

Social Networks with Unobserved Links

Arthur Lewbel
Boston College

Xi Qu
Shanghai Jiao Tong University

Xun Tang
Rice University

Original: July 2019, revised July 2022

1 Introduction

In many social and economic environments, an individual's behavior or outcome depends on both his own characteristics and on the behavior and characteristics of other individuals.

While often assumed in practice, the linear-in-means assumption is very unlikely to hold in many applications like classrooms, where peer and contextual effects are more likely to operate through actual friendships with varying strengths, instead of equal influence from all group members. We also show how to use our identification results to empirically test the linear-in-means assumption. We reject this assumption in the STAR data.

1.1. The Model. Let $y_i \in \mathbb{R}$ and $X_i \in \mathbb{R}^K$ denote the outcome and exogenous covariates, respectively, for an individual i . Each individual belongs to one of L groups, a.k.a. networks. Assume there are n_l individuals in each group $l \in \{1, \dots, L\}$. Each group l has an unobserved $n_l \times n_l$ adjacency matrix G_l , whose (i, j) -th component is either binary (equals 1 if i is linked to j , and 0 otherwise), or is a generic number (a weight) indicating the strength of the link between i and j .¹

The researcher only observes y_i and X_i for each individual i , and the identity of the group that each individual i belongs to. The researcher does not observe the adjacency matrices G_1, \dots, G_L . For example, suppose each group is an elementary school class, and each G_l describes a network of friendships or study partners among the students in class l . The researcher observes each student i 's test score y_i and the student's vector of demographic and other characteristics X_i . The researcher also observes which class (i.e., group) each student is in, but does not observe who is friends with whom, or who studies with whom, within each class. Instead of observing or modeling the adjacency matrices of each group (i.e., class), we only assume that there is an unknown distribution of latent adjacency matrices, from which each group's matrix G_l is drawn.

We assume a standard linear social network model:²

$$y_i = \alpha + G_l y_i + X_i \beta + G_l X_i \gamma + \epsilon_i, \quad (1)$$

where y_i and ϵ_i are $n_l \times 1$ vectors of outcomes and errors, respectively, α an $n_l \times 1$ vector of ones, and X_i an $n_l \times K$ matrix of covariates. Assume for now that the errors ϵ_i are i.i.d. and uncorrelated with X_i (these conditions can be relaxed). Our asymptotics are that the number of members n_l of each network l is fixed, but the total number of networks L goes to

If the adjacency matrix G_l were observed for each group l in the sample, then point identification and estimation of these parameters under general conditions would follow from existing methods in the literature. For example, one could use the linear instrumental variables estimator of Bramoullé, Djebbari and Fortin (2009), which uses data on friends of friends, i.e., $G_l^2 X_l$, as instruments for endogenous regressors $G_l y_l$.

1.2. Intuition for Identification and Estimation. To explain the intuition for our identification strategy, let us continue to use the example of students in a class. Begin by making the simplifying assumption that all classes are the same size, having n students per class (later, in Section 6.3, we describe multiple methods of generalizing our results to handle variation in group sizes).

Equation (1) says that each element of y_l (that is, each student's test score) is a linear function of the characteristics of that student, and of the test scores and characteristics of that student's friends. One could imagine trying to directly estimate these linear functions

simulations are in the appendix.

2 Literature Review

Standard estimators of social interactions models, like Lee (2007), Bramoullé, Djebbari and Fortin (2009), and Lin (2010) assume network links are reported in the data. One popular model that does not require observing the network is the “linear-in-means” model.

effects operate through the same adjacency matrix G_i . This assumption is standard in the literature whenever both peer and contextual effects are included in a model. See, e.g., Lee (2007), Bramoullé, Djebbari and Fortin (2009), and de Paula Rasul and Souza (2020). One paper that relaxes this assumption is Blume, Brock, Durlauf and Jayaraman (2015). This assumption is generally imposed because it would be difficult to distinguish from data the extent to which any observed link applies to peer effects versus to contextual effects. We are not aware of any data sets where such information has been collected. However, since our identification is intended precisely to cover situations where link data is not, or cannot, be observed, it is possible that our methods could be extended to cover such models. We discuss the possibility of extending our method to cover this case of multiple adjacency matrices within each group in Appendix E.

We conclude this literature review by noting a deep connection between identification of linear network models and identification of traditional structural systems of linear equations, going back to the rank and order conditions described by Koopmans (1949) and the Cowles foundation, and in more detail in Fisher (1966). First, consider the setting in de Paula, Rasul and Souza (2020), which is equation (1), but simplified by having $G_i = G$ and $n_i = n$.

