

Estimating network-mediated causal effects via spectral embeddings

Alex Hayes and Keith Levin

2022-07-07

University of Wisconsin-Madison

Setting



Network $A \in \{0, 1\}^{n \times n}$

Nodal covariates $\mathbf{W}_1, \dots, \mathbf{W}_n \in \mathbb{R}^p$

Nodal outcomes $Y_1, \dots, Y_n \in \mathbb{R}$

Network model: stochastic blockmodel (SBM)

Degree-corrected SBM



k “blocks” or communities

$\mathbf{Z}_i \in \{0, 1\}^k$ one-hot indicator of node i 's block

$\gamma_i \in [0, 1]$ node i 's popularity

$B \in [0, 1]^{k \times k}$ inter-block edge probabilities

$$\mathbb{P}[A_{ij} = 1 \mid \mathbf{Z}, \gamma] = \gamma_i \cdot \mathbf{Z}_i B \mathbf{Z}_j^T \cdot \gamma_j$$

Latent positions in the SBM

$$\mathbb{P}[A \mid \mathbf{Z}, B] = \underbrace{\begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}}_{\mathbf{z}} \underbrace{\begin{bmatrix} 0.6 & 0.02 \\ 0.02 & 0.6 \end{bmatrix}}_B \underbrace{\begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix}}_{\mathbf{z}^T}$$

$$\mathbb{P}[A \mid \mathbf{Z}, B] = \begin{bmatrix} 0.6 & 0.6 & 0.02 & 0.02 \\ 0.6 & 0.6 & 0.02 & 0.02 \\ 0.02 & 0.02 & 0.6 & 0.6 \\ 0.02 & 0.02 & 0.6 & 0.6 \end{bmatrix}$$

Intuition: incorporate latent positions into regression

$$\underbrace{\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}}_{\text{outcome}} = \underbrace{\begin{bmatrix} 2.02 & 0.81 \\ -0.04 & -1.83 \\ 0.76 & -0.58 \\ -0.34 & -0.50 \end{bmatrix}}_W \underbrace{\begin{bmatrix} 0.3 \\ 9.6 \end{bmatrix}}_{\beta_w} + \underbrace{\begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}}_Z \underbrace{\begin{bmatrix} 0.5 \\ 0.3 \end{bmatrix}}_{\beta_z} + \underbrace{\begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \varepsilon_4 \end{bmatrix}}_{\text{error}}$$

β_z describes impact of belonging to particular block

Problem: SBMs are not very expressive network models

Solution: random dot product graphs

$$A_{ij} \mid \mathbf{X} \sim \text{Bernoulli}\left(\mathbf{x}_i^T \mathbf{x}_j\right)$$
$$\mathbf{x}_i \sim F$$

F is a $k \ll n$ dimensional inner product distribution. Then

$$\mathbb{E}[A \mid \mathbf{X}] = \mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{S} \mathbf{U}^T$$

Define $\mathbf{X} = \mathbf{U} \mathbf{S}^{1/2} \in \mathbb{R}^{n \times k}$ to be the latent positions

Regression incorporating network principle components

Regression with rank- k truncated eigendecomposition of A

$$\mathbb{E}[A \mid \mathbf{X}] = \mathbf{X}\mathbf{X}^T = \mathbf{U}\mathbf{S}\mathbf{U}^T$$

$$\mathbb{E}[Y_i \mid \mathbf{W}_i, \mathbf{X}_i] = \beta_0 + \mathbf{W}_i\beta_w + \mathbf{X}_i\beta_x$$

β_x still roughly “effect of belonging to block Z_i with popularity γ_i ”, but β_x is rotationally unidentifiable

Regression incorporating network principle components

This has been done in [Le and Li \(2021\)](#)!

$$\mathbb{E}[Y_i | \mathbf{W}_i, A] = \mathbf{W}_i \beta_w + \xi + \alpha$$

Define S_k to be truncated eigenspace of $\mathbb{E}[A | \mathbf{X}]$ and
 $\mathcal{R} = \text{col}(\mathbf{W}) \cap S_k$. Require $\xi \in \mathcal{R}$, $\mathbf{W}_i \beta_w \perp \mathcal{R}$ and $\alpha \perp \mathcal{R}$

$$\mathbb{E}[Y_i | \mathbf{W}_i, A] = \underbrace{\mathbf{W}_i \beta_w}_{\text{covariate effect}} + \underbrace{\mathbf{W}_i \theta}_{\text{joint network covariate effect}} + \underbrace{\alpha}_{\text{network effect}}$$

