

A Sensitivity Analysis Framework for Causal Inference Under Interference

Matvey Ortyashov AmirEmad Ghassami

Department of Mathematics and Statistics, Boston University

Abstract

In many applications of causal inference, the treatment received by one unit may influence the outcome of another, a phenomenon referred to as interference. Although there are several frameworks for conducting causal inference in the presence of interference, practitioners often lack the data necessary to adjust for its effects. In this paper, we propose a weighting-based sensitivity analysis framework that can be used to assess the systematic bias arising from ignoring interference. Unlike most of the existing literature, we allow for the presence of unmeasured confounding, and show that the combination of interference and unmeasured confounding is a notable challenge to causal inference. We also study a third factor contributing to systematic bias: lack of transportability. Our framework enables practitioners to assess the impact of these three issues simultaneously through several easily interpretable sensitivity parameters that can reflect a wide range of intuitions about the data.

Keywords— causal inference, sensitivity analysis, interference, unmeasured confounding, transportability, spillover effects

1 Introduction

Researchers in causal inference frequently invoke the assumption of no interference, which states that the treatment received by any unit should only affect the outcome of that unit (Rubin, 1980). Due to the widespread adoption of this assumption, most literature ignores the presence of *spillover effects*, which occur when one unit’s treatment may “spill over,” and have a causal impact on (i.e., interfere with) another’s outcome. Ignoring spillover effects can result in systematic bias, a form of bias that cannot be mitigated by increasing the sample size. Concerns about the bias introduced by spillover effects are particularly relevant in fields where interference is common, such as the study of infectious diseases (VanderWeele et al., 2015).

While the role of ignored interference in creating systematic bias is widely recognized, few authors have explicitly explored the structure, the components, and the extent of this bias. As a result, notwithstanding a proliferation of recent research concerned with proper estimation strategies in the presence of observed interference (Liu et al., 2016; Forastiere et al., 2021; Lee et al., 2023; McNealis et al., 2024; Tortú et al., 2024; Papadogeorgou and Samanta, 2023; Forastiere et al., 2024), the subject of sensitivity analysis for ignored interference has hitherto been underexplored. Despite the limited amount of literature concerning the topic, sensitivity analyses can play an important role in the presence of spillover effects. In many settings, practitioners may incorrectly assume that there are no spillover effects within a given set of data or, due to cost or privacy concerns, may lack the ability to learn the structure of the interference present in the data. A sensitivity analysis framework would enable researchers facing these issues to evaluate how ignoring interference impacts their estimates and gauge the strength of interference necessary to change their findings.

In this paper, we study the structure of the bias that arises due to failing to adjust for various complexities introduced by the presence of interference. Unlike the majority of the existing literature on interference (with some exceptions, such as Chen et al. (2025), Wu and Franks (2025), and Khot et al. (2025)), we allow for the presence of unobserved confounders. Importantly, we will demonstrate that if interference is present, the impacts of unmeasured confounding on the bias differ from those that exist without interference, making the combination of unmeasured confounding and ignored interference a distinct challenge to causal inference. Using our bias decomposition, we propose a sensitivity analysis framework that parametrizes the bias through several easily interpretable sensitivity parameters that reflect a practitioner’s beliefs about factors such as the strength of interference effects and unmeasured confounding. Our analyses also account for another data complexity that can occur alongside interference: the presence of undefined potential outcomes.

Furthermore, we consider settings where the researcher aims to *transport* their causal findings from a reference domain to a target domain. For this case, we investigate a third source of bias: lack of transportability. Causal transportability, alongside several concepts closely related to it, such as external validity, generalizability, data fusion, and transfer learning, has garnered significant attention over the past several years (Bareinboim and Pearl, 2016; Mitra et al., 2022; Huang, 2024; Degtiar and Rose, 2023; Colnet et al., 2024; Vuong et al., 2025). As discussed by Buchanan et al. (2023) as well as Bhadra and Schweinberger (2025), though there is a wide range of considerations about transporting causal effects in the presence of interference (one of which is that transportability requires similar patterns of interference between populations, as noted by Hernán and VanderWeele (2011)), few authors have explicitly examined them. We study how our causal estimands change between different populations, and extend our sensitivity analysis to account for lack of transportability.

To the best of our knowledge, our paper is the first to develop a comprehensive sensitivity analysis framework that allows practitioners to simultaneously account for the bias that comes from ignored interference, unmeasured confounding, and lack of transportability. Although Forastiere et al. (2021) discuss the bias of an estimator that does not adjust for interference and unmeasured confounding, they do not present sensitivity parameters or offer a structured approach that can be employed by a practitioner wishing to calculate the bias of their naive estimator. To date, only VanderWeele et al. (2015) have developed a formal sensitivity analysis framework for unmeasured confounding in the presence of interference, though their methodology does not account for ignored interference or lack of transportability. Unlike their work, our method not only incorporates the aforementioned data complexities, but also does not require strong assumptions about the outcome generating process, the distribution of the unmeasured confounder, or the selection bias. However, if one is willing to make certain assumptions, we demonstrate that our bias decomposition can be significantly simplified. Our bias decomposition generalizes previous weighting-based results, such as those of Shen et al. (2011), and is related to the approaches of Hong et al. (2021) and Huang (2024).

The rest of this paper is organized as follows. Section 2 describes the problem setting. In Section 3, we provide the main theoretical results for bias arising from unmeasured confounding and interference. Section 4 extends the results of Section 3 to the task of causal transport. Section 5 extends the results of Sections 3 and 4 to undefined potential outcomes. We conclude in Section 6. All of the proofs are provided in the Appendix.

2 Problem Description

Consider a setting in which unit $i \in \mathcal{I}$ receives treatment $A_i \in \{0, 1\}$ and has an observed outcome Y_i . Without interference, unit i ’s outcome depends only on unit i ’s treatment; however, in the presence of interference, Y_i may also depend on the treatments received by other units (Rubin, 1974, 1990). We refer to those units whose treatments affect Y_i as the *neighbors* of unit i . We let \mathcal{N}_i denote the set of indices of unit i ’s neighbors, and let the vector $\mathbf{A}_{\mathcal{N}_i}$ represent their treatments. We do not impose any restrictions on which units can belong to \mathcal{N}_i . Thus, our approach allows for more flexibility than the common setting of

partial interference, in which units are partitioned into blocks, and interference can only occur within blocks, not across them (Sobel, 2006). We adopt the potential outcomes framework, where units' outcomes under different realizations of the treatments are considered distinct random variables (Rubin, 1974). Specifically, if A_i is set to a and $\mathbf{A}_{\mathcal{N}_i}$ is set to \mathbf{a} , we write the potential outcome of unit i as $Y_i^{(a,\mathbf{a})}$. Furthermore, we let \mathbf{X}_i and \mathbf{U}_i denote the observed and unobserved pre-treatment covariates respectively. For the context of causal transport, we use variable S_i to demarcate units in the reference population ($S_i = 1$) from those in the target population ($S_i = 2$). We assume that our observational data only contains units with $S_i = 1$, and that the two populations may differ in terms of covariates (as long as those covariates do not affect the outcome) and treatment generating mechanisms.

Under neighborhood interference, a unit with n_i neighbors may have 2^{n_i+1} distinct potential outcomes. Nevertheless, in a majority of applications, many of these potential outcomes would equal each other. Thus, we assume that there exists a scalar-valued function $g_i(\mathbf{A}_{\mathcal{N}_i}) = G_i$ which acts as a summary of the treatments received by unit i 's neighbors and reduces the number of potential outcomes we need to consider. This assumption is formalized below.

Assumption 1 (Outcome Exposure Mapping). *For all $i \in \mathcal{I}$, there exists a function $g_i : \{0, 1\}^{n_i} \rightarrow \mathbb{R}$ such that for any two different sets of treatments received by a unit's neighbors, \mathbf{a}_1 and \mathbf{a}_2 , we have that*

$$Y_i^{(a,\mathbf{a}_1)} = Y_i^{(a,\mathbf{a}_2)}$$

for $a \in \{0, 1\}$, as long as $g_i(\mathbf{a}_1) = g_i(\mathbf{a}_2)$. We assume, without loss of generality, that for all $i \in \mathcal{I}$, $g_i(\mathbf{A}_{\mathcal{N}_i}) \in \{0, 1, \dots, g_{\max}\}$ (i.e., $g_i(\mathbf{A}_{\mathcal{N}_i})$ is discrete). Moreover, we assume that g_i only depends on i through the dimension of the domain, and that $g_{\max} < \max_i 2^{n_i}$.

Effectively, Assumption 1 states that if $A_i = a$ and $G_i = g$, the potential outcome can be written as $Y_i^{(a,g)}$. The function g_i is a particular case of the *exposure mapping* described by Aronow and Samii (2017). We call function g_i the *outcome exposure mapping*, G_i the *neighborhood treatment*, and A_i the *personal treatment* of unit i .¹ Throughout the rest of this paper, we drop the subscript i and refer to generic instances of the aforementioned variables without the subscript.

Next, we propose a version of the consistency assumption, which relates the unit's observed outcome, Y , to its potential outcomes.

Assumption 2 (Consistency under Interference). *The potential outcome $Y^{(a)}$ satisfies $Y^{(a)} = Y^{(a,G)} = \sum_{g=0}^{g_{\max}} \mathbb{I}(G = g) \cdot Y^{(a,g)}$. Moreover, the observed outcome Y satisfies $Y = \sum_{a=0}^1 \mathbb{I}(A = a) \cdot Y^{(a)}$.*

Taken together, the two components of Assumption 2 imply that if a unit receives personal treatment a and neighborhood treatment g , its outcome equals $Y^{(a,g)}$, i.e., $Y = \sum_{a=0}^1 \sum_{g=0}^{g_{\max}} \mathbb{I}(A = a, G = g) \cdot Y^{(a,g)}$. In addition to Assumptions 1 and 2, we impose one ancillary constraint to develop tractable bias decompositions in Sections 3 and 4.

Assumption 3 (Well-Defined Potential Outcomes). *The potential outcomes $Y^{(a)}$ and $Y^{(a,g)}$ for $a \in \{0, 1\}$, $g \in \{0, 1, \dots, g_{\max}\}$ are well-defined for every unit.*

In general, Assumption 3 does not necessarily hold in the presence of interference, and thus, is not required for our bias decomposition. Nevertheless, we introduce Assumption 3 in order to avoid unnecessary complexities in notation that do not employ any new methodological tools and do not offer any additional insights into our main results. In Section 5, we explore cases in which Assumption 3 may be violated, and extend our sensitivity analysis to such settings.

¹Note that the exposure mapping assumption is not the only way to define potential outcomes in the presence of interference: alternative definitions, typically based on a given treatment allocation strategy, have also been posited in the literature (Tchetgen Tchetgen and VanderWeele, 2012; Liu et al., 2016; Lee et al., 2023).

