

# Estimating network-mediated causal effects via spectral embeddings

Alex Hayes<sup>1</sup>   Mark M. Fredrickson<sup>2</sup>   Keith Levin<sup>1</sup>

<sup>1</sup>Department of Statistics, University of Wisconsin-Madison

<sup>2</sup>Department of Statistics, University of Michigan

## Abstract

We consider the task of mediation analysis for network data, and present a model in which mediation occurs in a latent embedding space. Under this model, node-level interventions have causal effects on nodal outcomes, and these effects can be partitioned into a direct effect independent of the network, and an indirect effect induced by homophily.

## Motivating example: smoking in adolescent social networks

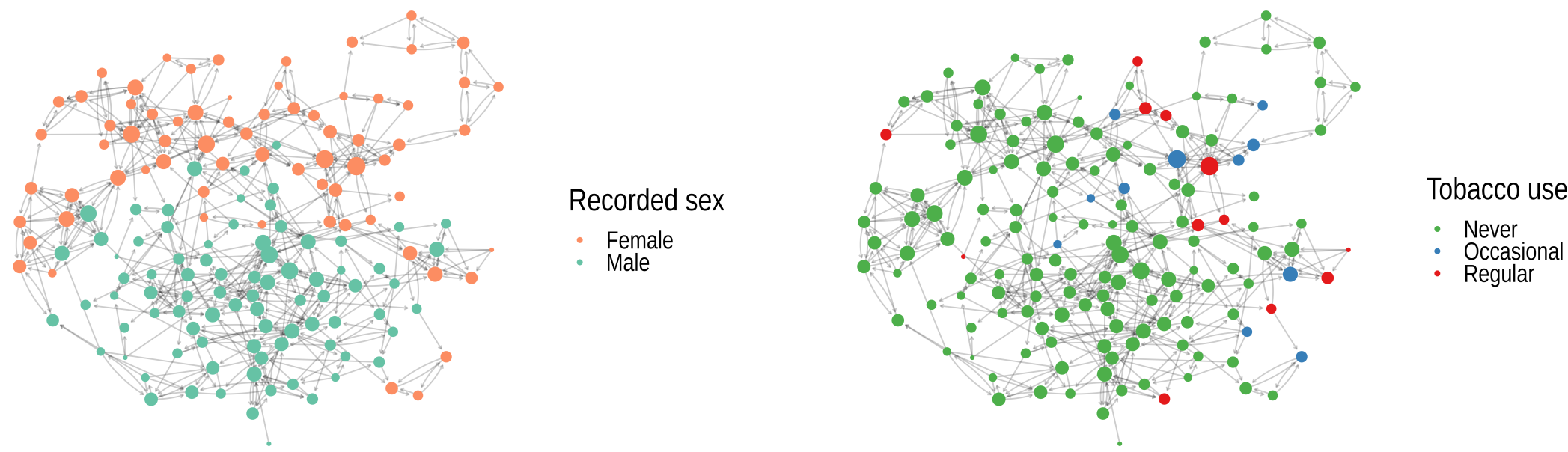


Figure 1. Directed friendships in a secondary school in Glasgow, reported in the Teenage Friends and Lifestyle Study (wave 1). Each node represents one student.

## Notation & inferential targets

We assume we have a (symmetric) network with nodes  $1, \dots, n$ .

Network adjacency matrix	$A \in \mathbb{R}^{n \times n}$
Edge $i \sim j$	$A_{ij} \in \mathbb{R}$
Treatment	$T_i \in \{0, 1\}$
Outcome	$Y_i \in \mathbb{R}$
Confounders	$C_i \in \mathbb{R}^p$
Friend group (latent)	$X_i \in \mathbb{R}^d$

