     n        A o    on*, 2021,    (534), 901–918.

**Gro       **, "The Logic of Graph Neural Networks," in "2021 36th Annual ACM/IEEE Symposium on Logic in Computer Science" IEEE 2021, pp. 1–17.

**H   X an `      on **, "Measuring Diffusion over a Large Network," *      o    o   no       u     o   o   n  *, 2024.

**Horn         t n  o      an ` H      t **, "Multilayer Feedforward Networks are Universal Approximators," *   u          o  *, 1989,

# GNNs for Network Confounding

ヽ  r    an ` B    a , "The Iteration Number of Colour Refinement," in "47th International Colloquium on Automata, Languages, and Prog

"Causal Inference for Spatial Treatments," $X$ $n$

$X$

**Zop** , "1-WL Expressiveness Is (Almost) All You Need," $X$ $n$ $X$ $5$ , 2022.