2.1 Estimands

In this paper, we have two primary parameters of interest: the *natural average main effect*, denoted by ϕ_1 , and the *transported natural average main effect*, denoted by ϕ_2 . In order to describe these parameters, we begin by considering the *controlled individual main effect*, $\tau(g)$:

$$\tau(g) := Y^{(1,g)} - Y^{(0,g)}.$$

$\tau(g)$ represents the individual causal effect of changing A when G is set to g , and can be viewed as a direct controlled effect (Robins and Greenland, 1992; VanderWeele, 2011). In addition to $\tau(g)$, for each unit, we can define the *natural individual main effect*, κ , which measures the causal impact of changing A while G equals the actual value of neighborhood treatment received by the unit:

$$\kappa := Y^{(1,G)} - Y^{(0,G)} = Y^{(1)} - Y^{(0)} = \sum_{g=0}^{g_{\max}} \mathbb{I}(G = g) \cdot \tau(g).$$

We can then define both of our parameters of interest as averages of κ , one taken over the reference population ($S = 1$), the other over the target population ($S = 2$):

$$\phi_s := \mathbb{E}[\underbrace{Y^{(1,G)} - Y^{(0,G)}}_{\kappa} | S = s] = \mathbb{E}[Y^{(1)} - Y^{(0)} | S = s] = \sum_{g=0}^{g_{\max}} \mathbb{E}[\mathbb{I}(G = g) \cdot \tau(g) | S = s].$$

As the average difference between $Y^{(1,G)}$ and $Y^{(0,G)}$, ϕ_s can be interpreted as measuring the impact of changing A from 0 to 1 while G is allowed to vary as it naturally does in the population where $S = s$.

Remark 1. Note that ϕ_2 is a generalization of a causal estimand commonly presented in the literature (Forastiere et al., 2021; Liu et al., 2023; Tortú et al., 2024; Zigler et al., 2025) and given by

$$\theta = \sum_{g=0}^{g_{\max}} p(G = g | S = 1) \cdot \mathbb{E}[\tau(g) | S = 1] = \sum_{g=0}^{g_{\max}} \mathbb{E}[\mathbb{I}(G = g) | S = 1] \cdot \mathbb{E}[\tau(g) | S = 1].$$

The formula above highlights that θ is a special case of ϕ_2 for populations where: (i) $\tau(g) \perp\!\!\!\perp G | S = 2$, (ii) $G \perp\!\!\!\perp S$, and (iii) $\mathbb{E}[\tau(g) | S = 1] = \mathbb{E}[\tau(g) | S = 2]$. Therefore, if $G \perp\!\!\!\perp \tau(g) | S = 1$, $\theta = \phi_1$. Otherwise, while ϕ_1 can be viewed as the average impact of changing A if G was allowed to vary as it does in the reference population, θ represents the average impact of changing A if the neighborhood treatment in the reference population was randomized, with probabilities given by $p(G = g | S = 1)$.

Remark 2. Though we have hitherto discussed parameters that measure the causal effect of changing the personal treatment A , we can also define **spillover effects**, which measure the causal impact of changing the neighborhood treatment G . Similarly to the natural individual main effect, we can define the **natural individual spillover effect** at $A = a$, denoted by $\gamma(a)$:

$$\gamma(a) = Y^{(a,G)} - Y^{(a,0)} = Y^{(a)} - Y^{(a,0)} = \sum_{g=0}^{g_{\max}} \left\{ \mathbb{I}(G = g) \cdot Y^{(a,g)} \right\} - Y^{(a,0)}.$$

$\gamma(a)$ measures the causal spillover effect for any given unit by comparing the potential outcome where $A = a$ and G is set to the actual neighborhood treatment received by that unit with the potential outcome where $A = a$ and G is set to 0. Note that the **individual total effect** can be defined in terms of the individual natural main and spillover effects:

$$\text{Individual Total Effect} = Y^{(1,G)} - Y^{(0,0)} = Y^{(1)} - Y^{(0,0)} = \kappa + \gamma(0)$$

While we introduce spillover effects in this section, our focus in this work remains developing a sensitivity analysis framework for ϕ_1 and ϕ_2 . Thus, we only discuss spillover effects insofar as they figure into our bias decompositions for the average natural main effects.

In settings where all of the relevant variables are observed, the requirements of conditional exchangeability and positivity are typically posited to allow for identification. However, since we wish to develop a framework that allows for the existence of unobserved confounders and lack of transportability, we will adjust these assumptions. Before doing so, we differentiate between several different categories of covariates in our setting. Namely, we write $\mathbf{X} = \mathbf{X}_{AY} \cup \mathbf{X}_{GY} \cup \mathbf{X}_{AG} \cup \mathbf{X}_{AS} \cup \mathbf{X}_{GS}$ and $\mathbf{U} = \mathbf{U}_{AY} \cup \mathbf{U}_{GY} \cup \mathbf{U}_{AG} \cup \mathbf{U}_{AS} \cup \mathbf{U}_{GS}$ where variables in $\{\mathbf{X}_{AY}, \mathbf{U}_{AY}\}$ are common causes of A and Y , variables in $\{\mathbf{X}_{GY}, \mathbf{U}_{GY}\}$ are common causes of \mathbf{A}_N and Y , variables in $\{\mathbf{X}_{AG}, \mathbf{U}_{AG}\}$ are common causes of A and \mathbf{A}_N , variables in $\{\mathbf{X}_{AS}, \mathbf{U}_{AS}\}$ are common causes of A and S , and variables in $\{\mathbf{X}_{GS}, \mathbf{U}_{GS}\}$ are common causes of \mathbf{A}_N and S . Using this notation, we posit the following modified versions of positivity and conditional exchangeability under interference.

Assumption 4 (Positivity). *For $a \in \{0, 1\}$ and $s \in \{1, 2\}$ we have, almost surely in (\mathbf{X}, \mathbf{U})*

$$p(A = a|S = s, \mathbf{X}, \mathbf{U}) > 0 \text{ and } p(S = s|\mathbf{X}, \mathbf{U}) > 0.$$

Assumption 5 (Weak Conditional Exchangeability). *For $a \in \{0, 1\}$, we have*

$$Y^{(a)} \perp\!\!\!\perp A|S, \tilde{\mathbf{X}}, \tilde{\mathbf{U}}$$

where $\tilde{\mathbf{X}} = \mathbf{X}_{AY} \cup \mathbf{X}_{AG} \cup \mathbf{X}_{AS}$ and $\tilde{\mathbf{U}} = \mathbf{U}_{AY} \cup \mathbf{U}_{AG} \cup \mathbf{U}_{AS}$.

Remark 3. *Assumption 5 highlights three additional challenges introduced by the combination of interference and lack of transportability.*

- *Under interference, one needs to condition on common causes of A and \mathbf{A}_N , even if these variables are not common causes of A and Y , and do not need to be adjusted for in the absence of interference.*
- *Since the reference and target populations may differ in terms of their treatment generating mechanisms, S also serves as a common cause of A and \mathbf{A}_N , and thus, must be adjusted for.*
- *Because the reference and target populations may differ in terms of their covariates (namely, causes of A and \mathbf{A}_N), there is a collider structure that necessitates adjusting for common causes of A and S .*

Note that because Assumption 5 allows for the existence of a set of unobserved covariates that ensures independence, it is significantly weaker than the version of exchangeability typically presented in the absence of unmeasured confounding.

Remark 4. *The combination of Assumptions 2 and 5 implicitly precludes most forms of dependence between A and \mathbf{A}_N . In Assumption 2, we assume that $Y^{(a)}$ equals $Y^{(a,G)}$, as opposed to $Y^{(a,G^{(a)})}$, a condition which would typically be violated by the existence of a direct causal pathway from A to \mathbf{A}_N . In Assumption 5, we assume that the common causes of A and \mathbf{A}_N are sufficient to close the confounding pathway between A and Y through \mathbf{A}_N ; however, this would not hold true in most settings in which direct causal pathways from \mathbf{A}_N to A exist.*

Under Assumptions 1-5, one could identify ϕ_s if they were to observe $\tilde{\mathbf{U}}$. One possible approach is to use inverse probability weighting (IPW), as in Proposition 1.

Proposition 1 (IPW Identification Formula). *Under Assumptions 1 through 5, ϕ_s , can be identified as*

$$\phi_s = \mathbb{E} \left[\frac{\mathbb{I}(A = 1) \cdot Y}{p(A = 1|S = s, \tilde{\mathbf{X}}, \tilde{\mathbf{U}})} | S = s \right] - \mathbb{E} \left[\frac{\mathbb{I}(A = 0) \cdot Y}{p(A = 0|S = s, \tilde{\mathbf{X}}, \tilde{\mathbf{U}})} | S = s \right].$$

Assumptions 1-5 are sufficiently strong to develop a bias decomposition for ϕ_1 . However, reasoning about ϕ_2 is inherently more challenging, as the practitioner does not observe any units from the target population. Thus, in order to develop a tractable bias decomposition for ϕ_2 , we assume that only the treatment generating mechanism, not the outcome generating mechanism, may differ between the two populations, and that the covariates which are the causes of Y do not differ between the two populations. Assumption 6 formalizes this requirement.

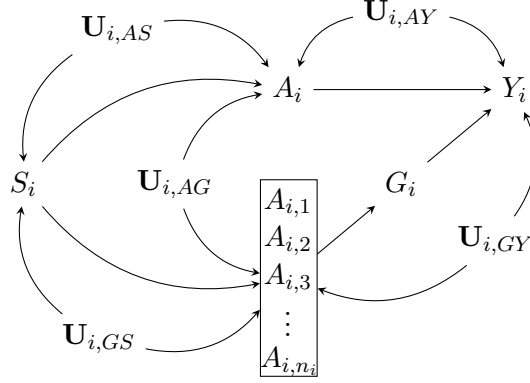


Figure 1: A directed acyclic graph consistent with our assumptions. \mathbf{X}_i is omitted for legibility, $\{A_{i,1}, A_{i,2}, \dots, A_{i,n_i}\} = \mathbf{A}_{\mathcal{N}_i}$. Note that this graph is meant for general illustration and there are other systems which satisfy our assumptions (e.g., the five categories of the covariates do not have to be disjoint).

Assumption 6 (Transportability Assumptions). *For any $a \in \{0, 1\}$ and $g \in \{0, 1, \dots, g_{max}\}$, we have*

- $Y^{(a,g)} \perp\!\!\!\perp \mathbf{A}_{\mathcal{N}} | \mathbf{X}_{GY}, \mathbf{U}_{GY}$,
- S only has a direct causal effect on A and G , and
- Y and S do not share common causes, or, in other words, $Y^{(a,g)} \perp\!\!\!\perp S$.

One graphical model consistent with our assumptions is depicted in Figure 1. While we introduce ϕ_1 and ϕ_2 together in this section, to provide sufficient detail for both bias decompositions, we investigate them separately in Sections 3 and 4 respectively.