Solves an identifiability issue we sweep under the rug

AddHealth data from Le and Li (2021)

- Data on 2,152 high school students (nodes)
- 7,986 self-reported friendships (edges)
- Outcome: measure of mental health
- Covariates: race, sex, grade in school

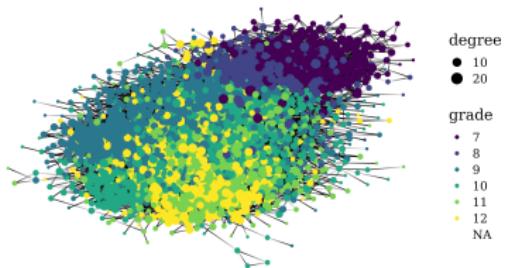


Figure 1: Grade based homophily in a high school social network.

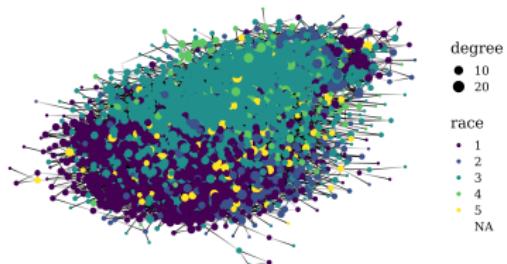


Figure 2: Race based homophily in a high school social network.

Results when applied to AddHealth data

- dimension k of latent space estimated to be 9
- $\dim(\mathcal{R})$ estimated to be zero, interpreted as no network-outcome confounding
- α : strong and significant network effect
- β_w : significant sex and grade effects conditional on network
- β_w : weak effect of race conditional on network
- OLS estimates strong effect of race unconditional on network

Results when applied to AddHealth data: a mystery

- dimension k of latent space estimated to be 9
- $\dim(\mathcal{R})$ estimated to be zero, interpreted as no network-outcome confounding
- α : strong and significant network effect
- β_w : significant sex and grade effects conditional on network
- β_w : weak effect of race conditional on network
- OLS estimates strong effect of race unconditional on network

Why does controlling for latent position in network make the race coefficient go away?

Spoiler

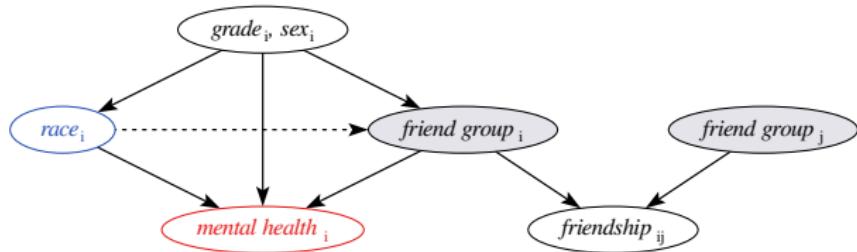


Figure 3: Results are consistent with this causal model, where race causes group membership and group membership causes mental health. Controlling for friend group (position in network) leads to overcontrol bias when estimating the average treatment effect of race

Causal reduced rank network regression



Network $A \in \{0, 1\}^{n \times n}$

Nodal covariates $\mathbf{W}_1, \dots, \mathbf{W}_n \in \mathbb{R}^{p+1}$

Nodal outcomes $Y_1, \dots, Y_n \in \mathbb{R}$

Partition $\mathbf{W}_i = (T_i, \mathbf{C}_i)$

Treatment $T_i \in \{0, 1\}$

Controls $\mathbf{C}_i \in \mathbb{R}^p$

No interference or contagion!