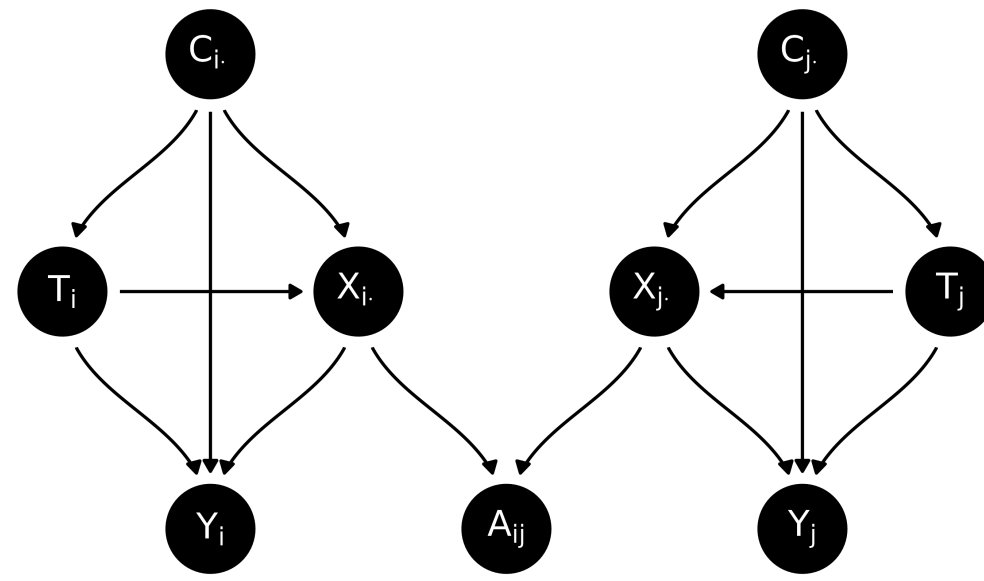


Figure 2. A directed acyclic graph (DAG) representing the causal pathways in a network with homophilous mediation, for node  $a$  in a network with two nodes called  $i$  and  $j$ .

We are interested in the causal effect of  $T_i$  on  $Y_i$  as mediated by the latent position  $X_i$ . More precisely, we want to estimate the *natural direct effect* and the *natural indirect effect*

$$\begin{aligned}\Psi_{\text{nde}}(t, t^*) &= \mathbb{E}[Y_i(t, X_i(t^*)) - Y_i(t^*, X_i(t^*))] \\ \Psi_{\text{nie}}(t, t^*) &= \mathbb{E}[Y_i(t, X_i(t)) - Y_i(t, X_i(t^*))]\end{aligned}$$

## Semi-parametric network model

Let  $A \in \mathbb{R}^{n \times n}$  be a random symmetric matrix, such as the adjacency matrix of an undirected graph. Let  $P = \mathbb{E}[A | X] = XX^T$  be the expectation of  $A$  conditional on  $X \in \mathbb{R}^{n \times d}$ , which has independent and identically distributed rows  $X_1, \dots, X_n$ . That is,  $P$  has  $\text{rank}(P) = d$  and is positive semi-definite with eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d > 0 = \lambda_{d+1} = \dots = \lambda_n$ . Conditional on  $X$ , the upper-triangular elements of  $A - P$  are independent  $(\nu_n, b_n)$ -sub-gamma random variables.

Examples: (degree-corrected) stochastic blockmodels, mixed-membership blockmodels, overlapping blockmodels, (weighted, noisily-observed) random dot product graphs, LDA, factor models, etc

The outcome regression functional is linear in  $T_i, C_i$ , and  $X_i$  and the mediator regression functional is linear in  $T_i, C_i$ , and  $T_i \cdot C_i$ :

$$\begin{aligned}\underbrace{\mathbb{E}[Y_i | T_i, C_i, X_i]}_{\mathbb{R}} &= \underbrace{\beta_0}_{\mathbb{R}} + \underbrace{T_i}_{\{0,1\}} \underbrace{\beta_t}_{\mathbb{R}} + \underbrace{C_i}_{\mathbb{R}^{1 \times p}} \underbrace{\beta_c}_{\mathbb{R}^p} + \underbrace{X_i}_{\mathbb{R}^{1 \times d}} \underbrace{\beta_x}_{\mathbb{R}^d}, & (\text{outcome model}) \\ \underbrace{\mathbb{E}[X_i | T_i, C_i]}_{\mathbb{R}^{1 \times d}} &= \underbrace{\theta_0}_{\mathbb{R}^{1 \times d}} + \underbrace{T_i}_{\{0,1\}} \underbrace{\theta_t}_{\mathbb{R}^{1 \times d}} + \underbrace{C_i}_{\mathbb{R}^{1 \times p}} \underbrace{\Theta_c}_{\mathbb{R}^{p \times d}} + \underbrace{T_i}_{\{0,1\}} \underbrace{C_i}_{\mathbb{R}^{1 \times p}} \underbrace{\Theta_{tc}}_{\mathbb{R}^{p \times d}}. & (\text{mediator model})\end{aligned}$$