3 Sensitivity Analysis for Unmeasured Confounding and Unobserved Interference

As discussed in Section 1, we examine two limitations to conducting causal inference in the presence of interference: failing to adjust for interference and failing to account for unobserved confounders. Both of these issues can introduce systematic bias into estimation; therefore, we examine both in tandem. To formalize, we consider the case where a practitioner tries to estimate the following functional, which, in the absence of interference and unmeasured confounding, equals ϕ_1 :

$$\psi = \mathbb{E} \left[\frac{\mathbb{I}(A = 1) \cdot Y}{p(A = 1 | S = 1, \mathbf{X}_{AY})} \middle| S = 1 \right] - \mathbb{E} \left[\frac{\mathbb{I}(A = 0) \cdot Y}{p(A = 0 | S = 1, \mathbf{X}_{AY})} \middle| S = 1 \right].$$

A typical estimator of ψ is the IPW estimator (Hirano and Imbens, 2001), which is given by

$$\hat{\psi} = \mathbb{E}_n \left[\frac{\mathbb{I}(A = 1, S = 1) \cdot Y}{\hat{p}(A = 1 | S = 1, \mathbf{X}_{AY}) \cdot \hat{p}(S = 1)} \right] - \mathbb{E}_n \left[\frac{\mathbb{I}(A = 0, S = 1) \cdot Y}{(1 - \hat{p}(A = 1 | S = 1, \mathbf{X}_{AY})) \cdot \hat{p}(S = 1)} \right].$$

In the formula above, \mathbb{E}_n is the empirical expectation operator, $\hat{p}(S = 1)$ is an estimate of $p(S = 1)$, and $\hat{p}(A = 1 | \mathbf{X}_{AY}, S = 1)$ is an estimator of $p(A = 1 | \mathbf{X}_{AY}, S = 1)$. Throughout the rest of this section, we assume that the nuisance function estimator and $\hat{p}(S = 1)$ are correctly specified and that the estimation

error of $\hat{\psi}$ for ψ is negligible. Thus, we can treat $\hat{\psi}$ as an unbiased estimator of ψ . Our goal is to study the bias that arises from using $\hat{\psi}$ instead of an unbiased estimator of ϕ_1 . Specifically, we are interested in the following quantity:

$$\text{Bias}_{\phi_1}(\hat{\psi}) = \psi - \phi_1.$$

The differences between the formula of ψ given above and the IPW identification formula of ϕ_1 highlight two sources of bias that result from the practitioner's use of the naive estimator.

First, as described in Section 2, because of the causal path between A and Y through \mathbf{A}_N , one needs to *condition on common causes* of A and A_N (\mathbf{X}_{AG} and \mathbf{U}_{AG}) as well as the common causes of A and S (\mathbf{X}_{AS} and \mathbf{U}_{AS}) for Assumption 5 to hold. However, if the practitioner is unaware of interference, they have no reason to include these common causes inside the propensity score model. The second source of bias arises from the fact that the naive propensity score model used in $\hat{\psi}$ does not account for unmeasured confounders of the relationship between A and Y (\mathbf{U}_{AY}). This becomes particularly problematic in settings with interference, since \mathbf{U}_{AY} may frequently contain both, individual-level covariates (e.g., a given school district's spending per student), as well as neighborhood-level covariates (e.g., the average spending per student of neighboring school districts). Thus, if the practitioner is ignoring interference, they may not only be missing common causes of A and \mathbf{A}_N , common causes of A and S , and individual-level unmeasured confounders, but are likely ignoring most of the neighborhood-level confounders for the relationship between A and Y . Our sensitivity analysis enables practitioners to address the bias resulting from all of these omissions.

3.1 Bias Decomposition and Sensitivity Parameters

To conduct our sensitivity analysis, we begin by comparing ψ with the unbiased IPW form of ϕ_1 (as presented in Proposition 1) through a multiplicative error in weights, given by

$$\varepsilon_a = \frac{p(A = a|S = 1, \tilde{\mathbf{X}}, \tilde{\mathbf{U}})}{p(A = a|S = 1, \mathbf{X}_{AY})}.$$

We refer to $p(A = a|S = 1, \tilde{\mathbf{X}}, \tilde{\mathbf{U}})$ as the *true propensity score* for $A = a$ and to $p(A = a|S = 1, \mathbf{X}_{AY})$ as the *pseudo-propensity score* for $A = a$. We call ε_a the multiplicative error in weights (MEW) score for $A = a$. ε_a forms the basis of our sensitivity analysis, as $\text{Bias}_{\phi_1}(\hat{\psi})$ depends on ε_a and the potential outcomes. We formalize this in Theorem 1.

Theorem 1 (Bias Decomposition for ϕ_1). *Under Assumptions 1 through 5, the bias of the naive estimator, $\hat{\psi}$, for the natural average treatment effect, ϕ_1 , can be written as*

$$\text{Bias}_{\phi_1}(\hat{\psi}) = \sum_{a \in \{0,1\}} (-1)^{1-a} \left(\rho_{Y^{(a,0)}, \varepsilon_a} \cdot \sigma_{Y^{(a,0)}} \cdot \sigma_{\varepsilon_a} \right) \quad (T_1)$$

$$+ \sum_{a \in \{0,1\}} (-1)^{1-a} \left(\rho_{\gamma^{(a)}, \varepsilon_a} \cdot \sigma_{\gamma^{(a)}} \cdot \sigma_{\varepsilon_a} \right), \quad (T_2)$$

where σ_X represents the standard deviation of arbitrary variable X conditional on $S = 1$, and $\rho_{W,Z}$ represents the correlation between arbitrary variables W and Z conditional on $S = 1$.

While the details of the proof are left for Appendix A, we provide some intuition behind Theorem 1 below. T_1 and T_2 each highlight a different aspect of the bias that is introduced by ignored interference and unmeasured confounding. Namely, the bias of the naive estimator for ϕ_1 is governed by a sum of covariances that reflect the relationship between MEW scores and two types of causal variables: the baseline potential

outcomes and the natural spillover effects. We write each of these covariances as a product of a correlation term and two standard deviation terms, which form the basis of our sensitivity analysis framework.

T_1 reflects the joint variability of MEW scores and the baseline potential outcomes, $Y^{(1,0)}$ and $Y^{(0,0)}$. We refer to $Y^{(1,0)}$ and $Y^{(0,0)}$ as the baseline potential outcomes since we view $G = 0$ as the reference neighborhood treatment. Three terms contribute to T_1 : the correlation between the baseline potential outcomes and MEW scores ($\rho_{Y^{(a,0)}, \varepsilon_a}$), the standard deviation of the baseline potential outcomes ($\sigma_{Y^{(a,0)}}$), and the standard deviation of MEW scores (σ_{ε_a}). If the variability of the baseline potential outcomes is high ($\sigma_{Y^{(a,0)}}$ is large), if the variance of the MEW scores is significant (σ_{ε_a} is large), and if MEW scores and potential outcomes are highly correlated ($|\rho_{Y^{(a,0)}, \varepsilon_a}|$ is close to 1), T_1 would contribute significantly to the bias.

T_2 is a measure of the joint variability of individual natural spillover effects ($\gamma(a)$) and MEW scores. T_2 contributes significantly to the bias if the following three conditions are met simultaneously: if there is a large amount of heterogeneity in the natural individual spillover effects ($\sigma_{\gamma(a)}^2$ is large), if there is significant variability in MEW scores ($\sigma_{\varepsilon_a}^2$ is large), and if there is a strong correlation between MEW scores and spillover effects ($|\rho_{\gamma(a), \varepsilon_a}|$ is close to 1).

Remark 5. *The presence of T_2 highlights the fact that under interference, even in the absence of common causes of A and A_N and common causes of A and S , the effects of unmeasured confounding on the bias still depend on the distribution of the individual natural spillover effects ($\gamma(a)$). Intuitively, T_2 adjusts for the fact that $Y^{(a)}$ is no longer equal to $Y^{(a,0)}$, but is now a function of G and several distinct variables, $Y^{(a,g)}$ for $a \in \{0, 1\}$ and $g \in \{0, 1, \dots, g_{max}\}$. Note that in the absence of interference, our bias decomposition reduces to the one proposed by Shen et al. (2011):*

$$\text{Bias}_{\phi_1}(\hat{\psi}) = \sum_{a \in \{0, 1\}} (-1)^{1-a} \left(\rho_{Y^{(a,0)}, \varepsilon'_a} \cdot \sigma_{Y^{(a,0)}} \cdot \sigma_{\varepsilon'_a} \right).$$

where $\varepsilon'_a = \frac{p(A=a | \mathbf{X}_{AY}, \mathbf{U}_{AY}, S=1)}{p(A=a | S=1, \mathbf{X}_{AY})}$.

3.2 Specifying the Sensitivity Parameters

Next, we provide several guidelines for specifying the parameters of our bias decomposition. Note that these guidelines are meant to be widely applicable, and that in many cases, practitioners may have certain domain knowledge that allows them to specify these parameters in a less conservative manner.

Specifying Parameters in T_1 As discussed previously, the first term reflects the relationship between MEW scores and the baseline potential outcomes. Six sensitivity parameters need to be specified in T_1 .

$\sigma_{Y^{(a,0)}}$ acts as a scaling factor in T_1 . If $Y^{(0,0)}$ and $Y^{(1,0)}$ (the baseline potential outcomes) are bounded between y_{\min}^{ref} and y_{\max}^{ref} , then by Popoviciu's inequality on variances, $\sigma_{Y^{(a,0)}}$ is bounded from above by $\frac{1}{2}(y_{\max}^{\text{ref}} - y_{\min}^{\text{ref}})$. We propose specifying an upper bound for both values of $\sigma_{Y^{(a,0)}}$ as a percentage of their maximum possible value through a single sensitivity parameter $\eta_{\text{baseline}} \in [0, 1]$:

$$\max_{a \in \{0, 1\}} \sigma_{Y^{(a,0)}} \leq \frac{\eta_{\text{baseline}}}{2} \cdot (y_{\max}^{\text{ref}} - y_{\min}^{\text{ref}}).$$

Tools like contour plots can be used to assess the sensitivity of $\hat{\psi}$ to a variety of values of η_{baseline} . If the practitioner has additional information about the baseline potential outcomes, they could specify different values of η_{baseline} for $\sigma_{Y^{(1,0)}}$ and $\sigma_{Y^{(0,0)}}$.