Interpreting regression coefficients

In observational settings: $\beta_t = 2$ implies that, on average, a one-unit increase in T_i is associated with a two-unit increase in Y_i

When controlling for all confounders: $\beta_t = 2$ is an estimate of the average treatment effect

β_t can also estimate other causal quantities ([VanderWeele, 2015](#))

Causal estimands

- Average treatment effect: how much the outcome Y would change on average if the treatment T were changed from $T = t$ to $T = t^*$

$$\Psi_{\text{ate}}(t, t^*) = \mathbb{E}[Y_t - Y_{t^*}]$$

- Controlled direct effect: how much the outcome Y would change on average if the mediator \mathbf{X} were fixed at level \mathbf{x} uniformly in the population, but the treatment were changed from $T = t$ to $T = t^*$

$$\Psi_{\text{cde}}(t, t^*, \mathbf{x}) = \mathbb{E}[Y_{t\mathbf{x}} - Y_{t^*\mathbf{x}}]$$

Causal estimands

- Natural direct effect: how much the outcome Y would change if the exposure T were set at level $T = t^*$ versus $T = t$ but for each individual the mediator \mathbf{X} were kept at the level it would have taken, for that individual, if T had been set to t^*

$$\Psi_{\text{nde}}(t, t^*) = \mathbb{E}[Y_t \mathbf{x}_{t^*} - Y_{t^*} \mathbf{x}_{t^*}]$$

- Captures the effect of the exposure on the outcome that would remain if we were to disable the pathway from the exposure to the mediator

Causal estimands

- Natural indirect effect: how much the outcome Y would change on average if the exposure were fixed at level $T = t^*$ but the mediator \mathbf{X} were changed from the level it would take if $T = t$ to the level it would take if $T = t^*$

$$\Psi_{\text{nie}}(t, t^*) = \mathbb{E}[Y_t \mathbf{x}_t - Y_t \mathbf{x}_{t^*}]$$

- Captures the effect of the exposure on the outcome that operates by changing the mediator

$$\Psi_{\text{ate}}(t, t^*) = \Psi_{\text{nde}}(t, t^*) + \Psi_{\text{nie}}(t, t^*)$$

Identifying assumptions for mediated effects

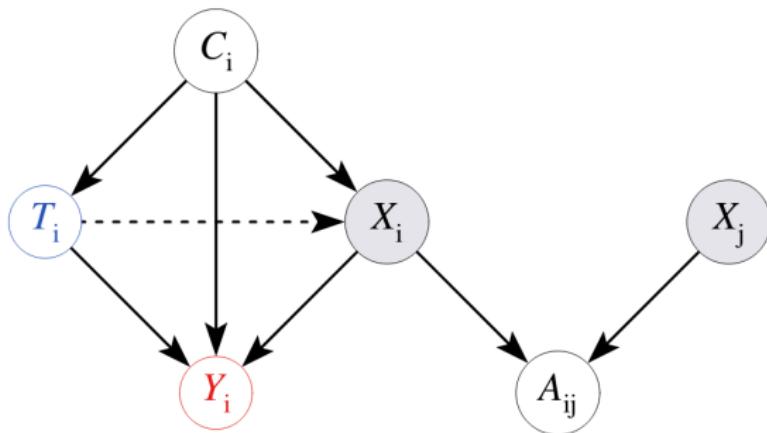


Figure 4: Directed acyclic graph (DAG) relating for node i in a reduced rank network regression model. Portions of the DAG corresponding to nodes $j \neq i$ are omitted. X_i and X_j are not observed.

Consequences for regressions

Suppose the latent positions \mathbf{X} are mediators. Then:

- Network regressions β equal the natural direct effect of T on Y , not the average treatment effect
- If we compute a separate regression, we can compute the natural indirect effect of T on Y (i.e. race causes group memberships causes mental health outcomes)
- These two effects add up to the average treatment effect
- Ordinary least squares ignoring the network structure can be used to estimate the average treatment effect

A causal re-interpretation of Le and Li (2021)'s results

Under mediation model:

- dimension k of latent space estimated to be 9
- $\dim(\mathcal{R})$ estimated to be zero, interpreted as no network-outcome confounding race, grade and sex do not perfectly cause membership in any friend group
- α : strong and causal effect of friend group on mental health
- β_w : significant sex and grade natural direct effects
- β_w : small natural direct effect of race
- OLS estimates large average treatment effect of race