The linear-in-means model, which corresponds to a G having all off-diagonal elements equal to $1/(n-1)$, suffers from the "reflection problem" as pointed out by Manski (1993). The reflection problem is a failure to obtain identification because of a violation of the rank

draws from some unknown distribution of possible networks. As explained below, our method requires these networks to be exogenous from the individual characteristics whose social effects are to be identified.

By convention in the literature, the diagonal entries in each G_l are all zeros, i.e., $G_{l,ii} = 0$ for $i = 1, \dots, n_l$. The off-diagonal entries $G_{l,ij} \in \mathbb{R}$ measure the strength of the link between individuals i and j , with $G_{l,ij} = 0$ signifying the absence of a link. The unobserved adjacency matrices G_1, \dots, G_L are assumed to be row-normalized. That is, given a group adjacency matrix G_l , the (i,j) -th component in the row-normalized version G_l is $G_{l,ij} = \frac{G_{l,ij}}{\sum_{j'=1}^{n_l} G_{l,ij'}}$.

Assumption 5 (Non-trivial effects) (i) For each $k < K$, the 2-by-2 matrix

$$\begin{pmatrix} \gamma_k & \delta_k \\ \delta_k & \gamma_k \end{pmatrix}$$

has full rank. (ii) $\delta_k \neq c1$ for any $c \in \mathbb{R}$, where δ_k is a matrix of reduced-form coefficients for the K -th regressor as defined in equation (5).

Part (i) of Assumption 5 rules out the pathological case where some pair of regressors have proportional contextual and peer effects. As long as one regressor has contextual and peer coefficients that are not proportional to those of any other regressor, we can reorder the columns of X to make that regressor be the K -th regressor to satisfy part (i). A sufficient but not necessary condition for part (i) is $\delta_k = 0$ (one of the regressors has no contextual effect) while γ_k , δ_k , and γ_k are all nonzero for all $k < K$. Part (ii) of Assumption 5 rules out another pathological case, where the K -th regressor of each individual i has identical marginal effects on its own expected outcome, but no impact on that of any other group member.

In addition to Assumptions 1 to 5, to obtain identification we will require some exclusion restrictions, to satisfy a rank condition. These are discussed at length in Section 4.1.

4 Identification

The first step of our identification strategy is to show how the reduced-form parameters relate to the structural components of our model. As we show below, $E(y_j | X)$ is linear in X . Hence the reduced-form parameters can be alternatively defined as the coefficients of X in this conditional expectation.

Lemma 1 Under Assumptions 1-4, the reduced-form parameters β_0 and β_k for $1 \leq k \leq K$, defined in (5), are identified.

The proof of lemma 1 is in Appendix A, but the intuition is as follows. Let y_i denote the outcome for individual i . By construction,

$$E(y_i | X) = \beta_0 + e_i E(M)X + e_i E(MG)X ; \tag{6}$$

where e_i is a $1 \times n$ unit-vector whose i -th component is 1. Observe that the right-hand side of (6) is linear in all Kn components of X , so $E(y_j | X)$ is linear in X . This equation holds because G and M are independent from X by Assumption 3, and $E(M' | X) =$

$E [ME(" j X; G) j X] = 0$ by Assumption 2. The equality in (6) also uses the fact that the row-normalization of G implies

$$M = \sum_{s=0}^{hX-1} (G)^s = \mathbf{1} \mathbf{1}' . \quad (7)$$

The second equality here holds because, by row-normalization, each row of M adds up to the same constant $\mathbf{1} \mathbf{1}'$.

In the reduced form of equation (6), the slope coefficient for the k -th regressor of individual j is $\sum_k e_i E(M) e_j^k + \sum_k e_i E(MG) e_j^k$. (Note that, for a generic $n \times n$ matrix Q , the product $e_i Q e_j^k$ returns the $(i; j)$ -th component in Q .) The full rank and the invertibility conditions in Assumption 4 guarantee the identification of these reduced-form coefficients. These identified vectors of regressor coefficients are then arranged into the K matrices of reduced-form coefficients β_k for $k = 1; \dots; K$.

Remark 1 The representation of $E(y j X)$ in (6) is consistent not only with the simultaneous social network model with complete information given by equation (1), but also with

size is moderately large.⁷ Otherwise, the researcher needs to take measures to estimate the reduced-form coefficients using limited data. For example, instead of requiring the sample size be large relative to the number of regressors in OLS, de Paula et al. (2020) impose a sparsity condition on the structural-form adjacency matrix, and then use a penalization approach to estimate the reduced-form interaction matrix. In contrast, we propose alternative ways to deal with such data deficiency using anatomy of partitioned regressions in Section 6.2(f) (y)-3894(nS(h)1 thrthrtrs9Td(Othr6c2(tZe(c)(o)11(n611(p)11(p)(t)8n)p12(c1(p)1r12(c1(p)12(tZe(c)1(p)(tc)r343(s-9T

be a scalar multiple of I in order for (9) to hold for $(a_k; b_k)$. Case 3: $a_k \notin a_k, b_k \notin b_k$. Then (11) requires $\alpha_k = \frac{b_k - b_k}{a_k - a_k} \alpha_k$, which is a scalar multiple of α_k . Again, this implies that in order for (9) to hold for $(a_k; b_k)$, α_k must be a scalar multiple of I . In each of these three cases, the implication of (11) contradicts part (ii) of Assumption 5. \square