σ_{ε_a} is a measure of variability in MEW scores. Note that the mean of ε_a is 1. If ε_a deviates significantly from this value, that would imply that the pseudo-propensity score for $A = a$ significantly overestimates (if

$\varepsilon_a < 1$) or underestimates (if $\varepsilon_a > 1$) the true propensity score (i.e., the effects of ignoring common causes and confounding are fairly strong). Following an argument analogous to the one presented by Shen et al. (2011), it can be proven that for $a \in \{0, 1\}$, $\sigma_{\varepsilon_a} \leq \sqrt{\mathbb{E}\left[\frac{1-p(A=a|S=1, \mathbf{X}_{AY})}{p(A=a|S=1, \mathbf{X}_{AY})} | S=1\right]}$ (this upper bound can be calculated from the observed data). Once again, unless additional information is available to the practitioner, we propose specifying an upper bound for both values of σ_{ε_a} as some proportion of the maximum standard deviation of the MEW scores through a single sensitivity parameter $\eta_\varepsilon \in [0, 1]$:

$$\sigma_{\varepsilon_a} \leq \eta_\varepsilon \cdot \sqrt{\mathbb{E}\left[\frac{1-p(A=a|S=1, \mathbf{X}_{AY})}{p(A=a|S=1, \mathbf{X}_{AY})} | S=1\right]} \text{ for } a \in \{0, 1\}.$$

Remark 6. *The choice of η_ε can be viewed from the perspective of previous sensitivity analyses by Rosenbaum (2002) and Tan (2006). Namely, instead of specifying η_ε , one could assume that ε_a is bounded between α_a^{-1} and α_a , where $\alpha_a \geq 1$ is a sensitivity parameter. In that case, by the Bhatia-Davis inequality, $\sigma_{\varepsilon_a} \leq \sqrt{\alpha_a - 2 + \alpha_a^{-1}}$. A choice of α_a which provides a non-trivial bound on σ_{ε_a} in light of the distribution of the pseudo-propensity scores can be viewed as corresponding to $\eta_\varepsilon < 1$. A value of 1 for α_a and 0 for η_ε corresponds to the setting in which the pseudo-propensity score is identical to the true propensity score (i.e., no confounding or interference).*

The two $\rho_{Y^{(a,0)}, \varepsilon_a}$ values describe the correlation between the baseline potential outcomes and MEW scores, and determine the sign of T_1 . If these correlations are zero, then T_1 does not contribute to the bias, despite the presence of unmeasured confounding and interference. Generally, a positive $\rho_{Y^{(a,0)}, \varepsilon_a}$ implies that units with higher values of $Y^{(a,0)}$ tend to have pseudo-propensity scores that are lower than their true propensity scores; meanwhile, a negative $\rho_{Y^{(a,0)}, \varepsilon_a}$ indicates that the pseudo-propensity scores of units with high values of $Y^{(a,0)}$ tend to overestimate their true propensity scores. If a practitioner does not have domain knowledge which allows them to specify the signs of the correlations, they can first set these signs so as to maximize bias, and then vary their magnitude from 0 to 1 by setting $\rho_{\text{baseline}} = |\rho_{Y^{(0,0)}, \varepsilon_0}| = |\rho_{Y^{(1,0)}, \varepsilon_1}|$.

Specifying Parameters in T_2 T_2 describes the relationship between the individual natural spillover effects and MEW scores. To determine the magnitude and sign of T_2 , four additional terms (beyond those present in T_1) need to be specified.

$\sigma_{\gamma(a)}$ measures the variability of natural spillover effects if G is allowed to vary as it does in the reference population. Note that this term is large if there is a significant number of units which receive neighborhood treatments other than 0 and if there is notable variability in causal contrasts of the form $Y^{(a,g)} - Y^{(a,0)}$. If the natural spillover effects ($\gamma(a)$) are both bounded between x_{\min}^{ref} and x_{\max}^{ref} for $a \in \{0, 1\}$, the worst-case bound for this term is $\frac{1}{2} \cdot (x_{\max}^{\text{ref}} - x_{\min}^{\text{ref}})$. In the absence of additional information, we suggest specifying an upper bound for both values of $\sigma_{\gamma(a)}$ as a single proportion of the maximum standard deviation of $\gamma(a)$ through the sensitivity parameter $\eta_\gamma \in [0, 1]$:

$$\max_{a \in \{0, 1\}} \sigma_{\gamma(a)} \leq \frac{\eta_\gamma}{2} \cdot (x_{\max}^{\text{ref}} - x_{\min}^{\text{ref}}).$$

$\rho_{\gamma(a), \varepsilon_a}$ determines the sign of T_2 . Generally, if the correlation is positive, then the pseudo-propensity scores of units with higher values of $\gamma(a)$ usually underestimate their true propensity scores, and if the correlation is negative, units with higher values of $\gamma(a)$ have pseudo-propensity scores that typically overestimate their true propensity scores. If domain knowledge does not allow a practitioner to specify these correlations, we propose setting the signs of $\rho_{\gamma(a), \varepsilon_a}$ so as to maximize bias, and adjusting their common magnitude $\rho_{\text{spillover}} = |\rho_{\gamma(0), \varepsilon_0}| = |\rho_{\gamma(1), \varepsilon_1}|$ from 0 to 1.

4 Sensitivity Analysis for Transportability

In this section, we consider settings in which the practitioner is interested in estimating the transported natural average treatment effect (ϕ_2) alongside the natural average treatment effect (ϕ_1). Because treatment assignment may differ between the target and reference populations (and our parameter of interest depends on the distribution of G), in general, $\phi_1 \neq \phi_2$. If the practitioner ignores the differing treatment assignments between the two populations, and assumes that their naive estimate is generalizable to the target population, they will incur another source of bias due to lack of transportability, in addition to the bias arising from ignoring interference and unmeasured confounding (Degtiar and Rose, 2023). Proposition 1 highlights that bias due to lack of transportability occurs since the practitioner's naive estimator, $\hat{\psi}$, only incorporates information from the reference population ($S = 1$), rather than the target population ($S = 2$).

To help practitioners interested in estimating the bias they incur due to unmeasured confounding, ignored interference, and lack of transportability, we decompose the bias of $\hat{\psi}$ for ϕ_2 as follows:

$$\text{Bias}_{\phi_2}(\hat{\psi}) = \psi - \phi_2 = (\psi - \phi_1) + (\phi_1 - \phi_2).$$

The second equality forms the basis of our bias decomposition, which is formalized below.

Theorem 2 (Bias Decomposition for ϕ_2). *Under Assumptions 1 through 6, the bias of the naive estimator, $\hat{\psi}$, for the transported natural average treatment effect, ϕ_2 , can be written as*

$$\text{Bias}_{\phi_2}(\hat{\psi}) = \underbrace{\left\{ \text{Bias}_{\phi_1}(\hat{\psi}) \right\}}_{(T_1+T_2)} + \underbrace{\sum_{g=0}^{g_{\max}} \left\{ \tilde{\rho}_{\tau(g),v(g)} \cdot \tilde{\sigma}_{\tau(g)} \cdot \tilde{\sigma}_{v(g)} + \mathbb{E}[\tau(g)] \cdot \left(p(G=g|S=1) - p(G=g|S=2) \right) \right\}}_{(T_3)},$$

where $v(g) = p(G=g|S=1, \mathbf{X}_{GY}, \mathbf{U}_{GY}) - p(G=g|S=2, \mathbf{X}_{GY}, \mathbf{U}_{GY})$, T_1 and T_2 are defined as in Theorem 1, while $\tilde{\sigma}_X$ and $\tilde{\rho}_{W,Z}$ are defined as in Theorem 1, but unconditionally (i.e., over both, the reference and target populations).

Theorem 2 underlines that the difference between ϕ_1 and ϕ_2 , given by T_3 , depends on two main factors for each level of g . The first is the covariance of the controlled individual main effects, $\tau(g)$, and $v(g)$, an important quantity measuring the difference in the nuisance function for G (based on \mathbf{X}_{GY} and \mathbf{U}_{GY}) between the reference and target populations. If $v(g)$ deviates significantly from 0, then the conditional probability of receiving treatment $G=g$ differs greatly between the two populations. At each level of g , the covariance of $\tau(g)$ and $v(g)$ contributes significantly to T_3 if the following three conditions hold simultaneously: if there is a large amount of heterogeneity in the controlled individual main effects ($\tilde{\sigma}_{\tau(g)}$ is large), if there is significant variability in $v(g)$ ($\sigma_{v(g)}$ is large), and if there is a strong correlation between $\tau(g)$ and $v(g)$ ($|\tilde{\rho}_{\tau(g),v(g)}|$ is close to 1). As before, we decompose this covariance term into a product of two standard deviations and a correlation.

The second factor affecting T_3 is the product of $\mathbb{E}[\tau(g)]$ and $p(G=g|S=1) - p(G=g|S=2)$. At each level of g , this term contributes significantly to T_3 if both, the average value of the controlled individual main effects is large, and if there is a strong imbalance in the marginal distribution of G across the reference and target population. This term reflects that our parameter of interest is a function of $\tau(g)$ and G .

In order to specify T_3 , we propose extending our previous methodology. If the controlled individual main effects ($\tau(g)$), like the natural spillover effects ($\gamma(a)$), are bounded between x_{\min}^{ref} and x_{\max}^{ref} for each value of g , then we can write

$$\mathbb{E}[\tau(g)] \in [x_{\min}^{\text{ref}}, x_{\max}^{\text{ref}}] \text{ for } g \in \{0, 1, \dots, g_{\max}\} \text{ and}$$

$$\max_{g \in \{0, 1, \dots, g_{\max}\}} \tilde{\sigma}_{\tau(g)} \leq \frac{\eta_{\tau}}{2} \cdot (x_{\max}^{\text{ref}} - x_{\min}^{\text{ref}}),$$

where $\eta_\tau \in [0, 1]$ represents a proportion of the maximum standard deviation of $\tau(g)$. The practitioner can specify a different parameter for each level of G if they have additional information allowing them to do so (as opposed to specifying a single value η_τ).

Next, we propose setting the sensitivity parameter $\beta \in [0, 1]$ such that for each value of g , $p(G = g|S = 1) - p(G = g|S = 2) \in [-\beta, \beta]$. β reflects a practitioner's beliefs about the imbalance in the neighborhood treatments between the two populations. This parameter not only bounds $p(G = g|S = 1) - p(G = g|S = 2)$, but can also be used to create an upper bound for the standard deviation of $v(g)$ through Popoviciu's inequality:

$$\tilde{\sigma}_{v(g)} \leq \beta.$$

Remark 7. If the practitioner has domain knowledge which allows them to determine all values of $p(G = g|S = 1) - p(G = g|S = 2) = \zeta(g)$, they can use the Bhatia-Davis inequality to create an upper bound for the standard deviation of $v(g)$ as $\tilde{\sigma}_{v(g)} \leq \eta_v \cdot \sqrt{1 - \zeta(g)^2}$ where $\eta_v \in [0, 1]$ is the proportion parameter reflecting the practitioner's belief about the variability of $v(g)$.

As before, in the absence of additional information, the practitioner can set the signs of the correlations $(\tilde{\rho}_{\tau(g), v(g)})$ so as to maximize bias, and assume they have a common magnitude, which we denote by $\rho_{\text{transport}}$.

The sensitivity analysis framework outlined above is broadly generalizable. Nevertheless, Theorem 2 can also be significantly simplified in many settings. We outline several such cases below.