Since $\Psi_{ate} = \Psi_{nde} + \Psi_{nie}$ and Ψ_{ate} is large and Ψ_{nde} is small, Ψ_{nie} must be large; there is large natural indirect effect of race on mental health

Semi-parametric causal identification

If the mediation DAG holds (non-parametric assumption!) and additionally

$$\underbrace{\mathbb{E}[Y \mid T, \mathbf{C}, \mathbf{X}]}_{\mathbb{R}} = \underbrace{\beta_0}_{\mathbb{R}} + \underbrace{t}_{\{0,1\}} \underbrace{\beta_t}_{\mathbb{R}} + \underbrace{\mathbf{c}}_{\mathbb{R}^{1 \times p}} \underbrace{\beta_c}_{\mathbb{R}^p} + \underbrace{\mathbf{x}}_{\mathbb{R}^{1 \times k}} \underbrace{\beta_x}_{\mathbb{R}^k}$$

$$\underbrace{\mathbb{E}[\mathbf{X} \mid T, \mathbf{C}]}_{\mathbb{R}^{1 \times k}} = \underbrace{\boldsymbol{\theta}_0}_{\mathbb{R}^{1 \times k}} + \underbrace{t}_{\{0,1\}} \underbrace{\boldsymbol{\theta}_t}_{\mathbb{R}^{1 \times k}} + \underbrace{\mathbf{c}}_{\mathbb{R}^{1 \times p}} \underbrace{\boldsymbol{\Theta}_c}_{\mathbb{R}^{p \times k}}$$

Then:

$$\Psi_{\text{cde}}(t, t^*, \mathbf{x}) = \Psi_{\text{nde}}(t, t^*) = (t - t^*) \beta_t$$

$$\Psi_{\text{nie}}(t, t^*) = (t - t^*) \boldsymbol{\theta}_t \beta_x.$$

Intuition behind these estimands: chain rule

$$\mathbb{E}[Y \mid T, \mathbf{C}, \mathbf{X}] = \beta_0 + t \beta_t + \mathbf{c} \beta_c + \mathbf{x} \beta_x$$

$$\mathbb{E}[\mathbf{X} \mid T, \mathbf{C}] = \boldsymbol{\theta}_0 + t \boldsymbol{\theta}_t + \mathbf{c} \boldsymbol{\Theta}_c$$

$$\begin{aligned}(t - t^*) \cdot \frac{\partial \mathbb{E}[Y \mid T, \mathbf{C}, \mathbf{X}]}{\partial t} &= (t - t^*) \cdot \frac{\partial}{\partial t} (\beta_0 + t \beta_t + \mathbf{c} \beta_c + \mathbf{x} \beta_x) \\&= (t - t^*) \cdot \left(\beta_t + \frac{\partial \mathbf{x}}{\partial t} \beta_x \right) \\&= \underbrace{(t - t^*) \beta_t}_{\text{direct effect}} + \underbrace{(t - t^*) \boldsymbol{\theta}_t \beta_x}_{\text{indirect effect}} \\&\qquad\qquad\qquad \underbrace{\phantom{(t - t^*) \beta_t + (t - t^*) \boldsymbol{\theta}_t \beta_x}_{\text{total effect}}}\end{aligned}$$

Constructing purpose-built causal estimators

We want to estimate

$$\Psi_{\text{cde}}(t, t^*, \mathbf{x}) = \Psi_{\text{nde}}(t, t^*) = (t - t^*) \beta_t$$

$$\Psi_{\text{nie}}(t, t^*) = (t - t^*) \boldsymbol{\theta}_t \beta_x.$$

Standard to fit two regressions and multiply coefficients to estimate indirect effect ([VanderWeele and Vansteelandt, 2014](#))

Adjacency spectral embedding

Lemma (Lyzinski et al. (2015), Lemma 5)

Suppose that $(A, \mathbf{X}) \sim \text{RDPG}(F, n)$. Then, letting $\widehat{\mathbf{X}}_i \in \mathbb{R}^d$ denote the i -th row of $\widehat{\mathbf{X}}$, there exists a universal constant C and a sequence of orthogonal matrices $Q_n \in \mathbb{R}^{k \times k}$ such that eventually,

$$\max_{i \in [n]} \|Q_n \widehat{\mathbf{X}}_i - \mathbf{X}_i\| \leq \frac{C \log n}{\sqrt{n}}.$$

This occurs even if A is observed with sub-gamma noise (Levin et al., 2022).