The reduced-form coefficients α_0 and α_k are identified by Lemma 1. Therefore, for each $k \in K - 1$;

where θ and d are column vectors that stack $\theta^{(s)}$ and $d^{(s)}$ respectively for $s = 1; \dots; S$; and

where all off-diagonal elements of $G^{(s)}$ equal $1/(n^{(s)} - 1)$. The reflection problem shows that in this model, even if $G^{(s)}$ were known, the structural parameters would not be identified without additional restrictions. Since our model includes this linear-in-means model as a special case, we must require at least as many additional restrictions for identification.⁹

There are two types of rank restrictions that are most natural to impose. The first type are exclusion restrictions, which consist of assuming that some elements of either β or γ equal zero (like the exclusion restrictions commonly used to identify linear simultaneous systems of equations). Graham and Hahn (2005) use such exclusion restrictions to identify the linear-in-means model.¹⁰ To illustrate, suppose $K = 3$ and $S = 1$. In this case it suffices to assume that one regressor X_k has no contextual effect ($\beta_k^{(1)} = 0$) and a non-zero direct effect ($\gamma_k^{(1)} \neq 0$), while another regressor $X_{k'}$ has no direct effect ($\gamma_{k'}^{(1)} = 0$) and a non-zero contextual effect ($\beta_{k'}^{(1)} \neq 0$). More generally, with $K \geq 3$, R has full rank generically if R is defined by the exclusion restrictions that there exist $k, k' < K$ with $\beta_k = 0, \gamma_{k'} = 0$ and $\gamma_k \neq 0; \beta_{k'} \neq 0$. So essentially, we get identification if one regressor has no contextual effects and another has no direct effects. In contrast, restricting two regressors to both have no contextual effects but nonzero individual effects would not suffice to make R full rank (this turns out to be a case where the order condition would be satisfied but the rank condition is not).

Since it would be unusual for covariates to have contextually no effects, we would need to assume that

still does not provide enough restrictions for identification (note that increasing S from 1 to 2 increased the number of required restrictions). However, if we impose one exclusion restriction, such as assuming that one contextual effect (i.e., one element of γ) equals zero, and we impose the constraint that $\gamma^{(1)} \neq \gamma^{(2)}$, then that provides enough restrictions to generically satisfy Theorem 1.

Note that the requirement that $\gamma^{(1)} \neq \gamma^{(2)}$ can be tested in this case, since, by equation (16), $\gamma^{(1)} \neq \gamma^{(2)}$ if and only if $m^{(1)} \neq m^{(2)}$.

The assumption that α and β do not vary by environment in this example can be relaxed. For example, if the direct effects α are the same across groups but the contextual effects vary, so $\gamma^{(1)} \neq \gamma^{(2)}$, then the full rank condition required for identification will still hold generically if one of the regressors has no contextual effect in either environment, that is, if one element in $\gamma^{(1)}$ and $\gamma^{(2)}$ equals zero.

For our empirical application in Section 7, we analyze students' math test scores. In that application, we assume two environments corresponding to small ($s = 1$) and large ($s = 2$) class sizes. For identification we allow α to vary by class size while fixing β and γ . This generalizes the models using class size variation to estimate constant peer effects (e.g., Boozer and Cacciola (2001) and Graham (2008)). We then need one additional exclusion restriction. For this we assume that a student's number of days of absence from school has an impact on his own test score but not on those of other classmates, so the element of γ corresponding to days of absence is set to zero. This exclusion restriction is motivated by the fact that common specifications of student outcomes in the empirical literature typically do

further rank restrictions. The two approaches proposed in this section could precisely serve this purpose. For example, if the model imposes no contextual effects, i.e., $\gamma_k = 0$ for $k = 1; 2; 3$, we can uniquely solve for $(\beta; \alpha_1; \alpha_2; \alpha_3)$ from the linear system (14) provided the coefficient matrix, after dropping the last three rows, has full rank (four). Alternatively, we can accommodate contextual effects but exploit the presence of multiple environments to add rank restrictions by adopting the second approach proposed above. We note that these additional required rank restrictions may in practice impose strong additional assumptions on the model.