Corollary 1. If the marginal distribution of G is equal between the reference and target populations, T_3 can be simplified to

$$T_3 = \sum_{g=0}^{g_{\max}} \left\{ \tilde{\rho}_{\tau(g), v(g)} \cdot \tilde{\sigma}_{\tau(g)} \cdot \tilde{\sigma}_{v(g)} \right\}.$$

If the neighborhood treatment is randomized in both, the reference population and the target population, T_3 can be simplified to

$$T_3 = \sum_{g=0}^{g_{\max}} \left\{ \mathbb{E}[\tau(g)] \cdot \left(p(G = g|S = 1) - p(G = g|S = 2) \right) \right\}.$$

Proposition 2. Under Assumptions 1 through 6, the bias of $\hat{\psi}$ for ϕ_2 reduces to the bias of $\hat{\psi}$ for ϕ_1 provided that Y is generated according to the model

$$Y = f_0(\mathbf{X}, \mathbf{U}) + f_1(A, \mathbf{X}, \mathbf{U}) + f_2(G, \mathbf{X}, \mathbf{U}) + \epsilon(A, G),$$

where $\epsilon(A, G)$ is a zero-mean noise term uncorrelated with any of the confounders.

Remark 8. As discussed previously, θ is a special case of ϕ_2 , and our bias decomposition can be extended to account for θ . Namely, if the parameter of interest is θ , T_3 becomes

$$T_3 = \sum_{g=0}^{g_{\max}} \left\{ \tilde{\rho}_{\tau(g), \pi(g)} \cdot \tilde{\sigma}_{\tau(g)} \cdot \tilde{\sigma}_{\pi(g)} \right\}.$$

where $\pi(g) = p(G = g|\mathbf{X}_{GY}, \mathbf{U}_{GY}, S = 1)$. Analysis similar to the one done in Proposition 2 shows that if the parameter of interest is θ , then $T_3 = 0$ if Y is generated by

$$Y = f_0(\mathbf{X}, \mathbf{U}) + f_1(A, \mathbf{X}, \mathbf{U}) + f_2(G, \mathbf{X}, \mathbf{U}) + f_3(A, G) + \epsilon(A, G).$$

5 Bias Decomposition in the Presence of Undefined Potential Outcomes

In earlier sections of this paper, we considered settings in which all potential outcomes are well-defined for all units. While Assumption 3, which encodes this condition, is widely applicable, there may be certain settings in which it is unreasonable. Concerns about violations of Assumption 3 are most relevant for applications where the number of neighbors (which determines the domain of the exposure mapping function) is viewed as fixed. For instance, in a geographic context, if G represents the number of a neighboring countries adopting a certain policy, defining the potential outcome $Y^{(1,2)}$ for a country that has only one neighbor could be seen as conceptually incoherent. Moreover, in general, potential outcomes of the form $Y^{(a,g)}$ are undefined for units without neighbors (Forastiere et al., 2021; Kim, 2025). Due to these issues, practitioners may wish to define a different number of potential outcomes for each unit to better match their philosophical intuitions about causality.

In this section, we extend our earlier results to account for this complexity. Without loss of generality, we will assume that $G \in \{0, 1, 2, \dots, g_{max}\}$ represents the number of treated neighbors a unit has. For the sake of notational convenience, we also make the following assumptions:

- $Y^{(a,g)}$ for $a \in \{0, 1\}$ is well-defined for any unit with g or more neighbors, and is undefined for any unit with fewer than g neighbors.
- N records the number of neighbors a unit has, and the maximum value of N is g_{max} (i.e., $N \in \{0, 1, 2, \dots, g_{max}\}$).
- The potential outcomes $Y^{(1,0)}$ and $Y^{(0,0)}$ are well-defined for units with $N = 0$ (i.e., they behave like units that have neighbors none of whom are treated).
- $V_g = \mathbb{I}(N \geq g)$ tracks units for which $Y^{(a,g)}$ is well-defined. Note that if $V_{g'} = 1$, then $V_{\tilde{g}} = 1$ for any $\tilde{g} < g'$.

In order to extend our bias decomposition to this setting, the assumptions and causal estimands outlined in Section 2 need to be modified. First, we introduce a new version of consistency which relates the observed outcome to the potential outcomes in a manner which depends on a unit's number of neighbors.

Assumption 7 (Consistency with Undefined Potential Outcomes). *Given that $N = n$, the potential outcome $Y^{(a)}$ satisfies $Y^{(a)} = \sum_{g=0}^n \mathbb{I}(G = g) \cdot Y^{(a,g)}$ and the observed outcome Y satisfies $Y = \sum_{a=0}^1 \mathbb{I}(A = a) \cdot Y^{(a)}$.*

Since potential outcomes are only well-defined conditionally under this framework, the causal estimands need to be re-defined. We begin by considering a measure of the impact of changing A , the *controlled individual main effect*,

$$\tau(g) := \begin{cases} \text{N.A.} & \text{if } V_g = 0 \\ Y^{(1,g)} - Y^{(0,g)} & \text{if } V_g = 1. \end{cases}$$

Here, $\tau(g)$ still represents the causal impact of changing A from 0 to 1 while g is fixed. However, this quantity is now only well-defined for units who have enough neighbors to receive neighborhood treatment g . The other causal estimand measuring the impact of changing A , the *natural individual main effect*, κ , now also depends on the number of neighbors a unit has. Namely, if $N = n$, then $\kappa = Y^{(1)} - Y^{(0)} = \sum_{g=0}^n \mathbb{I}(G = g) \cdot Y^{(a,g)}$. Our previous parameters of interest, the *natural average main effect*, ϕ_1 , and the *transported natural average main effect*, ϕ_2 , can be expressed as conditional expectations of κ across the reference and target populations

respectively, as

$$\begin{aligned}
\phi_s &= \sum_{n=0}^{g_{\max}} \mathbb{E}[\underbrace{Y^{(1)} - Y^{(0)}}_{\kappa} | S = s, N = n] \cdot p(N = n | S = s) \\
&= \sum_{n=0}^{g_{\max}} \sum_{g=0}^n \mathbb{E}[\mathbb{I}(G = g) \cdot \tau(g) | S = s, N = n] \cdot p(N = n | S = s) \\
&= \sum_{g=0}^{g_{\max}} \sum_{n=g}^{g_{\max}} \mathbb{E}[\mathbb{I}(G = g) \cdot \tau(g) | S = s, N = n] \cdot p(N = n | S = s) \\
&= \sum_{g=0}^{g_{\max}} \mathbb{E}[\mathbb{I}(G = g) \cdot \tau(g) | S = s, V_g = 1] \cdot p(V_g = 1 | S = s).
\end{aligned}$$

Next, we modify our positivity and exchangeability assumptions to account for undefined potential outcomes.

Assumption 8 (Positivity with Undefined Potential Outcomes). *For any $a \in \{0, 1\}$, any $s \in \{1, 2\}$, and any $n, g \in \{0, 1, 2, \dots, g_{\max}\}$, we have, almost surely in (\mathbf{X}, \mathbf{U}) ,*

- $p(A = a | S = s, N = n, \mathbf{X}, \mathbf{U}) > 0$,
- $p(S = s, N = n | \mathbf{X}, \mathbf{U}) > 0$, and
- $p(S = s, V_g = 1 | \mathbf{X}, \mathbf{U}) > 0$.

Assumption 9 (Weak Conditional Exchangeability with Undefined Potential Outcomes). *Given that $N = n$ where $n \in \{0, 1, 2, \dots, g_{\max}\}$, for any $a \in \{0, 1\}$,*

$$Y^{(a)} \perp\!\!\!\perp A | \tilde{\mathbf{X}}, \tilde{\mathbf{U}}, S, N = n.$$

Finally, we posit one new condition, which allows us to deal with the difficulties introduced by the fact that the potential outcome $Y^{(a)}$ for $a \in \{0, 1\}$ is defined differently depending on a unit's number of neighbors.

Assumption 10 (Independence Assumption for Undefined Potential Outcomes). *The two nuisance functions for the personal treatment, the true propensity score and the pseudo-propensity score, do not depend on N . In other words, for $a \in \{0, 1\}$, $s \in \{1, 2\}$, and $n \in \{0, 1, 2, \dots, g_{\max}\}$, we have $p(A = a | S = s, N = n, \tilde{\mathbf{X}}, \tilde{\mathbf{U}}) = p(A = a | S = s, \tilde{\mathbf{X}}, \tilde{\mathbf{U}})$ and $p(A = a | S = s, N = n, \mathbf{X}_{AY}) = p(A = a | S = s, \mathbf{X}_{AY})$.*

Using the assumptions outlined above, we can extend our previous formula providing identification (if the unmeasured confounders were observed) to our new causal estimands.

Proposition 3 (IPW Identification Formula with Undefined Potential Outcomes). *Under Assumptions 1, 7, 8, 9, and 10, ϕ_s can be identified by the formula given in Proposition 1.*

Note that while our identification formula does not change, we can no longer say that $\mathbb{I}(A = a) \cdot Y$ is equivalent to $\mathbb{I}(A = a) \cdot Y^{(a)}$ unless we condition on $N = n$ (since $Y^{(a)}$ is no longer a well-defined quantity over the whole population). While our existing assumptions suffice to develop a bias decomposition for ϕ_1 , we once again need an additional set of conditions to develop a bias decomposition for ϕ_2 .

Assumption 11 (Transportability Assumptions with Undefined Potential Outcomes). *For any $a \in \{0, 1\}$ and any $g \in \{0, 1, 2, \dots, g_{\max}\}$, we have*

1. $Y^{(a,g)} \perp\!\!\!\perp \mathbf{A}_{\mathcal{N}} | V_g = 1, \mathbf{X}_{GY}, \mathbf{U}_{GY}$,

2. S only has a direct causal effect on A and G , and
3. Y and S do not share common causes.

Using the conditions outlined above, we can develop a new bias decomposition for settings with undefined potential outcomes.

Theorem 3 (Bias Decomposition in Presence of Undefined Potential Outcomes). *Under Assumptions 1, 7, 8, 9, and 10, the bias of a naive estimator, $\hat{\psi}$ for the natural average treatment effect, ϕ_1 , can be written as*

$$\begin{aligned} \text{Bias}_{\phi_1}(\hat{\psi}) &= \sum_{a \in \{0,1\}} (-1)^{1-a} \cdot \text{Cov}\left(Y^{(a,0)}, \varepsilon_a | S = 1\right) \\ &+ \sum_{a \in \{0,1\}} \sum_{n=0}^{g_{\max}} (-1)^{1-a} \cdot \text{Cov}\left(\gamma(a), \varepsilon_a | S = 1, N = n\right) \cdot p(N = n | S = 1). \end{aligned}$$

Meanwhile, under Assumptions 1, 7, 8, 9, 10, and 11, the bias of $\hat{\psi}$ for the transported natural main effect, ϕ_2 , can be written as

$$\text{Bias}_{\phi_2}(\hat{\psi}) = \left\{ \text{Bias}_{\phi_1}(\hat{\psi}) \right\} + \sum_{g=0}^{g_{\max}} \left\{ \text{Cov}(\tau(g), v(g) | V_g = 1) + \mathbb{E}[\tau(g) | V_g = 1] \cdot \left(p(G = g | S = 1) - p(G = g | S = 2) \right) \right\},$$

where $v(g) = p(G = g | S = 1, V_g = 1, \mathbf{X}_{GY}, \mathbf{U}_{GY}) \cdot p(V_g = 1 | S = 1) - p(G = g | S = 2, V_g = 1, \mathbf{X}_{GY}, \mathbf{U}_{GY}) \cdot p(V_g = 1 | S = 2)$.