Our estimator: plug ASE into ordinary least squares

Use least squares (with robust standard errors) to estimate the regression coefficients, then plug into causal estimators

$$\mathbf{W} = \begin{bmatrix} 1 & T & \mathbf{C} \end{bmatrix} \in \mathbb{R}^{n \times (1+1+p)}.$$

For estimates $\hat{\mathbf{X}}$ of \mathbf{X} , possibly equal to \mathbf{X} itself, we estimate θ_0 , θ_t , and Θ_c

$$\begin{bmatrix} \hat{\theta}_0(\hat{\mathbf{X}}) \\ \hat{\theta}_t(\hat{\mathbf{X}}) \\ \hat{\Theta}_c(\hat{\mathbf{X}}) \end{bmatrix} = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \hat{\mathbf{X}}.$$

Our idea: plug ASE into ordinary least squares

$$\mathbf{w}(\hat{\mathbf{X}}) = \begin{bmatrix} 1 & \boldsymbol{\tau} & \mathbf{C} & \hat{\mathbf{X}} \end{bmatrix} \in \mathbb{R}^{n \times (1+1+p+k)}$$

Then

$$\begin{bmatrix} \hat{\beta}_0(\hat{\mathbf{X}}) \\ \hat{\beta}_t(\hat{\mathbf{X}}) \\ \hat{\beta}_c(\hat{\mathbf{X}}) \\ \hat{\beta}_x(\hat{\mathbf{X}}) \end{bmatrix} = \left[\mathbf{w}(\hat{\mathbf{X}})^T \mathbf{w}(\hat{\mathbf{X}}) \right]^{-1} \mathbf{w}(\hat{\mathbf{X}})^T \mathbf{Y}.$$

Our idea: plug ASE into ordinary least squares

Plug estimates into standard product-of-coefficients estimator

$$\begin{aligned}\hat{\Psi}_{\text{cde}}(\hat{\mathbf{X}}) &= \hat{\Psi}_{\text{nde}}(\hat{\mathbf{X}}) = (t - t^*) \hat{\beta}_t(\hat{\mathbf{X}}) \\ \hat{\Psi}_{\text{nie}}(\hat{\mathbf{X}}) &= (t - t^*) \hat{\theta}_t(\hat{\mathbf{X}}) \hat{\beta}_x(\hat{\mathbf{X}})\end{aligned}$$

Theoretical results

Key tuning parameter: have to pick dimension of latent space k .
If you get it right, then...

Theorem (informal)

Asymptotically, regression coefficients using $\hat{\mathbf{X}}$ and \mathbf{X} converge to the same distribution under generic low-rank models for A with i.i.d. sub-gamma noise

Corollary (informal)

Regression coefficients based on $\hat{\mathbf{X}}$ are asymptotically normal and converge at \sqrt{n} -rates.

Corollary (informal)

$\hat{\Psi}_{cde}(\hat{\mathbf{X}})$ and $\hat{\Psi}_{nie}(\hat{\mathbf{X}})$ are asymptotically normal and converge at \sqrt{n} -rates.

Rotational unidentifiability of mediator coefficients

There exists some unknown orthogonal rotation Q such that

$$\sqrt{n} \begin{pmatrix} \widehat{\boldsymbol{\theta}}_0(\widehat{\mathbf{X}}) Q^T - \boldsymbol{\theta}_0 \\ \widehat{\boldsymbol{\theta}}_t(\widehat{\mathbf{X}}) Q^T - \boldsymbol{\theta}_t \\ \widehat{\boldsymbol{\Theta}}_c(\widehat{\mathbf{X}}) Q^T - \boldsymbol{\Theta}_c \end{pmatrix} \rightarrow \text{Normal}(0, \Sigma_{\theta})$$