4.2 Individual labels

Define the label of an individual in a group I to be the row of Y_I and X_I where that

ability, one could simply randomly label individuals from 1 to n in each group. However,

5 Estimation

To estimate the structural parameters of our model, we use a sample of outcomes and regressors over random networks $(y_i; X_i)_{i=1,2,\dots,L}$. Assume that across $i = 1; \dots; L$, $(y_i; G_i; X_i; \theta_i)$

function in (17) depends on $\hat{\alpha}_k^{(s)}$ smoothly. As $L \rightarrow \infty$, this objective function converges in probability, uniformly over the parameter space, to its limit where $\hat{\alpha}_k^{(s)}$ is replaced by $\alpha_k^{(s)}$. Lemma 2 implies this limit is uniquely minimized at the actual $(\alpha_k^{(s)}; \beta_k^{(s)})$. By a standard argument for the consistency of extremum estimators, $(\hat{\alpha}_k^{(s)}; \hat{\beta}_k^{(s)})$ converges in probability to $(\alpha_k^{(s)}; \beta_k^{(s)})$ for each s and k . Note that $\alpha_k^{(s)}$ and $\beta_k^{(s)}$ consist of known constants, $a_k^{(s)}$, $b_k^{(s)}$, and $m_k^{(s)}$ for $k \in K$ and $s \in S$. It then follows from the Slutsky Theorem that $\hat{\theta}$ is consistent for θ .

In Appendix A, we also explain why $\hat{\theta}$ is \sqrt{L} -convergent and asymptotically normal. Essentially, this result comes from the parametric convergence of OLS regression coefficients, and application of the delta method.

6 Extensions

6.1 Group-level variables and group fixed effects

The identification and estimation methods in Sections 4 and 5 can be readily extended to accommodate group-level regressors. Suppose each group l has a row vector of group-level characteristics $z_l \in \mathbb{R}^P$. For example these could be attributes of the teacher when each group is an elementary school class.

For the moment, consider just a single environment, so $S = 1$ and the s superscript is omitted. Including group level effects the structural model becomes

$$y_l = \alpha + G_l y_l + z_l \gamma + X_l \beta + G_l X_l \delta + \epsilon_l,$$

with $\gamma \in \mathbb{R}^P$ being a column vector of additional coefficients. One could interpret γ as a source of "correlated effects". Let Assumption 1, 2 and 3 hold with X_l replaced by $(X_l; z_l)$, and let part (i) of Assumption 4 hold with $X_l = (1; z_l)$.

Now if we have multiple environments, then run the above reduced form regressions separately for each environment s as before, but now including z_i as additional regressors. We may then identify and estimate θ from $\beta_0^{(s)}; \beta_k^{(s)}$ for $s = S$ and $R\theta = c$ as before, and estimate each $\alpha^{(s)}$ using $\hat{\alpha}^{(s)} = \hat{\alpha}^{(s)}(1 - \hat{\alpha}^{(s)})$.

Finally, this procedure can be further extended to accommodate unobserved group-level fixed effects (denoted δ_i). Essentially, we can remove these fixed effects by applying group-level demeaning of the outcomes to the reduced form, prior to recovering the structural parameters. Specifically, the method consists of replacing the dependent variables y in the first-stage reduced-form regressions with demeaned outcomes $y - \bar{y}$, and following the same steps as before to estimate the structural parameters θ . Then, we can recover the remaining parameters α and β by plugging the estimates for θ into the non-demeaned reduced form in (18), and applying an exogeneity and location normalization assumption that $E(\delta_i | z_i; X_i; G_i) = 0$. Details of this procedure are provided in Appendix F.

6.2 Dimension reduction

Again, begin by considering the case of only one environment, so s superscripts can be dropped. In the first-step regressions of

matrices β_k for $k = 1, \dots, K$. Then, given these β_k matrices, one can proceed as before to estimate the model.

With multiple environments ($S > 1$), the above regressions would be run separately in each environment, before proceeding to the later steps of identification and estimation as before. Either of the above dimension reduction methods may be especially useful in applications with multiple environments, where the number of groups in some environments s could be small relative to $Kn^{(s)}$. We adopt the second approach to estimate reduced form coefficients in our application.

6.3 Variation in group sizes

Our identification and estimation method assumes that all groups within each environment s have the same group size $n^{(s)}$. But with K individual characteristics in X , this requires observing enough groups of size $n^{(s)}$ (meaning that $L^{(s)}$, the number of groups in environment

7.1 Data description

We observe a cohort of students who were in kindergarten in 1985-1986. Seventy-nine public schools were selected to participate in the project, representing various geographic locations (inner city, urban, suburban or rural). Students and teachers were randomly assigned to classes with varying sizes of 13 to 25 students.¹⁶ Note that our estimator neither requires nor directly exploits this random assignment; however, random assignment does make some of our assumptions more plausible. An example is the dimension reduction discussed in Section 6.2.