Theorem 3 can serve as the basis of a sensitivity analysis framework developed according to the principles discussed in Sections 3 and 4. Note that if the potential outcome $Y^{(1)}$ is well-defined for all units (without conditioning on N) and if $V_g = 1$ for any $g \in \{0, 1, 2, \dots, g_{\max}\}$, Theorem 3 is equivalent to Theorem 1 by the law of total expectation. Thus, many of our results from the previous sections of this paper can be extended to applications with undefined potential outcomes.

6 Conclusion

We introduced a sensitivity analysis framework, which, unlike existing frameworks, can be used to explore the effects of unmeasured confounding, omitted interference, and lack of transportability simultaneously. Our method does not require strict parametric assumptions about the data generating mechanism and provides practitioners with the ability to specify bias through several easily interpretable sensitivity parameters. These parameters are flexible and can integrate a wide array of perspectives informed by domain knowledge. We also investigated several special cases under which our bias decomposition can be significantly simplified. Finally, we extended our methodology to account for the additional data complexity of undefined potential outcomes, which can arise in the presence of interference.

References

- Aronow, P. M. and Samii, C. (2017). Estimating average causal effects under general interference, with application to a social network experiment. *The Annals of Applied Statistics*, 11(4):1912 – 1947.
- Bareinboim, E. and Pearl, J. (2016). Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences*, 113(27):7345–7352.
- Bhadra, S. and Schweinberger, M. (2025). Causal inference under network interference. *arXiv preprint arXiv:2508.06808*.

- Buchanan, A. L., Katenka, N., Lee, Y., Wu, J., Pantavou, K., Friedman, S. R., Halloran, M. E., Marshall, B. D., Forastiere, L., and Nikolopoulos, G. K. (2023). Methods for assessing spillover in network-based studies of hiv/aids prevention among people who use drugs. *Pathogens*, 12(2):326.
- Chen, W., Cai, R., Qiao, J., Yan, Y., and Hernández-Lobato, J. M. (2025). Causal effect estimation under networked interference without networked unconfoundedness assumption. *arXiv preprint arXiv:2502.19741*.
- Colnet, B., Mayer, I., Chen, G., Dieng, A., Li, R., Varoquaux, G., Vert, J.-P., Josse, J., and Yang, S. (2024). Causal inference methods for combining randomized trials and observational studies: a review. *Statistical Science*, 39(1):165–191.
- Degtiar, I. and Rose, S. (2023). A review of generalizability and transportability. *Annual Review of Statistics and Its Application*, 10(1):501–524.
- Forastiere, L., Airoidi, E. M., and Mealli, F. (2021). Identification and estimation of treatment and interference effects in observational studies on networks. *Journal of the American Statistical Association*, 116(534):901–918.
- Forastiere, L., Del Prete, D., and Sciabolazza, V. L. (2024). Causal inference on networks under continuous treatment interference. *Social Networks*, 76:88–111.
- Hernán, M. A. and VanderWeele, T. J. (2011). Compound treatments and transportability of causal inference. *Epidemiology*, 22(3):368–377.
- Hirano, K. and Imbens, G. W. (2001). Estimation of causal effects using propensity score weighting: An application to data on right heart catheterization. *Health Services and Outcomes Research Methodology*, 2(3):259–278.
- Hong, G., Yang, F., and Qin, X. (2021). Did you conduct a sensitivity analysis? a new weighting-based approach for evaluations of the average treatment effect for the treated. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 184(1):227–254.
- Huang, M. Y. (2024). Sensitivity analysis for the generalization of experimental results. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 187(4):900–918.
- Khot, A., Oprescu, M., Schröder, M., Kagawa, A., and Luo, X. (2025). Spatial deconfounder: Interference-aware deconfounding for spatial causal inference. *arXiv preprint arXiv:2510.08762*.
- Kim, B. (2025). Estimating spillover effects in the presence of isolated nodes. *Spatial Economic Analysis*, pages 1–15.
- Lee, T., Buchanan, A. L., Katenka, N. V., Forastiere, L., Halloran, M. E., Friedman, S. R., and Nikolopoulos, G. (2023). Estimating causal effects of hiv prevention interventions with interference in network-based studies among people who inject drugs. *The Annals of Applied Statistics*, 17(3):2165.
- Liu, J., Ye, F., and Yang, Y. (2023). Nonparametric doubly robust estimation of causal effect on networks in observational studies. *Stat*, 12(1):e549.
- Liu, L., Hudgens, M. G., and Becker-Dreps, S. (2016). On inverse probability-weighted estimators in the presence of interference. *Biometrika*, 103(4):829–842.
- McNealis, V., Moodie, E. E. M., and Dean, N. (2024). Revisiting the effects of maternal education on adolescents’ academic performance: Doubly robust estimation in a network-based observational study. *Journal of the Royal Statistical Society Series C: Applied Statistics*, 73(3):715–734.

- Mitra, N., Roy, J., and Small, D. (2022). The future of causal inference. *American Journal of Epidemiology*, 191(10):1671–1676.
- Papadogeorgou, G. and Samanta, S. (2023). Spatial causal inference in the presence of unmeasured confounding and interference. *arXiv preprint arXiv:2303.08218*.
- Robins, J. M. and Greenland, S. (1992). Identifiability and exchangeability for direct and indirect effects. *Epidemiology*, 3(2):143–155.
- Rosenbaum, P. R. (2002). *Observational studies*. Springer series in statistics. Springer, New York, 2nd ed. edition.
- Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66(5):688.
- Rubin, D. B. (1980). Randomization analysis of experimental data: The fisher randomization test comment. *Journal of the American Statistical Association*, 75(371):591–593.
- Rubin, D. B. (1990). [On the application of probability theory to agricultural experiments. Essay on principles. Section 9.] Comment: Neyman (1923) and causal inference in experiments and observational studies. *Statistical Science*, 5(4):472–480.
- Shen, C., Li, X., Li, L., and Were, M. C. (2011). Sensitivity analysis for causal inference using inverse probability weighting. *Biometrical Journal*, 53(5):822–837.
- Sobel, M. E. (2006). What do randomized studies of housing mobility demonstrate? causal inference in the face of interference. *Journal of the American Statistical Association*, 101(476):1398–1407.
- Tan, Z. (2006). A distributional approach for causal inference using propensity scores. *Journal of the American Statistical Association*, 101(476):1619–1637.
- Tchetgen Tchetgen, E. J. and VanderWeele, T. J. (2012). On causal inference in the presence of interference. *Statistical Methods in Medical Research*, 21(1):55–75.
- Tortú, C., Crimaldi, I., Mealli, F., and Forastiere, L. (2024). Estimating causal effects of multi-valued treatments accounting for network interference: Immigration policies and crime rates. *Sociological Methods & Research*, 53(4):1794–1828.
- VanderWeele, T. J. (2011). Controlled direct and mediated effects: definition, identification and bounds. *Scandinavian Journal of Statistics*, 38(3):551–563.
- VanderWeele, T. J., Tchetgen Tchetgen, E. J., and Halloran, M. E. (2015). Interference and sensitivity analysis. *Statistical science: a review journal of the Institute of Mathematical Statistics*, 29(4):687.
- Vuong, Q., Metcalfe, R. K., Ling, A., Ackerman, B., Inoue, K., and Park, J. J. (2025). Systematic review of applied transportability and generalizability analyses: A landscape analysis. *Annals of Epidemiology*.
- Wu, J. and Franks, A. (2025). A latent factor panel approach to spatiotemporal causal inference. *arXiv preprint arXiv:2509.10974*.
- Zigler, C., Liu, V., Mealli, F., and Forastiere, L. (2025). Bipartite interference and air pollution transport: estimating health effects of power plant interventions. *Biostatistics*, 26(1):kxae051.

Appendices

A Proofs

Proof of Proposition 1. We will begin by considering $\mathbb{E}[Y^{(1)}|S = s]$. Note that:

$$\begin{aligned}
\mathbb{E}[Y^{(1)}|S = s] &= \mathbb{E}\left[\mathbb{E}[Y^{(1)}|S = s, \tilde{\mathbf{X}}, \tilde{\mathbf{U}}]|S = s\right] \\
&= \mathbb{E}\left[\mathbb{E}[Y^{(1)}|S = s, A = 1, \tilde{\mathbf{X}}, \tilde{\mathbf{U}}]|S = s\right] \\
&= \mathbb{E}\left[\mathbb{E}[Y|S = s, A = 1, \tilde{\mathbf{X}}, \tilde{\mathbf{U}}]|S = s\right] \\
&= \int_{\tilde{\mathbf{x}}, \tilde{\mathbf{u}}, y} y \cdot p(y|S = s, A = 1, \tilde{\mathbf{x}}, \tilde{\mathbf{u}}) \cdot p(\tilde{\mathbf{x}}, \tilde{\mathbf{u}}|S = s) \\
&= \int_{\tilde{\mathbf{x}}, \tilde{\mathbf{u}}, y} y \cdot \frac{p(y, A = 1, \tilde{\mathbf{x}}, \tilde{\mathbf{u}}|S = s)}{p(A = 1|S = s, \tilde{\mathbf{x}}, \tilde{\mathbf{u}})} \\
&= \int_{\tilde{\mathbf{x}}, \tilde{\mathbf{u}}, y, a} \mathbb{I}(A = 1) \cdot y \cdot \frac{p(y, a, \tilde{\mathbf{x}}, \tilde{\mathbf{u}}|S = s)}{p(A = 1|S = s, \tilde{\mathbf{x}}, \tilde{\mathbf{u}})} \\
&= \mathbb{E}\left[\frac{\mathbb{I}(A = 1) \cdot Y}{p(A = 1|S = s, \tilde{\mathbf{x}}, \tilde{\mathbf{u}})}|S = s\right].
\end{aligned}$$

An analogous argument for $\mathbb{E}[Y^{(0)}|S = s]$ proves Proposition 1. \square