Rotational unidentifiability of outcome coefficients

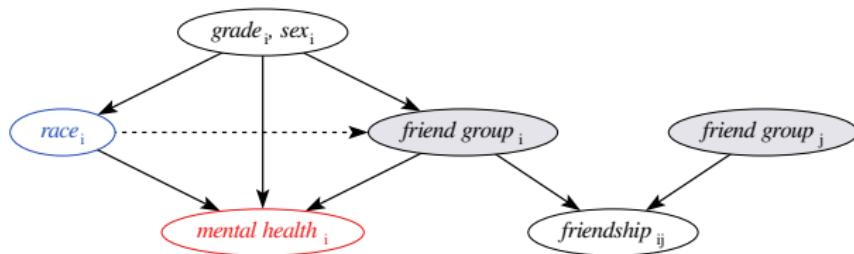
There exists some unknown orthogonal rotation Q (same as in last slide) such that

$$\sqrt{n} \begin{pmatrix} \hat{\beta}_0(\hat{\mathbf{X}}) - \beta_0 \\ \hat{\beta}_t(\hat{\mathbf{X}}) - \beta_t \\ \hat{\beta}_c(\hat{\mathbf{X}}) - \beta_c \\ Q \hat{\beta}_x(\hat{\mathbf{X}}) - \beta_x \end{pmatrix} \rightarrow \text{Normal}(0, \Sigma_{\beta})$$

Rotational unidentifiability of β_x and θ_t cancel each other out in causal estimator for $\Psi_{\text{nie}}(t, t^*)$!

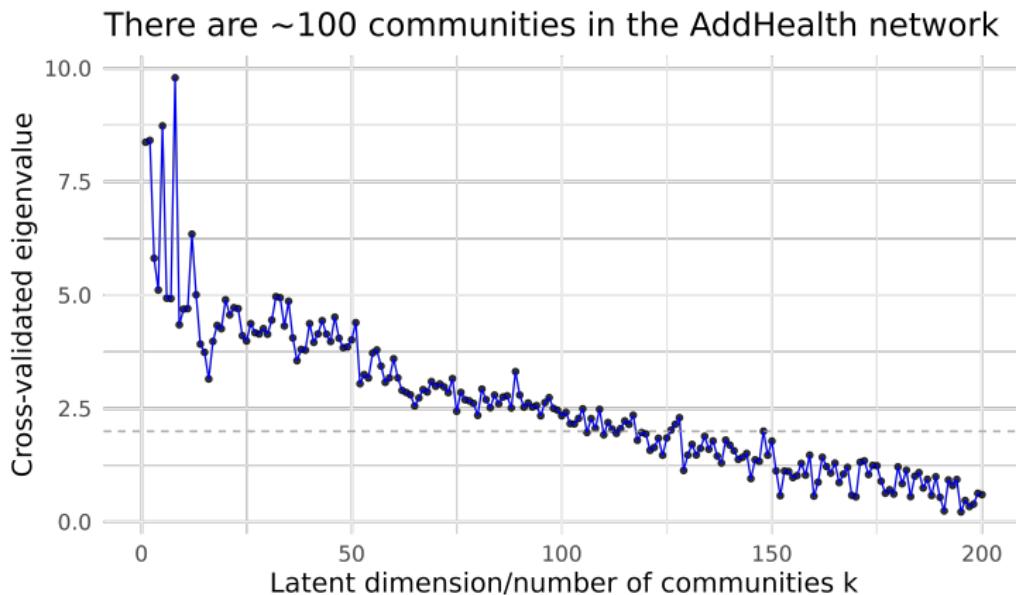
One last look at the AddHealth data

Recall



Choosing the rank of the network

- Use cross-validated eigenvalues by [Chen et al. \(2021\)](#)
- Check sensitivity of results to choice of k