Our sample consists of 258 classes that had at least 15 but no more than 25 students each. The total number of students in the sample is 5,189. We partition the classes in the sample into $S = 2$ environments: smaller classes with 15-20 students, and larger classes with 21-25 students according to the original design of the project. In each class, we order the students by their dates of birth, and use this ordering to label individual students. Table 7.1 reports summary statistics of the students' math test scores in the second and third grade (t_2 and t_3) and other individual-level or class-level variables to be used in our empirical analysis. These include a student's number of days of absence from school (abs), students' self-reported motivation scores (mot

the literature, is that the students enrolled in smaller classes had already developed better math skills than their peers in larger classes before the beginning of the third grade.

Table 7.1. Summary Statistics

	Small class size (122 classes)				Large class size (136 classes)			
	mean	median	std dev	range	mean	median	std dev	range
t3	620.7	618.0	40.88	[487.0, 774.0]	616.6	616.0	40.15	[510.0, 774.0]
t2	0.077	0.287	0.936	[-5.902, 1.042]	-0.029	0.287	1.023	[-6.355, 1.042]
abs	6.743	5.000	6.643	[0, 59]	6.902	5.000	6.429	[0, 55]
mot	49.29	50.00	3.990	[17, 59]	49.14	50.00	4.013	[18, 60]
tec	13.30	13.00	8.416	[0, 36]	14.19	14.00	9.079	[0, 38]

Notes: t3: raw scores for 3rd grade math; t2: standardized scores for 2rd grade math (using overall mean and std dev across all classes); abs: days of absence; mot: self-reported motivation score; tec: teacher experience (in # yrs).

Table 7.2. Test of Equal Means
(small vs. large classes)

	p-value		p-value
t3	0.001	abs	0.402
t2	< 0.001	mot	0.404

7.2 Econometric specification

Our model, corresponding to equation (1), is

$$t_{3,i} = \alpha^{(s)} + \sum_j G_{ij}^{(s)} t_{3,j} + \beta_1 \text{abs}_{l,i} + \beta_2 \text{mot}_{l,i} + \beta_3 t_{2,i} + \gamma^{(s)} \text{tec}_i + \sum_j G_{ij}^{(s)} \text{mot}_{l,j} + \beta_3 \sum_j G_{ij}^{(s)} t_{2,j} + \epsilon_{l,i},$$

where i and j are indices (labels) for individual students, l is an index for class, and (s) is the environment index. Each summation \sum_j is over all students in the same class l as student i . For each pair i and j , $G_{ij}^{(s)}$ is the row-normalized unobserved zero or nonzero link between the members labeled i and j in class l , in environment s . The coefficients to be estimated are peer effects $\gamma^{(s)}$, direct effects ($\beta_1; \beta_2; \beta_3$), contextual effects ($\beta_2; \beta_3$), intercepts $\alpha^{(s)}$, and correlated effects $\gamma^{(s)}$ (this last is the marginal impact of teacher experience, a group-level covariate).

The rank restrictions we have imposed for identification are as follows. First, this specification allows abs to have a direct effect ($\beta_1 \neq 0$) but no contextual effects ($\beta_1 = 0$). That is, a student's absence from school affects his own test scores, but has no impact on his classmates other than through peer effects. This is an exclusion restriction. Other covariates mot (self-reported motivation score) and t_2 (Grade 2 math score) are not restricted, and so can have both direct and contextual effects. Our second rank restriction is that we assume the individual effects and contextual effects are the same in the two environments, small and large class sizes (which is why β_2 and β_3 do not have s superscripts above). All other structural parameters, i.e., the intercept $\alpha^{(s)}$, the peer effect $\gamma^{(s)}$, and the correlated effect $\gamma^{(s)}$, are permitted to differ between small ($s = 1$) vs large ($s = 2$) classes. These con-

7.3 Estimation results

Table 7.3 reports our structural coefficient estimates. Standard errors are calculated using $B = 1000$ bootstrap samples, each of which is constructed by drawing classes from the original sample with replacement.

Estimates of peer effects are statistically significant and positive in both small and large classes, with the estimated coefficient being 0.85 and 0.92 respectively. A t-test for the equality of peer effects in small and large classes rejects the null of equality at the 1% level. The magnitudes of our estimates are comparable to earlier findings that used the same data but very different methodologies. For example, using a linear-in-means specification (with average class size of students in the previous year as an instrument) Boozer and Cacciola (2001) estimate the peer effects to be 0.86 for the second grade and 0.92 for the third grade. Defining links to be a simple function of measured social distance and employing some variance restrictions, our estimates are 0.85 and 0.92 for the second and third grades, respectively. Standard errors are 0.02 and 0.01 for the second and third grades, respectively.