Proof of Theorem 1. We will begin by rewriting the first component of ψ , $\mathbb{E}\left[\frac{\mathbb{I}(A=1) \cdot Y}{p(A=1|S=1, \mathbf{X}_{AY})}|S = 1\right]$. Note that:

$$\begin{aligned}
\mathbb{E}\left[\frac{\mathbb{I}(A = 1) \cdot Y}{p(A = 1|S = 1, \mathbf{X}_{AY})}|S = 1\right] &= \mathbb{E}\left[\frac{\mathbb{I}(A = 1) \cdot Y^{(1)}}{p(A = 1|S = 1, \mathbf{X}_{AY})}|S = 1\right] \\
&= \mathbb{E}\left[\mathbb{E}\left[\frac{\mathbb{I}(A = 1) \cdot Y^{(1)}}{p(A = 1|S = 1, \mathbf{X}_{AY})}|S = 1, \tilde{\mathbf{X}}, \tilde{\mathbf{U}}, Y^{(1)}\right]|S = 1\right] \\
&= \mathbb{E}\left[\frac{Y^{(1)}}{p(A = 1|S = 1, \mathbf{X}_{AY})} \cdot \mathbb{E}[\mathbb{I}(A = 1)|S = 1, \tilde{\mathbf{X}}, \tilde{\mathbf{U}}, Y^{(1)}]|S = 1\right] \\
&= \mathbb{E}\left[\frac{Y^{(1)}}{p(A = 1|S = 1, \mathbf{X}_{AY})} \cdot \mathbb{E}[\mathbb{I}(A = 1)|S = 1, \tilde{\mathbf{X}}, \tilde{\mathbf{U}}]|S = 1\right] \\
&= \mathbb{E}\left[Y^{(1)} \cdot \frac{p(A = 1|S = 1, \tilde{\mathbf{X}}, \tilde{\mathbf{U}})}{p(A = 1|S = 1, \mathbf{X}_{AY})}|S = 1\right] = \mathbb{E}[Y^{(1)} \cdot \varepsilon_1|S = 1].
\end{aligned}$$

Next, note that since $\mathbb{E}[\varepsilon_1|S = 1] = 1$, we can write:

$$\begin{aligned}
\mathbb{E}\left[\frac{\mathbb{I}(A = 1) \cdot Y}{p(A = 1|S = 1, \mathbf{X}_{AY})}|S = 1\right] - \mathbb{E}[Y^{(1)}|S = 1] &= \mathbb{E}[Y^{(1)} \cdot \varepsilon_1|S = 1] - \mathbb{E}[Y^{(1)}|S = 1] \\
&= \mathbb{E}[Y^{(1)} \cdot (\varepsilon_1 - 1)|S = 1] \\
&= \text{Cov}(Y^{(1)}, \varepsilon_1|S = 1) \\
&= \text{Cov}(Y^{(1,0)}, \varepsilon_1|S = 1) + \text{Cov}(\gamma(1), \varepsilon_1|S = 1)
\end{aligned}$$

We can go through a similar argument to show the following:

$$\mathbb{E}[Y^{(0)}|S=1] - \mathbb{E}\left[\frac{\mathbb{I}(A=0) \cdot Y}{p(A=0|S=1, \mathbf{X}_{AY})} | S=1\right] = -\text{Cov}(Y^{(0,0)}, \varepsilon_0 | S=1) - \text{Cov}(\gamma(0), \varepsilon_0 | S=1).$$

Combining both of the above results proves Theorem 1. \square

Proof of Theorem 2. We have already shown that $\psi - \phi_1 = T_1 + T_2$. Thus, we only have to show that $\phi_1 - \phi_2 = T_3$. Note that we can write $\mathbb{E}[Y^{(1)} - Y^{(0)}|S=s]$ as follows:

$$\begin{aligned} \mathbb{E}[Y^{(1)} - Y^{(0)}|S=s] &= \mathbb{E}\left[\frac{\mathbb{I}(S=s) \cdot (Y^{(1)} - Y^{(0)})}{p(S=s)}\right] \\ &= \sum_{g=0}^{g_{\max}} \mathbb{E}\left[\frac{\mathbb{I}(S=s, G=g) \cdot \tau(g)}{p(S=s)}\right]. \end{aligned}$$

Next, by Assumption 6, we have the following two relationships: (i) $S, G \perp\!\!\!\perp \tau(g) | \mathbf{X}_{GY}, \mathbf{U}_{GY}$, and (ii) $S \perp\!\!\!\perp \mathbf{X}_{GY}, \mathbf{U}_{GY}$. Thus, we can write:

$$\begin{aligned} \mathbb{E}[Y^{(1)} - Y^{(0)}|S=s] &= \sum_{g=0}^{g_{\max}} \mathbb{E}\left[\mathbb{E}\left[\frac{\mathbb{I}(S=s, G=g) \cdot \tau(g)}{p(S=s)} | \tau(g), \mathbf{X}_{GY}, \mathbf{U}_{GY}\right]\right] \\ &= \sum_{g=0}^{g_{\max}} \mathbb{E}\left[\frac{\tau(g)}{p(S=s)} \cdot \mathbb{E}\left[\mathbb{I}(S=s, G=g) | \tau(g), \mathbf{X}_{GY}, \mathbf{U}_{GY}\right]\right] \\ &= \sum_{g=0}^{g_{\max}} \mathbb{E}\left[\frac{\tau(g)}{p(S=s)} \cdot \mathbb{E}\left[\mathbb{I}(S=s, G=g) | \mathbf{X}_{GY}, \mathbf{U}_{GY}\right]\right] \\ &= \sum_{g=0}^{g_{\max}} \mathbb{E}\left[\frac{\tau(g)}{p(S=s)} \cdot p(S=s, G=g | \mathbf{X}_{GY}, \mathbf{U}_{GY})\right] \\ &= \sum_{g=0}^{g_{\max}} \mathbb{E}\left[\frac{\tau(g)}{p(S=s)} \cdot p(G=g | S=s, \mathbf{X}_{GY}, \mathbf{U}_{GY}) \cdot p(S=s | \mathbf{X}_{GY}, \mathbf{U}_{GY})\right] \\ &= \sum_{g=0}^{g_{\max}} \mathbb{E}\left[\tau(g) \cdot p(G=g | S=s, \mathbf{X}_{GY}, \mathbf{U}_{GY})\right] \\ &= \sum_{g=0}^{g_{\max}} \text{Cov}(\tau(g) \cdot p(G=g | S=s, \mathbf{X}_{GY}, \mathbf{U}_{GY})) \\ &\quad + \sum_{g=0}^{g_{\max}} \mathbb{E}[\tau(g)] \cdot \mathbb{E}[p(G=g | S=s, \mathbf{X}_{GY}, \mathbf{U}_{GY})] \\ &= \sum_{g=0}^{g_{\max}} \text{Cov}(\tau(g) \cdot p(G=g | S=s, \mathbf{X}_{GY}, \mathbf{U}_{GY})) \\ &\quad + \sum_{g=0}^{g_{\max}} \mathbb{E}[\tau(g)] \cdot p(G=g | S=s) \end{aligned}$$

Utilizing the above argument for $S=1$ and $S=2$ completes the proof. \square

Proof of Proposition 2. First, notice that since we assumed that G is discrete, for $g \in \{1, 2, \dots, g_{\max}\}$ and $s \in \{1, 2\}$, we have $p(G=g|S=s) = p(G \leq g|S=s) - p(G \leq g-1|S=s)$ and $p(G=g|S=s, \mathbf{X}_{GY}, \mathbf{U}_{GY}) =$

$p(G \leq g|S = s, \mathbf{X}_{GY}, \mathbf{U}_{GY}) - p(G \leq g-1|S = s, \mathbf{X}_{GY}, \mathbf{U}_{GY})$. Thus, for $s \in \{1, 2\}$,

$$\sum_{g=0}^{g_{max}} \text{Cov}\left(\tau(g), p(G = g|S = s, \mathbf{X}_{GY}, \mathbf{U}_{GY})\right) = \sum_{g=1}^{g_{max}} \text{Cov}\left(\tau(g-1) - \tau(g), p(G \leq g-1|S = s, \mathbf{X}_{GY}, \mathbf{U}_{GY})\right).$$

Similarly, for $s \in \{1, 2\}$, we have

$$\sum_{g=0}^{g_{max}} \mathbb{E}[\tau(g)] \cdot p(G = g|S = s) = \sum_{g=1}^{g_{max}} \left\{ \mathbb{E}[\tau(g-1) - \tau(g)] \cdot p(G \leq g-1|S = s) \right\} + \mathbb{E}[\tau(g_{max})].$$

Next, note that if $Y = f_0(\mathbf{X}, \mathbf{U}) + f_1(A, \mathbf{X}, \mathbf{U}) + f_2(G, \mathbf{X}, \mathbf{U}) + \epsilon(A, G)$, for any $g \in \{1, 2, \dots, g_{max}\}$, $\text{Cov}(\tau(g-1) - \tau(g), p(G \leq g-1|S = s, \mathbf{X}_{GY}, \mathbf{U}_{GY})) = 0$ and $\mathbb{E}[\tau(g-1) - \tau(g)] = 0$. Thus, $T_3 = \mathbb{E}[\tau(g_{max})] - \mathbb{E}[\tau(g_{max})] = 0$, completing the proof. \square

Proof of Proposition 3. Similar to our proof of Proposition 1, we will begin by considering $\mathbb{E}[Y^{(1)}|S = s, N = n]$. Note that:

$$\begin{aligned} \mathbb{E}[Y^{(1)}|S = s, N = n] &= \mathbb{E}\left[\mathbb{E}[Y^{(1)}|S = s, N = n, \tilde{\mathbf{X}}, \tilde{\mathbf{U}}]|S = s, N = n\right] \\ &= \mathbb{E}\left[\mathbb{E}[Y^{(1)}|S = s, N = n, A = 1, \tilde{\mathbf{X}}, \tilde{\mathbf{U}}]|S = s, N = n\right] \\ &= \mathbb{E}\left[\mathbb{E}[Y|S = s, N = n, A = 1, \tilde{\mathbf{X}}, \tilde{\mathbf{U}}]|S = s, N = n\right] \\ &= \int_{\tilde{\mathbf{x}}, \tilde{\mathbf{u}}, y} y \cdot p(y|S = s, N = n, A = 1, \tilde{\mathbf{x}}, \tilde{\mathbf{u}}) \cdot p(\tilde{\mathbf{x}}, \tilde{\mathbf{u}}|S = s, N = n) \\ &= \int_{\tilde{\mathbf{x}}, \tilde{\mathbf{u}}, y} y \cdot \frac{p(y, A = 1, \tilde{\mathbf{x}}, \tilde{\mathbf{u}}|S = s, N = n)}{p(A = 1|S = s, N = n, \tilde{\mathbf{x}}, \tilde{\mathbf{u}})} \\ &= \int_{\tilde{\mathbf{x}}, \tilde{\mathbf{u}}, y, a} \mathbb{I}(A = 1) \cdot y \cdot \frac{p(y, a, \tilde{\mathbf{x}}, \tilde{\mathbf{u}}|S = s, N = n)}{p(A = 1|S = s, N = n, \tilde{\mathbf{x}}, \tilde{\mathbf{u}})} \\ &= \mathbb{E}\left[\frac{\mathbb{I}(A = 1) \cdot Y}{p(A = 1|S = s, N = n, \tilde{\mathbf{X}}, \tilde{\mathbf{U}})}|S = s, N = n\right] \\ &= \mathbb{E}\left[\frac{\mathbb{I}(A = 1) \cdot Y}{p(A = 1|S = s, \tilde{\mathbf{X}}, \tilde{\mathbf{U}})}|S = s, N = n\right]. \end{aligned}$$