Mediated causal effects in the AddHealth data

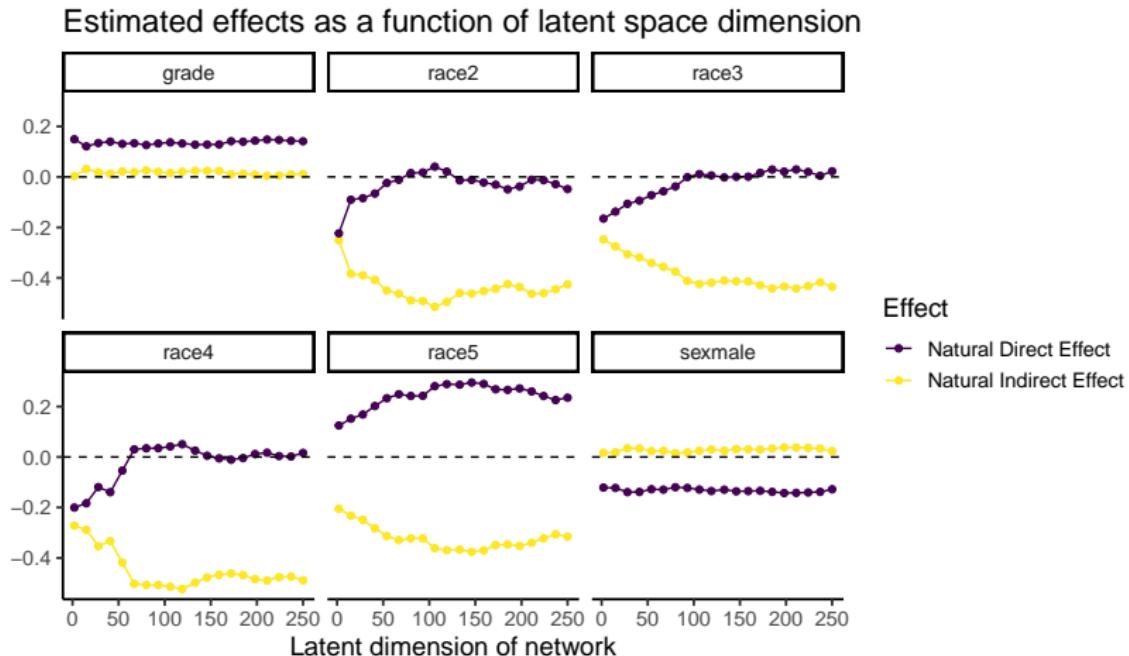


Figure 5: Point estimates of natural direct and indirect effects in the AddHealth dataset as a function of varying embedding dimension k .

Thank you! Questions?

Identifying assumptions as a system of assignments

For independent latent $u_{i,c}, u_{i,t}, u_{i,x}, u_{i,y}, u_{i,j} \sim \text{Uniform}[0, 1]$
(used for inverse probability transforms), previous DAG equivalent
to assuming