7.4 Specification tests

Table 7.4: Tests for Over-identification

	p-values
low disp.	0.569
high disp.	0.358

Table 7.5: Wald Test Statistics for Linear-in-Means (d.f. = 29)

	small class (p-val)	large class (p-val)
low disp.	79.915 (<.001)	63.874 (<.001)
high disp.	45.112 (.028)	61.061 (<.001)

Table 7.6: CMD Test Statistics for Poisson Random Network (d.f. = 3)

	small class (p-val)	large class (p-val)
low disp.	49.880 (<.001)	171.327 (<.001)
high disp.	36.954 (<.001)	101.636 (<.001)

Table 7.7: Differences in Test Scores under the Linear-in-Means Network

	Est. mean	p-val
small, low disp	6.054	0.105
large, low disp	-9.596	0.060
small, high disp	5.810	0.184
large, high disp	-6.405	0.239

Notes: Est. mean : average difference in class means of grade three math scores in a network with equal weights on all friends.

Table 7.8: Impact of Counterfactual Peer Effects

	Est. mean	p-val
small, low disp	16.198	0.003
large, low disp	-11.637	0.001
small, high disp	2.954	0.620
large, high disp	-5.301	0.187

Notes: Est. mean : average difference in class means of grade three math scores when peer effects in small and large classes are swapped in a network with equal weights on all friends.

In the linear-in-means specification, for every group I in each environment s , the adjacency matrix $G_I^{(s)}$ is constant (the same for all I) with all off-diagonal elements taking the exact same value. With the s superscript dropped for simplicity, this implies that, for each individual characteristic k ,

$$\frac{\partial \beta_k}{\partial \alpha} = (I - G)^{-1} (\alpha I + \beta_k G) = \frac{1}{1 - \alpha} G (\alpha I + \beta_k G).$$

This in turn means that all the off-diagonal components in $\frac{\partial \beta_k}{\partial \alpha}$ must be identical. We calculate Wald test statistics using a 6 × 6 leading principal minor of the reduced form co

a large number of simulated draws r) of the simulated model-implied marginal effects ($\hat{G}_r(p)$) $^{-1}(\hat{I}_k + \hat{G}_r(p))$. We define the distance between these two matrices as a weighted sum of the differences in average diagonal and off-diagonal components, respectively. We estimate p by minimizing $\hat{Q}(p)$. This objective function would asymptotically converge to

Appendix

A. Proofs

Proof of Lemma 1. The outcome of each individual i in group l is

$$y_{l;i} = \mathbf{X}_l^0 \boldsymbol{\beta}_{l;i} + \varepsilon_{l;i}$$

where $\varepsilon_{l;i} = \mathbf{M}_{l;ri} \boldsymbol{\eta}_i$ with $\mathbf{M}_{l;ri}$ being the i -th row in \mathbf{M}_l , and $\boldsymbol{\eta}_i$ is a $(Kn + 1)$ -by-1 random vector:

$$\boldsymbol{\eta}_i = [\eta_{0i}; (\eta_{1i} \mathbf{M}_{l;ri} + \eta_{1i} \mathbf{M}_{l;ri} \mathbf{G}_l); \dots; (\eta_{Ki} \mathbf{M}_{l;ri} + \eta_{Ki} \mathbf{M}_{l;ri} \mathbf{G}_l)]'$$

with η_{ki} being the k -th components in $\boldsymbol{\eta}_i$. Recall that the joint distribution of $(y_i; \mathbf{X}_i)$ is directly identified in the data-generating process (DGP) under Assumption 1. By construction, for each individual i ,

$$E[\mathbf{X}_i y_{l;i}] = E[\mathbf{X}_i \mathbf{X}_i^0 \boldsymbol{\beta}_{l;i}] + E[\mathbf{X}_i \varepsilon_{l;i}] = E[\mathbf{X}_i \mathbf{X}_i^0] E(\boldsymbol{\eta}_i),$$

where the second equality holds because of the exogeneity of $(\mathbf{G}; \mathbf{X})$ in Assumption 2, and the independence between \mathbf{G} and \mathbf{X} in Assumption 3. Under the non-singularity of $E[\mathbf{X}_i \mathbf{X}_i^0]$ in Assumption 4-(i), we can recover $E(\boldsymbol{\eta}_i)$ from the joint distribution of $(y_i; \mathbf{X}_i)$ as

$$E(\boldsymbol{\eta}_i) = \left(E[\mathbf{X}_i \mathbf{X}_i^0] \right)^{-1} E[\mathbf{X}_i y_{l;i}]$$

for each $i = 1; 2; \dots; n$. Rearranging the components in $E(\boldsymbol{\eta}_i)$, we identify $\beta_0 = (1 - \gamma)$ and $\beta_k = E[\mathbf{M}_l(\eta_{ki} \mathbf{I} + \eta_{ki} \mathbf{G}_l)]$ for each $k = 1; \dots; K$. \square