Next, we can write

$$\begin{aligned} \sum_{n=0}^{g_{max}} \mathbb{E}[Y^{(1)}|S = s, N = n] \cdot p(N = n|S = s) &= \sum_{n=0}^{g_{max}} \mathbb{E}\left[\frac{\mathbb{I}(A = 1) \cdot Y}{p(A = 1|S = s, \tilde{\mathbf{X}}, \tilde{\mathbf{U}})}|S = s, N = n\right] \cdot p(N = n|S = s) \\ &= \mathbb{E}\left[\frac{\mathbb{I}(A = 1) \cdot Y}{p(A = 1|S = s, \tilde{\mathbf{X}}, \tilde{\mathbf{U}})}|S = s\right]. \end{aligned}$$

The last equality uses the law of total expectation. Note that we did not use the law of total expectation while dealing with the potential outcome $Y^{(1)}$ since $Y^{(1)}$ is only well-defined conditionally on $N = n$. An analogous argument for $\mathbb{E}[Y^{(0)}|S = s, N = n]$ proves Proposition 3. \square

Proof of Theorem 3. First, we will prove our result for Bias_{ϕ_1} . We can begin by writing:

$$\mathbb{E}\left[\frac{\mathbb{I}(A = 1) \cdot Y}{p(A = 1|S = 1, \mathbf{X}_{AY})}|S = 1\right] = \sum_{n=0}^{g_{max}} \mathbb{E}\left[\frac{\mathbb{I}(A = 1) \cdot Y}{p(A = 1|S = 1, \mathbf{X}_{AY})}|S = 1, N = n\right] \cdot p(N = n|S = 1).$$

We can re-write this quantity as:

$$\begin{aligned}
&= \sum_{n=0}^{g_{\max}} \mathbb{E} \left[\frac{\mathbb{I}(A=1) \cdot Y^{(1)}}{p(A=1|S=1, \mathbf{X}_{AY})} | S=1, N=n \right] \cdot p(N=n|S=1) \\
&= \sum_{n=0}^{g_{\max}} \mathbb{E} \left[\mathbb{E} \left[\frac{\mathbb{I}(A=1) \cdot Y^{(1)}}{p(A=1|S=1, \mathbf{X}_{AY})} | S=1, N=n, \tilde{\mathbf{X}}, \tilde{\mathbf{U}}, Y^{(1)} \right] | S=1, N=n \right] \cdot p(N=n|S=1) \\
&= \sum_{n=0}^{g_{\max}} \mathbb{E} \left[\frac{Y^{(1)}}{p(A=1|S=1, \mathbf{X}_{AY})} \cdot \mathbb{E} \left[\mathbb{I}(A=1) | S=1, N=n, \tilde{\mathbf{X}}, \tilde{\mathbf{U}}, Y^{(1)} \right] | S=1, N=n \right] \cdot p(N=n|S=1) \\
&= \sum_{n=0}^{g_{\max}} \mathbb{E} \left[\frac{Y^{(1)}}{p(A=1|S=1, \mathbf{X}_{AY})} \cdot \mathbb{E} \left[\mathbb{I}(A=1) | S=1, N=n, \tilde{\mathbf{X}}, \tilde{\mathbf{U}} \right] | S=1, N=n \right] \cdot p(N=n|S=1) \\
&= \sum_{n=0}^{g_{\max}} \mathbb{E} \left[\frac{Y^{(1)}}{p(A=1|S=1, \mathbf{X}_{AY})} \cdot p(A=1|S=1, N=n, \tilde{\mathbf{X}}, \tilde{\mathbf{U}}) | S=1, N=n \right] \cdot p(N=n|S=1) \\
&= \sum_{n=0}^{g_{\max}} \mathbb{E} \left[Y^{(1)} \cdot \varepsilon_1 | S=1, N=n \right] \cdot p(N=n|S=1).
\end{aligned}$$

Note that by Assumption 10, $\varepsilon_1 = \frac{p(A=1|S=1, \tilde{\mathbf{X}}, \tilde{\mathbf{U}})}{p(A=1|S=1, \mathbf{X}_{AY})} = \frac{p(A=1|S=1, N=n, \tilde{\mathbf{X}}, \tilde{\mathbf{U}})}{p(A=1|S=1, N=n, \mathbf{X}_{AY})}$, and thus $\mathbb{E}[\varepsilon_1 | S=1, N=n] = \mathbb{E}[\varepsilon_1 | S=1] = 1$. Recall that we are interested in the following quantity:

$$\mathbb{E} \left[\frac{\mathbb{I}(A=1) \cdot Y}{p(A=1|S=1, \mathbf{X}_{AY})} | S=1 \right] - \sum_{n=0}^{g_{\max}} \left(\mathbb{E}[Y^{(1)} | S=1, N=n] \cdot p(N=n|S=1) \right).$$

We can write the above expression as:

$$\begin{aligned}
&= \sum_{n=0}^{g_{\max}} \left(\mathbb{E} \left[Y^{(1)} \cdot (\varepsilon_1 - 1) | S=1, N=n \right] \right) \cdot p(N=n|S=1) \\
&= \sum_{n=0}^{g_{\max}} \left(\mathbb{E} \left[Y^{(1,0)} \cdot (\varepsilon_1 - 1) | S=1, N=n \right] \right) \cdot p(N=n|S=1) \\
&+ \sum_{n=0}^{g_{\max}} \left(\mathbb{E} \left[\gamma(1) \cdot (\varepsilon_1 - 1) | S=1, N=n \right] \right) \cdot p(N=n|S=1) \\
&= \mathbb{E} \left[Y^{(1,0)} \cdot (\varepsilon_1 - 1) | S=1 \right] \\
&+ \sum_{n=0}^{g_{\max}} \left(\mathbb{E} \left[\gamma(1) \cdot (\varepsilon_1 - 1) | S=1, N=n \right] \right) \cdot p(N=n|S=1) \\
&= \text{Cov}(Y^{(1,0)}, \varepsilon_1 | S=1) + \sum_{n=0}^{g_{\max}} \text{Cov}(\gamma(1), \varepsilon_1 | S=1, N=n) \cdot p(N=n|S=1).
\end{aligned}$$

Note that while we can use the law of total expectation for $Y^{(1,0)}$ (since, in our setting, this quantity is well-defined for all units), we cannot do the same for $\gamma(a)$, as it is defined differently depending on a unit's number of neighbors. An analogous argument for $\sum_n \mathbb{E}[Y^{(0)} | S=1, N=n] \cdot p(N=n|S=1)$ proves the first part of Theorem 3.

Next, we will prove our result for Bias_{ϕ_2} . Note that under our assumptions, we can write $\mathbb{E}[\mathbb{I}(G=g) \cdot \tau(g) | S=$

$s, V_g = 1]$ as follows:

$$\begin{aligned}
\mathbb{E}[\mathbb{I}(G = g) \cdot \tau(g) | S = s, V_g = 1] &= \mathbb{E} \left[\frac{\mathbb{I}(G = g, S = s) \cdot \tau(g)}{p(S = s | V_g = 1)} | V_g = 1 \right] \\
&= \mathbb{E} \left[\mathbb{E} \left[\frac{\mathbb{I}(G = g, S = s) \cdot \tau(g)}{p(S = s | V_g = 1)} | V_g = 1, \mathbf{X}_{GY}, \mathbf{U}_{GY}, \tau(g) \right] | V_g = 1 \right] \\
&= \mathbb{E} \left[\frac{\tau(g)}{p(S = s | V_g = 1)} \cdot \mathbb{E} [\mathbb{I}(G = g, S = s) | V_g = 1, \mathbf{X}_{GY}, \mathbf{U}_{GY}, \tau(g)] | V_g = 1 \right] \\
&= \mathbb{E} \left[\frac{\tau(g)}{p(S = s | V_g = 1)} \cdot \mathbb{E} [\mathbb{I}(G = g, S = s) | V_g = 1, \mathbf{X}_{GY}, \mathbf{U}_{GY}] | V_g = 1 \right] \\
&= \mathbb{E} \left[\frac{\tau(g)}{p(S = s | V_g = 1)} \cdot p(G = g, S = s | V_g = 1, \mathbf{X}_{GY}, \mathbf{U}_{GY}) | V_g = 1 \right].
\end{aligned}$$

Using Assumption 11, we also have:

$$\begin{aligned}
\sum_{g=0}^{g_{\max}} \mathbb{E}[\mathbb{I}(G = g) \cdot \tau(g) | S = s, V_g = 1] \cdot p(V_g = 1 | S = s) &= \sum_{g=0}^{g_{\max}} \mathbb{E} \left[\frac{\tau(g)}{p(S = s | V_g = 1)} \cdot p(G = g, S = s | V_g = 1, \mathbf{X}_{GY}, \mathbf{U}_{GY}) | V_g = 1 \right] \cdot p(V_g = 1 | S = s) \\
&= \sum_{g=0}^{g_{\max}} \mathbb{E} \left[\tau(g) \cdot p(G = g | S = s, V_g = 1, \mathbf{X}_{GY}, \mathbf{U}_{GY}) | V_g = 1 \right] \cdot p(V_g = 1 | S = s) \\
&= \sum_{g=0}^{g_{\max}} \text{Cov}(\tau(g), p(G = g | S = s, V_g = 1, \mathbf{X}_{GY}, \mathbf{U}_{GY}) | V_g = 1) \cdot p(V_g = 1 | S = s) \\
&\quad + \mathbb{E}[\tau(g) | V_g = 1] \cdot \mathbb{E}[p(G = g | S = s, V_g = 1, \mathbf{X}_{GY}, \mathbf{U}_{GY}) | V_g = 1] \cdot p(V_g = 1 | S = s) \\
&= \sum_{g=0}^{g_{\max}} \text{Cov}(\tau(g), p(G = g | S = s, V_g = 1, \mathbf{X}_{GY}, \mathbf{U}_{GY}) | V_g = 1) \cdot p(V_g = 1 | S = s) \\
&\quad + \mathbb{E}[\tau(g) | V_g = 1] \cdot p(G = g, V_g = 1 | S = s) \\
&= \sum_{g=0}^{g_{\max}} \text{Cov}(\tau(g), p(G = g | S = s, V_g = 1, \mathbf{X}_{GY}, \mathbf{U}_{GY}) | V_g = 1) \cdot p(V_g = 1 | S = s) \\
&\quad + \mathbb{E}[\tau(g) | V_g = 1] \cdot p(G = g | S = s).
\end{aligned}$$

Utilizing the above argument for $S = 1$ and $S = 2$ completes the proof. \square