Under these moment assumptions, and DAG of Figure 2, letting  $\mu_c$  denote the mean of  $C_i$ , we have the following identification result:

$$\begin{aligned}\Psi_{\text{nde}}(t, t^*) &= (t - t^*) \beta_t, \text{ and} \\ \Psi_{\text{nie}}(t, t^*) &= (t - t^*) \theta_t \beta_x + (t - t^*) \mu_c \Theta_{tc} \beta_x.\end{aligned}$$

## Estimation challenge: friend groups $X$ unknown!

The *adjacency spectral embedding* (ASE) of  $A$  is well-known to be a good estimate of  $X$  under a broad, semi-parametric class of network models. Given a network with adjacency matrix  $A$ , the  $d$ -dimensional ASE is defined as

$$\hat{X} = \hat{U} \hat{S}^{1/2} \in \mathbb{R}^{n \times d},$$

where  $\hat{U} \hat{S} \hat{U}^T$  is the rank- $d$  truncated singular value decomposition of  $A$ . Under a suitably well-behaved network model, if  $d$  is correctly specified or consistently estimated, there is some  $d \times d$  orthogonal matrix  $Q$  such that

$$\max_{i \in [n]} \left\| \hat{X}_i - X_i Q \right\| = o_p(1).$$

## Estimation: plug in $\hat{X}$ for $X$

Let  $\hat{D} = \begin{bmatrix} 1 & T & C & \hat{X} \end{bmatrix} \in \mathbb{R}^{n \times (2+p+d)}$  and  $L = \begin{bmatrix} 1 & T & C & T \cdot C \end{bmatrix} \in \mathbb{R}^{n \times (2p+2)}$ . We estimate  $\beta_w$  and  $\beta_x$  via ordinary least squares as follows

$$\begin{bmatrix} \hat{\beta}_w \\ \hat{\beta}_x \end{bmatrix} = \left( \hat{D}^T \hat{D} \right)^{-1} \hat{D}^T Y.$$

Similarly, we estimate  $\Theta$  via ordinary least squares as

$$\hat{\Theta} = \left( L^T L \right)^{-1} L^T \hat{X}.$$

To estimate  $\Psi_{\text{nde}}$  and  $\Psi_{\text{nie}}$ , we combine regression coefficients from the network regression models

$$\begin{aligned}\hat{\Psi}_{\text{cde}} &= \hat{\Psi}_{\text{nde}} = (t - t^*) \hat{\beta}_t & \text{and} \\ \hat{\Psi}_{\text{nie}} &= (t - t^*) \hat{\theta}_t \hat{\beta}_x + (t - t^*) \cdot \hat{\mu}_c \cdot \hat{\Theta}_{tc} \hat{\beta}_x,\end{aligned}$$

where  $\hat{\mu}_c$  is the sample mean of  $C_i$ .

## Theory

Under a suitable network model and moment bounds on the regression errors, there exists a sequence of orthogonal matrices  $\{Q_n\}_{n=1}^\infty$  such that

$$\begin{aligned}\sqrt{n} \hat{\Sigma}_{\text{vec}(\Theta)}^{-1/2} \left( \text{vec} \left( \hat{\Theta} Q_n^T \right) - \text{vec}(\Theta) \right) &\rightarrow \mathcal{N}(0, I_{pd}), \text{ and} \\ \sqrt{n} \hat{\Sigma}_\beta^{-1/2} \left( \begin{bmatrix} \hat{\beta}_w - \beta_w \\ Q_n \hat{\beta}_x - \beta_x \end{bmatrix} \right) &\rightarrow \mathcal{N}(0, I_d).\end{aligned}$$

Further,

$$\begin{aligned}\sqrt{n} \hat{\sigma}_{\text{nde}}^2 \left( \hat{\Psi}_{\text{nde}} - \Psi_{\text{nde}} \right) &\rightarrow \mathcal{N}(0, 1), \text{ and} \\ \sqrt{n} \hat{\sigma}_{\text{nie}}^2 \left( \hat{\Psi}_{\text{nie}} - \Psi_{\text{nie}} \right) &\rightarrow \mathcal{N}(0, 1),\end{aligned}$$

where  $\hat{\sigma}_{\text{nde}}^2$  and  $\hat{\sigma}_{\text{nie}}^2$  are derived via the Delta method.

## Results applied to Glasgow data