Proof of Theorem 2. The estimators for reduced form coefficients in Step 1 are OLS estimators for slope coefficients in a regression. Thus under Assumptions 1-3 and 4-(i), $\hat{\beta}_k \xrightarrow{p} \beta_k, \hat{\eta}_k \xrightarrow{p} \eta_k$ for all $k = 1; \dots; K - 1$.

in Step 2 converges in probability to its population counterpart uniformly over $(a_k; b_k)$. That is, for all $k \in K$,

$$\sup_{a_k, b_k} \sum_{i,j} e_i(a_k \wedge_k + b_k \wedge_k - 1) e_j^0{}^2 - \sum_{i,j} e_i(a_k \wedge_k + b_k \wedge_k - 1) e_j^0{}^2 \xrightarrow{P} 0:$$

By Lemma 2, the limit function $\sum_{i,j} e_i(a_k \wedge_k + b_k \wedge_k - 1) e_j^0{}^2$

second step does not introduce additional sampling errors. A useful result for practitioners is that the first-step estimation precision can be enhanced using the dimension-reduction methods explained in 6.2. For example, in the current simulation example, the dimension-reduction method replaces $n = 10$ regressions on $n = K = 30$ explanatory variables with $n = 100$ regressions on K

and $\mathbf{x}_{k;ri}(\underline{n})$ denotes the i -th row of the $\underline{n} \times \underline{n}$ matrix $\mathbf{x}_k(\underline{n})$ and $\mathbf{0}$ a row vector of $(\underline{n} \times \underline{n})$ zeros.

Let $p(\cdot)$ denote the probability mass for n_i in the population. It then follows that for all $i = 1; \dots; \underline{n}$,

$$E(\mathbf{X}_i \mathbf{y}_{i;i}) = E_h(\mathbf{X}_i \mathbf{X}_i^0) [p_i(n) \mathbf{x}_i(n) + p(\underline{n}) \mathbf{x}_i(\underline{n})]$$

$$\Rightarrow E[\mathbf{x}_i(n_i)] = E(\mathbf{X}_i \mathbf{X}_i^0)^{-1} E(\mathbf{X}_i \mathbf{y}_{i;i}).$$

Thus $E[\mathbf{x}_k(n_i)]$, with n_i integrated out as a random variable, are identified and consistently estimable for $k = 1; 2; \dots; K$. Assuming $\mathbf{x}_k; \mathbf{y}_k; \mathbf{z}_k$ are the same for small and large classes, one can then proceed and apply the method in Section 4 to estimate the structural parameters of social effects. We use this method to balance group sizes within the environments of small or large classes in our application.

D. Dependent networks

In practice, the formation of links on a network may depend on individual characteristics in the data. We now discuss how to generalize our estimator to deal with this dependence.

Begin by considering a single environment s , where all groups within the environment have the same size n , and we omit the environment superscript. This procedure can be applied separately for each environment in the data to obtain reduced form coefficients, which would then be combined to obtain the structural parameters as in Theorems 1 and 2. Partition individual characteristics into two parts $\mathbf{X}_i = (\mathbf{X}_i^a; \mathbf{X}_i^e)$. Let \mathbf{X}_i^e denote an $n \times K_e$ matrix of excluded characteristics, i.e., covariates that affect outcomes but not link formation; let \mathbf{X}_i^a denote an n -by- K_a matrix that affect individuals' outcomes, link formation decisions, or both. For example, in our empirical application, we let \mathbf{X}_i^e be students' days of absence from school and test scores from previous years. This assumes friendships are independent of test scores conditional on observed demographics such as proximity of age. If we observe all variables that jointly determine network formation and outcomes, then our method can be applied after conditioning on \mathbf{X}_i^a .

There is a large and growing literature on network formation. To just name a few, Graham (2017), Hsieh, König, and Liu (2020), Hsieh, Lee, and Boucher (2020), Leung (2015), Leung (2020), and Sheng (2020) explicitly model how the links are formed as an equilibrium outcome. As stated in Graham (2019), "Ultimately, of course, the goal is to study the formation of networks and their consequences jointly, but such an integrated treatment remains largely aspirational at this stage". Our focus in this paper is on peer effects with unobserved links, so we simply adopt the conditional independence to deal with potential endogeneity in network formation.

