

Research proposal

# Enhancing Causal Inference with Network Information

**Polina Revina (8461465)**

**Supervisors: Erik-Jan van Kesteren, Javier Garcia-Bernardo**

*Methodology and Statistics for the Behavioural, Biomedical and Social Sciences  
Utrecht University*

**October 6, 2024**

**Word count: 750**

Candidate journals: Statistical Science; Sociological Methods & Research

# 1 Background

In many experimental and observational studies, interference among units is present, meaning that the treatment of one unit affects its neighbours' outcomes (Cox, 1958). This happens due to exposure transmission through social interaction or physical proximity of units. For instance, students who received tutoring support can interact with other students not assigned to this program and transmit the knowledge obtained during tutoring (Forastiere, Airoldi, & Mealli, 2021). As was demonstrated by Sobel, using methods that don't account for interference or wrongly imposing this assumption can lead to an entirely wrong conclusion about the effect (Sobel, 2006).

Different strategies were proposed in the literature to deal with interference effects. One possibility is to impose an additional assumption on the form that interference is taking. We can impose a *partial interference* assumption, meaning that the population can be divided into clusters and units can interfere within clusters but not across them. This assumption is sometimes coupled with *stratified interference*, proposed by Hundgens and Halloran, meaning that the interference is only affected by the proportion of other units with treatment in the group (Hudgens & Halloran, 2008).

Regarding experiment strategies, one possible design is to use clusters as randomisation units, another is a sequential two-stage randomisation, where in the first stage initial treatment is assigned and based on their response the final treatment is assigned in the second stage. Regarding the inference, there are also some options in the literature. Saeveje et al extended the Horvitz-Thompson estimator to the presence of a partial interference but its clear use in observational network data is lacking (Sävje, Aronow, & Hudgens, 2021). Van der Laan (Van der Laan, 2014) proposed a TMLE estimator and Ogburn et al (Ogburn, Sofrygin, Diaz, & Van der Laan, 2024) extended it to allow for dependence based on contamination. Forastiere et al (Forastiere et al., 2021) developed a new method based on propensity score that balances individual and neighbourhood covariates, which no longer require partial interference assumption.

However, current methods are more focused on experimental design than observational studies. There is still a limited understanding of the degree of bias of causal inference estimators in observational studies under different exposure scenarios and network structures. Existing methods often rely on the partial interference assumption, which may not reflect more complex scenarios where interference can occur across groups.

## 2 Research plan

This project aims to enhance causal inference methods for observational data by incorporating network information. In this research, I aim to identify causal effect estimation bias through a simulation study and to explore methods that can correct this bias. Thus, the research question is "*What is the bias of causal inference estimators that doesn't account for interference in the presence of it?*".

### Data generation process

In the first step, I'll generate synthetic data that simulates how the outcome and the treatment spread through the network. To generate a network I'll use the Barabasi-Albert model which is commonly used for modelling real-world systems, for instance, friendship groups or high school networks. It generates networks with a power-law degree distribution, resulting in some nodes having a much larger number of connections than others.

Data is simulated based on the network obtained in the previous step using spatial autocorrelation models, which model the effect of network dependencies. Under these models, each individual’s outcome can be presented as (i) weighted average responses of neighbouring units, (ii) a set of covariates and (iii) independent noise. The set of covariates includes two individual-level covariates, a measure of an individual characteristic and degree of centrality, and one network-level covariate, network efficiency. This allows the simulation of the relational effect from having treated neighbours, positional from having a certain position in the network and structural effect from the overall network structure.

In modelling the contagion mechanism, I assume a simple contagion, meaning that one connection to the treated is enough for the effect transmission. As the treatment of the unit depends on its covariates, it will be modelled using logistic regression.

## Causal effect estimation

Several causal inference methods will be applied to data from the previous step to estimate effect and bias from true causal effect. I chose matching, which is a statistical technique used in observational studies to estimate causal effects by pairing units with similar characteristics across treatment and control groups.

## New methods

The final step is dedicated to exploring methods to correct this bias by utilizing network information, namely graph neural networks.

## Software

All computations will be done in Python and the following libraries: Networkx for network generation, Pymatch for matching implementation, PyG and DGL for graph machine learning.

## References

- Cox, D. R. (1958). Planning of experiments.
- Forastiere, L., Airoidi, E. M., & 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.
- Hudgens, M. G., & Halloran, M. E. (2008). Toward causal inference with interference. *Journal of the American Statistical Association*, 103(482), 832–842.
- Ogburn, E. L., Sofrygin, O., Diaz, I., & Van der Laan, M. J. (2024). Causal inference for social network data. *Journal of the American Statistical Association*, 119(545), 597–611.
- Sävje, F., Aronow, P., & Hudgens, M. (2021). Average treatment effects in the presence of unknown interference. *Annals of statistics*, 49(2), 673.
- 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.
- Van der Laan, M. J. (2014). Causal inference for a population of causally connected units. *Journal of Causal Inference*, 2(1), 13–74.