$$\mathbf{C}_i \leftarrow f_C(u_{i,c})$$

$$T_i \leftarrow f_T(\mathbf{C}_i, u_{i,t})$$

$$\mathbf{X}_i \leftarrow f_X(T_i, \mathbf{C}_i, u_{i,x})$$

$$Y_i \leftarrow f_Y(\mathbf{X}_i, T_i, \mathbf{C}_i, u_{i,y})$$

$$A_{ij} \leftarrow f_A(\mathbf{X}_i, \mathbf{X}_j, u_{i,y})$$

where f_C, f_T, f_X, f_Y, f_A are arbitrary

Interventions on a network

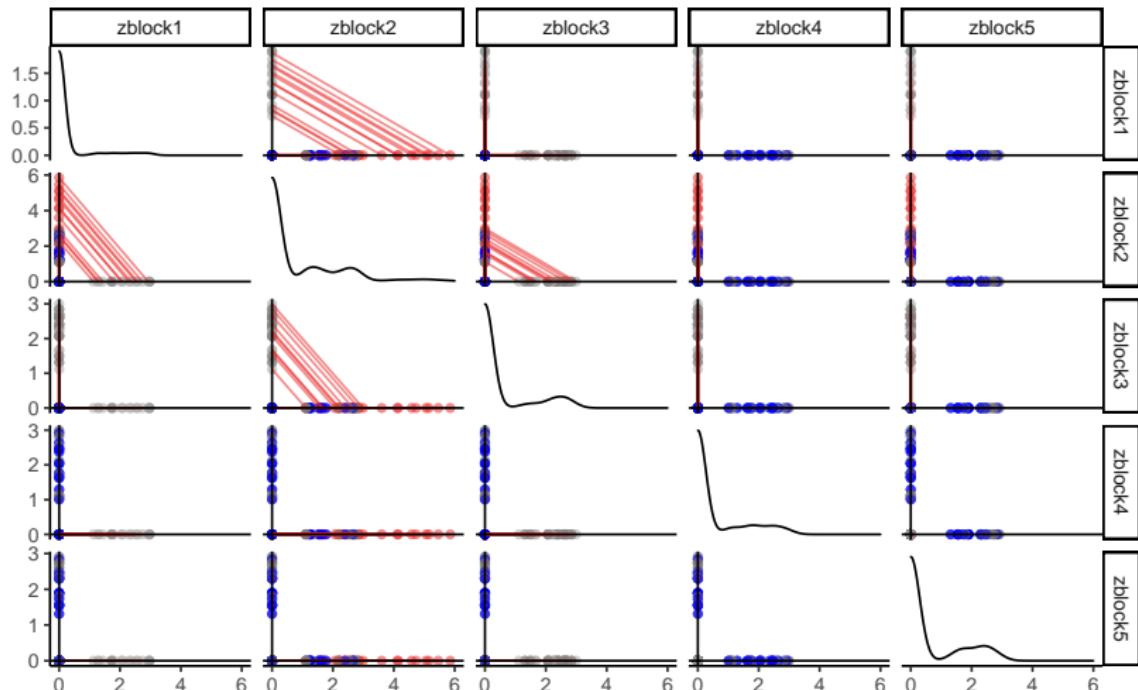


Figure 6: Canonical intervention when \mathbf{C} is highly informative.

Interventions on a network

$$\mathbb{E}[\mathbf{X} \mid T, \mathbf{C}] = \underbrace{\theta_0}_{\mathbb{R}^{1 \times k}} + \underbrace{t}_{\{0,1\}} \underbrace{\theta_t}_{\mathbb{R}^{1 \times k}} + \underbrace{\mathbf{c}}_{\mathbb{R}^{1 \times p}} \underbrace{\Theta_c}_{\mathbb{R}^{p \times k}}, + \underbrace{t}_{\{0,1\}} \underbrace{\mathbf{c}}_{\mathbb{R}^{1 \times p}} \underbrace{\Theta_{tc}}_{\mathbb{R}^{p \times k}}$$

In Figure 6, \mathbf{C} are latent parameters for a DC-SBM and
 $\theta_0 = \vec{0}, \theta_t = \vec{0}, \Theta_c = I_k$ and

$$\Theta_{tc} = \begin{bmatrix} -1 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Abstract

The last several years have seen a renewed and concerted effort to incorporate network data into standard tools for regression analysis, and to make network-linked data legible to practicing scientists. Thus far, this literature has primarily developed tools to infer associative relationships between nodal covariates and network structure. In contrast, we augment a statistical model for network regression with counterfactual assumptions and show how causal effects on a network can be partitioned into a direct effect that is uninfluenced by the network, and an indirect effect that is induced by homophily. The method is a conceptually straightforward integration of random dot product models for networks into the well-known product-of-coefficients mediation estimator.

References

Chen, F., S. Roch, K. Rohe, and S. Yu (2021, August). Estimating Graph Dimension with Cross-validated Eigenvalues. [arXiv:2108.03336 \[cs, math, stat\]](https://arxiv.org/abs/2108.03336).

Le, C. M. and T. Li (2021). Linear regression and its inference on noisy network-linked data. [arXiv:2007.00803 \[stat\]](https://arxiv.org/abs/2007.00803).

Levin, K., A. Lodhia, and E. Levina (2022). Recovering shared structure from multiple networks with unknown edge distributions. [Journal of Machine Learning Research 23](https://www.jmlr.org/papers/v23/levin22a.html), 1–48.

Lyzinski, V., D. Sussman, M. Tang, A. Athreya, and C. Priebe (2015, January). Perfect Clustering for Stochastic Blockmodel Graphs via Adjacency Spectral Embedding. [arXiv:1310.0532 \[stat\]](https://arxiv.org/abs/1310.0532).

VanderWeele, T. (2015). [Explanation in Causal Inference: Methods for Mediation and Interaction](#).

VanderWeele, T. and S. Vansteelandt (2014, January). Mediation

Analysis with Multiple Mediators. Epidemiologic methods 2(1),
95–115.