Suppose network formation is given by $G_i = (X_i^a; u_i)$, which does not involve X_i^e . The reduced form is:

$$E(y_{ij}|X_i) = \sum_{k=1}^K h_{X_i} M_i(\beta_k I + \beta_k G_i) X_{i,ck} + M_i E(\epsilon_{ij}|X_i; G_i) dF(G_{ij}|X_i), \quad (19)$$

where $X_{i,ck}$ denotes the k -th column in X_i as before. Assume (i) ϵ_i is independent of X_i^e conditional on $(X_i^a; u_i)$ and (ii) u_i is independent of X_i^e conditional on X_i^a . These conditions allow the unobserved errors ϵ_i and u_i to be correlated conditional on X_i^a . Under these assumptions, $E(M_{ij}|X_i)$ and $E(M_i G_{ij}|X_i)$ is a function of X_i^a but not X_i^e , and

$$\sum_{k=1}^K M_i E(\epsilon_{ij}|X_i; G_i) dF(G_{ij}|X_i) = \sum_{k=1}^K M_i E(\epsilon_{ij}|X_i^a)$$

Again we start with the case of a single environment where all groups have identical size n , and we suppress the group subscript l throughout this section to simplify notation. Let G and W be two possibly different n -by- n adjacency matrices. For each group, peer effects and contextual effects operate through two different adjacency matrices G and W

has a unique solution

$$\begin{pmatrix} a_{jk} \\ b_{jk} \end{pmatrix} = \begin{pmatrix} v_{:j} & v_{:J} \\ j & J \end{pmatrix}^{-1} \begin{pmatrix} x_{:k} \\ 0 \end{pmatrix}. \quad (23)$$

Proof of Lemma E.1. It is straightforward to check that $(a_{jk}; b_{jk})$ defined in (23) solves (22). To see that this is a unique solution, suppose there exists $(\tilde{a}_{jk}; \tilde{b}_{jk}) \neq (a_{jk}; b_{jk})$ such that (22) holds with $(a_{jk}; b_{jk})$ replaced by $(\tilde{a}_{jk}; \tilde{b}_{jk})$, and

$$\begin{pmatrix} v_{:j} & v_{:J} \\ j & J \end{pmatrix}^{-1} \begin{pmatrix} \tilde{a}_{jk} & \tilde{b}_{jk} \\ a_{jk} & b_{jk} \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \neq 0,$$

where the inequality follows from the rank condition in (20). It then follows that

$$(\tilde{a}_{jk} - a_{jk}) \begin{pmatrix} v_{:j} \\ j \end{pmatrix} + (\tilde{b}_{jk} - b_{jk}) \begin{pmatrix} v_{:J} \\ J \end{pmatrix} = E(\beta_1 M + \beta_2 MW) = 0. \quad (24)$$

The last equality is ruled out by (21). \square

Lemma E.1 provides an analog to Lemma (1). It may then be possible to combine these equality constraints with rank restrictions like exclusions and multiple environments to construct a corresponding extension of Theorem 1 to attain identification of this extended model.

F. Group-level fixed effects

Our identification strategy can be extended to allow for group-level unobserved heterogeneity, i.e., group-level fixed effects. First, we note that if the group-level unobserved heterogeneity is mean independent from the group and individual-level covariates in $(z; X)$ (corresponding to the usual assumption in random effects models), then the estimation method described in Section 6.1 can be directly applied, because in this case the conditional mean of y given $(z; X)$ is as specified in equation (18).

Now, consider instead the more general fixed effects model. We now have the reduced-form

$$y = M(X + GX + \eta) + \frac{1}{1} + \frac{z}{1} + \frac{\$}{1},$$

where η is still the intercept, z are observed group characteristics and $\$$ is the unobserved group heterogeneity (fixed effects). Let $D = I - C$, where C is an n -by- n matrix of identical entries $1/n$, so that Dy returns the within transformation of y . Then under the assumptions that $E(\eta|X; G) = 0$ and $G \perp X$, a within transformation leads to

$$Dy = DM(X + GX + \eta) \quad E(Dy|X) = E(DM)X + E(DMG)X.$$

Thus we can write the reduced-form coefficients for the k -th characteristic from a regression using the within transformation as

$$\tilde{\gamma}_k = E(\tilde{\gamma}_k DM + \tilde{\gamma}_k DMG) = DE[M(\tilde{\gamma}_k I + \tilde{\gamma}_k G)].$$

Assume the rank condition in Assumption 5-(i) holds and that

$$\tilde{\gamma}_k \notin cD \text{ for any } c \in \mathbb{R}. \quad (25)$$

This condition can in principle be checked directly using the identifiable $\tilde{\gamma}_k$. It can then be established that the following system

$$a_k \tilde{\gamma}_k + b_k \gamma_k$$

Fisher, F. 1966. The identification problem in econometrics. Huntington, N.Y.: Krieger Publishing Company.

Tibshirani, R. 1996. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, 58(1), 267-288.

Whitmore, D. 2005. Resource and peer impacts on girls academic achievement: evidence from a randomized experiment. *American Economic Review*, 95(2), 199–203.

Wooldridge, J. 2010. *Econometric analysis of cross section and panel data*. 2nd Ed. MIT Press.